



1 **Benchmarking data-driven inversion methods for the estimation of** 2 **local CO₂ emissions from XCO₂ and NO₂ satellite images**

3 Diego Santaren¹, Janne Hakkarainen², Gerrit Kuhlmann³, Erik Koene³, Frédéric Chevallier¹, Iolanda
4 Ialongo², Hannakaisa Lindqvist², Janne Nurmela², Johanna Tamminen², Laia Amorós², Dominik
5 Brunner³ and Grégoire Broquet¹

6 ¹Laboratoire des Sciences du Climat et de l'Environnement, LSCE/IPSL, CEA-CNRS-UVSQ, Université Paris-Saclay, Gif-
7 sur-Yvette, France

8 ²Finnish Meteorological Institute, Helsinki, Finland

9 ³Swiss Federal Laboratories for Materials Science and Technology (EMPA), Dübendorf, Switzerland

10 *Correspondence to:* diego.santaren@lsce.ipsl.fr

11 **Abstract.**

12 The largest anthropogenic emissions of carbon dioxide (CO₂) come from local sources such as cities and power plants. The
13 upcoming Copernicus CO₂ Monitoring Mission (CO2M) will provide satellite images of the CO₂ and NO₂ plumes associated
14 with these sources at a resolution of 2 km × 2 km and with a swath of 250 km. These images could be exploited with
15 atmospheric plume inversion methods to estimate local CO₂ emissions at the time of the satellite overpass and the
16 corresponding uncertainties. To support the development of the operational processing of satellite column-average XCO₂ and
17 NO₂ imagery, this study evaluates “data-driven inversion methods”, i.e., computationally light inversion methods that
18 directly process information from satellite images, local winds and meteorological data, without resorting to computationally
19 expensive dynamical atmospheric transport models. We have designed an objective benchmarking exercise to analyse and
20 compare the performance of five different data-driven inversion methods: two implementations with different complexity for
21 the cross-sectional flux approach (CSF and LCSF) and one implementation for the Integrated Mass Enhancement (IME), the
22 Divergence (Div) and the Gaussian Plume model inversion (GP) approaches. This exercise is based on pseudo-data
23 experiments with simulations of synthetic “true” emissions, meteorological and concentration fields, and CO2M
24 observations in a domain of 750 km × 650 km centred on Eastern Germany over 1-year. The performance of the methods is
25 quantified in terms of accuracy in the single-image (from individual images) or annual average (from the full series of
26 images) emission estimates and in terms of number of instant estimates for the city of Berlin and 15 power plants in this
27 domain. Several ensembles of estimations are conducted, using different scenarios for the available synthetic datasets. These
28 ensembles are used to analyse the sensitivity of the performance to the loss of data due to cloud cover, to the uncertainty in
29 the wind or to the added value of simultaneous NO₂ images. The GP and the LCSF methods generate the most accurate
30 estimates from individual images with similar Interquartile Ranges (IQR) in the deviations between the emission estimates
31 and the true emissions between ~20% and ~60% for all scenarios. When taking the cloud cover into account, these methods



32 produce respectively 274 and 318 instant estimates from the ~500 daily images that cover significant portions of the plumes
33 from the sources. Filtering the results based on the associated uncertainty estimates can improve the statistics of the IME and
34 CSF methods, but at the cost of a large decrease in the number of estimates. Due to a reliable estimation of uncertainty and
35 thus a suitable selection of estimates, the CSF method achieves similar if not better statistics of accuracy for instant estimates
36 compared to the GP and LCSF methods after filtering. In general, the performances for retrieving single-image estimates are
37 improved when, in addition to XCO₂ data, collocated NO₂ data are used to characterise the structure of plumes. With respect
38 to the estimates of annual emissions, the root mean square errors (RMSE) are for the most realistic benchmarking scenario
39 20% (GP), 27% (CSF), 31% (LCSF), 55% (IME) and 79% (Div). This study suggests that the Gaussian plume and/or the
40 cross-sectional approaches are currently the most efficient tools to provide estimates of CO₂ emissions from satellite images
41 and their relatively light computational cost will enable analysis of the massive amount of data provided by future missions
42 of satellite XCO₂ imagery.

43 1 Introduction

44 The satellite imagery of CO₂ column-averaged dry air mole fractions (XCO₂) has been identified as an essential
45 component of a future atmospheric observing system to monitor anthropogenic CO₂ emissions, and in particular to detect
46 and monitor hotspot atmospheric plumes and thus emissions, in order to verify emission reductions or assess national
47 budgets (Ciais et al., 2015; Pinty et al., 2017). The Copernicus Anthropogenic CO₂ Monitoring (CO2M) mission was
48 designed to meet these objectives with a constellation of two to three Low Earth Orbit (LEO) satellites flying in a sun-
49 synchronous low-earth orbit crossing the Equator around 11:30 local time. Each satellite will carry an imaging spectrometer
50 providing images of XCO₂ and of NO₂ tropospheric column densities along a 250 km wide swath with a resolution of 2 km ×
51 2 km (Sierk et al., 2019). Current satellite missions, like Sentinel-5 Precursor (Sentinel-5P) and the third Orbiting Carbon
52 Observatory (OCO-3, when targeting specific sources in its Snapshot Area Map -SAM- mode), already deliver NO₂ column-
53 density and XCO₂ images, albeit, for the former, at a resolution coarser than CO2M, and for the latter, over areas and at a
54 frequency much smaller than with CO2M. Upcoming missions, such as Global Observing SATellite for Greenhouse gases
55 and Water cycle (GOSAT-GW, Kasahara et al., 2020), MicroCarb (in its “city-mode”, Pascal et al., 2017) and TANGO
56 (Landgraf et al., 2020), are expected to increase the amount of CO₂ and NO₂ images of the plumes from emission hotspots.

57 Operational services are being developed such as the Copernicus capacity for anthropogenic CO₂ emissions monitoring
58 and verification support (CO2MVS, Pinty et al., 2017; Janssens-Maenhout et al., 2020), to process these XCO₂ and NO₂
59 images for the monitoring of emissions in a systematic and global way at spatial and time scales that are relevant for
60 policymakers and to support emission mitigation actions. Plume inversion systems are used to derive estimates of the CO₂
61 emissions from local sources using satellite images of the corresponding atmospheric plumes. One of the key elements of
62 operational services will thus be standard plume inversion methods providing precise and reliable data in an automated and
63 fast manner. Various plume inversion approaches and implementations are now regularly used to process the existing



64 spaceborne atmospheric plumes images (Varon et al., 2018; Zheng et al. 2020; Kuhlmann et al., 2021; Nassar et al., 2021;
65 Jacob et al., 2022; Hakkarainen et al., 2023a). Therefore, there is a need to benchmark in a quantitative way the plume
66 inversion methods for the estimation of local emissions of CO₂, and more generally of greenhouse gases and pollutants.

67 Monitoring anthropogenic CO₂ emissions of point sources or cities from satellite XCO₂ images is challenging as
68 corresponding column-average enhancements are often small compared with the local fluctuations of the “background” CO₂
69 field due to biogenic CO₂ fluxes and to neighbour anthropogenic sources, and with the typical level of errors in the XCO₂
70 retrievals (Buchwitz et al., 2013). Despite this challenge, the potential of CO₂ imagers to estimate anthropogenic emissions
71 has been demonstrated with observing system simulation experiments (OSSEs) using synthetic data, for power plants
72 (Bovensmann et al., 2010), cities (Pillai et al., 2016; Broquet et al., 2018; Wang et al., 2020) and in a more general way, at
73 local to national scales (Santaren et al., 2021). Furthermore, several studies have shown that the joint analysis of co-located
74 NO₂ satellite observations strongly enhances the skill to detect the XCO₂ enhancement plumes from sources in XCO₂
75 images, and consequently to estimates the corresponding CO₂ emissions (Reuter et al., 2019; Kuhlmann et al., 2021). NO₂
76 observations are indeed characterised by a better signal-to-noise ratio and a generally small and low-amplitude background
77 field, due to the relatively short lifetime of nitrogen oxides (NO_x).

78 CO₂ emissions of large point sources and cities can be estimated from satellite images by plume inversion systems
79 integrating the observations with dynamical transport model simulations of atmospheric CO₂ concentrations (e.g., Broquet et
80 al., 2018; Ye et al., 2020; Santaren et al., 2021). In principle, the use of such dynamical models could support the analysis of
81 the 3D dynamical patterns of the observed plume and thus the accuracy of the inversion. They could also support the
82 derivation of the spatial distribution of the emissions within cities, and of the temporal variation of the emissions
83 corresponding to a plume in the hours preceding each satellite overpass. However they can be strongly impacted by
84 modelling errors which become critical at local scale, when trying to model plumes from emission hotspots over a few tens
85 to a few hundreds of kilometres (Brunner et al., 2023). Furthermore, their computational burden hampers their use for a
86 global and routine coverage of the sources in an operational context. *Data-driven plume inversion methods* appear to be
87 currently more suitable for such wide-scale applications (Ehret et al., 2022). These are computationally light inversion
88 methods that directly process information from satellite images and local winds and meteorological data (typically from
89 operational weather analyses), without resorting to dynamical atmospheric transport models.

90 The main data-driven approaches for estimating local emissions based on satellite images of plumes that have been tested
91 and analysed in a significant number of studies are: 1) the Integrated Mass Enhancement (IME) approach, which relates the
92 total mass of plumes to the corresponding emissions; it has been used for retrieving CH₄ emissions from airborne
93 observations (Frankenberg et al., 2016) or from fine-scale satellite data (Varon et al., 2018); 2) the Gaussian plume approach
94 which extracts emissions from the fit of plume shapes by Gaussian functions and was applied for instance to estimate power
95 plant CO₂ emissions from OCO-2 satellite data (Nassar et al. 2017; 2021); 3) the cross-sectional flux approach which infers
96 emissions from the fluxes passing through cross-sections of the plumes and whose potential to estimate CO₂ emissions of
97 power plants with CO₂ and NO₂ satellite imagery data was assessed, for instance, by Kuhlmann et al. (2021); 4) the



98 divergence (Div) approach, which derives emissions from the application of the divergence operator to fields of fluxes and
99 which was originally designed to estimate nitrogen oxide (NO_x) emissions from NO_2 data provided by the TROPOMI
100 satellite imagery (e.g. Beirle et al., 2019; 2021) and was more recently adapted to the quantification of CO_2 emissions
101 (Hakkarainen et al., 2022). Contrarily to the other methods of this study, the Div method only produces annual estimates
102 from average fields extracted from multiple images.

103 Against this background, the aim of this study is to benchmark these four data driven plume inversion approaches for the
104 monitoring of CO_2 emission hotspots with CO_2M images. We present a benchmarking framework to objectively evaluate
105 and compare the performance of different implementations of the four data-driven approaches (Sect. 2.1) to estimate CO_2
106 local emissions from such satellite data. For this purpose, we use one year of synthetic satellite observations closely
107 mimicking those expected from the upcoming CO_2M mission (Sect. 2.2) that were generated in the ESA funded
108 SMARTCARB project from high-resolution atmospheric transport simulations (e.g. Brunner et al., 2019; Kuhlmann et al.,
109 2020). The emissions of the city of Berlin and 15 large power plants are estimated from these synthetic satellite data and the
110 ability of the different inversion methods is assessed by comparing their estimates to the corresponding *true* values used by
111 the atmospheric transport model. Performances of the different inversion approaches are evaluated for 1) single-image
112 estimates that are retrieved from daily images (Sect. 3) and, 2) annual estimates that are computed from the inversion of one
113 year of data (Sect. 4). Furthermore, performances are analysed for different scenarios regarding the data used by the
114 inversions, where the impacts of considering the cloud cover in the data, the uncertainties in the wind and the use of
115 collocated NO_2 data are assessed. Finally, results are discussed by analysing 1) the potential of ensemble approaches that
116 would gather different inversion methods and, 2) the trade-off between overall accuracy and number of estimates when the
117 cases are filtered based on the uncertainties in the estimates computed by the plume inversion methods (Sect. 5).

118 **2 Data and methods**

119 **2.1 Data-driven inversion methods**

120 Five different emission quantification methods are evaluated in this study: (1) the integrated mass enhancement method
121 (IME), (2) the cross-sectional flux (CSF) method, (3) the light cross-sectional flux (LCSF) method, (4) the Gaussian plume
122 (GP) method and (5) the divergence (Div) method. More precisely, what is studied here are specific configurations of certain
123 methods as is the case for the CSF and LCSF “methods” which are derived from the same general approach. But, hereinafter
124 we will refer to these configurations as methods to avoid weighing down the text. The general approaches have been widely
125 used and described in previous papers such as Varon et al. (2018) and Beirle et al. (2019, 2021). The specific
126 implementations of the CSF and Div methods tested here have been used extensively by the authors in previous studies
127 (Kuhlmann et al., 2019, 2020, 2021 and Hakkarainen et al., 2022). They have been slightly upgraded in the course of this
128 benchmarking exercise to improve their stability, accuracy, and capability of running in a fully automated way. Details of the
129 methods are presented in an accompanying study by Kuhlmann et al. (2023). Further details about the theory of the Div



130 method and its application are given in Koene et al. (2023) and Hakkarainen et al. (2022, 2023b). All algorithms and tools
131 used in this work have been integrated into a Python library for *data-driven emission quantification* (ddeg), which has been
132 made publicly available and is described in Kuhlmann et al. (2023). We provide below a short description of these methods
133 with an emphasis on their relative advantages and limitations and on the way they estimate uncertainty. The main features of
134 the methods are summarised in Table 1 and illustrated in Figure 1 and Figure A1. Table 1 also lists the computation times of
135 the methods calculated for the same inversion example using the same hardware. As the methods have all been implemented
136 in the same Python package, the timings are directly comparable.

137 All methods except the Div method can provide estimates derived from individual satellite images. The Div approach as
138 implemented here is based on the averaging of information contained within multiple images and hence typically delivers
139 annual estimates. We will hereinafter refer to the IME, CSF, LCSF and GP methods as single-image methods. These
140 methods share a common algorithmic sequence that starts with identifying clusters of enhancements above a background in
141 satellite images. Subsequently, these clusters are assigned to plumes from specific known sources, and finally, the emissions
142 of the corresponding sources are estimated. The plume detection combines the first two stages and can be used to discern
143 plumes from unreported sources; however the ability of the different approaches to detect unknown point sources has not
144 been studied here, as the primary focus is to analyse their potential to detect and process plumes of known sources from
145 CO₂M-like satellite images (see Sect. 2.2). Moreover, as previously mentioned, a benefit of the CO₂M mission is the
146 availability of co-registered XCO₂ and NO₂ columns, which can further benefit the plume detection and emission
147 quantification steps.

148 Obtaining the column enhancements over the background can be achieved with different thresholding techniques as
149 detailed below. When it comes to NO₂, the global background field is insignificant but in the case of CO₂, its amplitude is
150 important and can vary significantly in space and time due to biogenic and other anthropogenic fluxes surrounding the
151 sources of interest and due to gradients in the background. Another common feature is the need for defining an effective
152 wind speed, which describes the average mass transport of CO₂ within the plumes. This a major challenge as wind speed
153 varies with altitude whereas satellite images contain integrated column measurements with no vertical resolution.
154 Additionally, the horizontal resolutions of wind products are generally different from those of satellite images. To address
155 these limitations, the methods determine effective winds in a more or less sophisticated manner.

156 Finally, all methods have implemented some quality control on their estimates. These checks are more or less restrictive
157 depending on the methods and may filter out, for example, cases with overlapping plumes originating from neighbouring
158 sources. Further details are provided in Kuhlmann et al. (2023). Of particular note is the fact that our implementation of the
159 GP method discards values that are below 1/4 or beyond 4 times the “true” values averaged one hour before the satellite
160 overpass (10:00 to 11:00 UTC); this filtering stabilises the otherwise underdetermined inversion. Unlike the other methods,
161 the GP method thus uses a priori information about the source strength, which artificially improves its performance.



162 2.1.1 Cross-sectional flux (CSF) inversion method

163 The cross-sectional flux inversion method has been used in many studies such as for example the determination of CH₄
164 emissions of point sources from high-resolved satellite data for which its superiority over other methods has been
165 demonstrated within the framework of the study of Varon et al. (2018). In brief, this method calculates the fluxes through
166 single or multiple cross-sections of the plumes as the product of effective winds and integrals of column mass enhancements
167 along plume transects (line densities). Under the assumption of steady-state conditions, these fluxes are equivalent to the
168 emissions. The CSF method used in this study has been used by Kuhlmann et al. (2020, 2021) for the estimation of CO₂
169 emissions from CO₂ and NO₂ images. These studies have demonstrated that the inclusion of NO₂ observations significantly
170 increases the number and precision of the estimates.

171 The plume detection module of the CSF approach determines in a first stage the CO₂ or NO₂ pixels that are significantly
172 enhanced above the background with a statistical z-test (Kuhlmann et al., 2021). To perform this, a Gaussian kernel to
173 average local observations values is applied and the background field is at this stage computed by applying a median filter.
174 The parameters defining the z-test were carefully assessed in order to get enough valid pixels to describe a plume while
175 avoiding false detections (Kuhlmann et al. 2019). The detected pixels are then grouped by a labelling algorithm and assigned
176 to a source. Finally, a curve representing the centerlines of the plume is fitted to the detected pixels.

177 For the quantification of CO₂ emissions, the CSF method groups the detected plume pixels into sub-polygons along the
178 curved plume, whose width equals ~5 km (2-3 pixels of CO₂M data). All detected pixels within a sub-polygon are used to
179 construct a single estimate of the line density. Following Reuter et al. (2019), the CSF method assumes that the plume
180 transect follows a Gaussian behaviour, after removing the background signal with a normalised convolution. To obtain the
181 line densities, the integration of the fitted Gaussian functions does not require any additional computation as the line
182 integrals are simply equal to the amplitude parameters of the fitted Gaussian functions. Then, in order to be converted into
183 fluxes, line densities are multiplied by effective winds which are the horizontal winds at the corresponding source locations
184 and times of the satellite overpasses, vertically weighted by the GNFR-A/SNAP-1 emission profile (Brunner et al., 2019).

185 Finally, the CO₂ emission of a given source retrieved from a given satellite image is computed by averaging the CO₂
186 estimated fluxes of all the sub-polygons describing the plume downstream of the source. The uncertainty in the emission
187 estimate is then computed by propagation of the uncertainties in the line densities computation and in the wind; the
188 uncertainties in the line densities are extracted from the standard deviation of the sub-polygon estimates and capture mostly
189 satellite data noise through uncertainty in the Gaussian fitting.

190 When NO₂ data are used in conjunction with CO₂, detections of plumes are first performed for NO₂, while the CO₂ and
191 NO₂ enhancements are fitted simultaneously by Gaussian functions that share the same mean (or central location) and the
192 same standard deviation. Thus, the fit of CO₂ enhancements takes advantage of the better signal-to-noise ratio of NO₂ data
193 by better constraining the parameters of the Gaussian functions, which provides more accurate estimates of CO₂ line
194 densities and hence CO₂ emissions.



195 **2.1.2 Light cross-sectional flux (LCSF) inversion method**

196 The light cross-sectional flux method shares the same theoretical foundations as the CSF method, but its implementation
197 is largely different. It is derived from the method originally developed by Zheng et al. (2020) to estimate the CO₂ emissions
198 of cities and industrial areas in China that produce atmospheric plumes clearly detectable in transects of OCO-2 data which
199 are characterised by a resolution of few km² and by a narrow swath about 10 km wide. This method has been applied to the
200 routine and automatic estimation of isolated clusters of CO₂ emissions worldwide (Chevallier et al., 2020) and to study the
201 temporal variability of the emissions based on several years of OCO-2 and OCO-3 data (Chevallier et al., 2022). The method
202 has undergone significant modifications for this comparative study, where the location of the emission sources is known, in
203 order to fully harness the potential of high-resolution satellite imagery.

204 For a given source and satellite overpass, the LCSF method performs a simple detection of the plume by extracting from
205 the satellite image an area which is 100 km wide in across-wind (perpendicular) direction and which extends downwind the
206 source over a distance equal to the distance travelled by the wind in one hour. The method then selects the pixels of the
207 extracted area where XCO₂ or NO₂ enhancements – simply defined as the difference between data values and the average
208 data of the area – are greater than the spatial variability, i.e. the standard deviation of the data contained within the area.

209 The quantification of the source emission is then performed on each selected enhancement by extracting again a 100 km
210 wide across-wind area centred at the enhancements and extending 10 km (~5 CO₂M pixels) downwind from the
211 enhancements. The sums of a linear term accounting for large scale variations in the background fields and a Gaussian
212 function describing the plume cross-section perpendicular to the wind direction are then fitted to the data contained within
213 these areas. The plume detection and fitting of the enhancements can be carried out in the same way when NO₂ data are
214 available. And, standard deviations and means of the Gaussian functions fitted with NO₂ data are then used for fitting CO₂
215 enhancements; CO₂ data constrain in this case only the amplitudes of the CO₂ Gaussian functions. This allows transferring
216 information derived from NO₂ data when estimating CO₂ emissions from CO₂ data.

217 CO₂ line densities are, as for the CSF method, derived from the Gaussian functions fitted with CO₂ data and converted
218 into emission estimates by the multiplication of an effective wind. For the LCSF method, this effective wind is extracted at
219 the location of the enhancements and at an altitude above ground of 100 m, as preliminary tests have shown that extracting
220 winds at the altitude of 100 m yields, for the LCSF approach, better inversion results compared to other altitudes or
221 alternative methods of computing the effective winds.

222 Finally, under steady-state atmospheric conditions, the cross-sectional CO₂ flux derived at each selected enhancement is
223 equivalent to the upwind source emissions. Therefore, as several enhancements belonging to a same atmospheric signature of
224 a source are generally processed, the algorithm produces multiple individual estimates of the source emission; the estimate
225 computed by the method for a given source and from a given image is then computed as the median value of these individual
226 estimates; the use of the median helping to reduce the impact of outliers. Moreover, uncertainties in the individual estimates



227 provided by the LCSF method are computed by propagation of the errors derived by the fitting algorithm when generating
228 the line densities; uncertainties in the final estimates are finally the median of these uncertainties.

229 **2.1.3 Gaussian plume (GP) inversion method**

230 The Gaussian plume inversion approach assumes that observed plumes can be described with Gaussian plume models. This
231 approach has been widely used such as for example in the determination of CH₄ point source emissions (Varon et al., 2018),
232 the use of OCO-2 data to quantify CO₂ emissions from power plants (Nassar et al., 2017), or in a framework to estimate at
233 the global scale CO₂ emissions from large cities and point sources (Wang et al., 2020). Compared to previous Gaussian
234 plume inversions, the GP inversion method used in this work allows the Gaussian plume model (like the CSF method) to
235 handle curved plumes (see Sect 3.2.1 in Hakkarainen et al., 2023b).

236 The detection of plumes, i.e. of the CO₂ or NO₂ enhancements from the background, is carried out using the same
237 algorithm as for the CSF method. Then, the inversion uses a Levenberg-Marquardt least-squares optimization to find the
238 optimal parameters of the Gaussian functions fitting the enhancements and, of the Bézier curves describing the centre lines
239 of the plumes (Hakkarainen et al., 2023b). If NO₂ data and CO₂ data are simultaneously available, then the Gaussian plume
240 model is first fitted to the NO₂ observations and the optimised parameters regarding the plume shape are subsequently used
241 as first guesses for the fitting to CO₂ observations. These derived parameters are constrained to remain close to the optimised
242 parameters obtained from the fitting of NO₂ data. Finally, the uncertainties in the Gaussian plume estimates are obtained by
243 propagation of the uncertainties in the fitted parameters for the wind speed and for the source strength.

244 To ensure the convergence of the minimization algorithm, first-guessed values of the fitted parameters need to be
245 carefully prescribed: parameters of the centre-line curves, for example, are initialised from the curves retrieved by the plume
246 detection algorithm, and the initial wind speed is calculated as in the CSF method (see Sect. 2.1.1). Most importantly, the
247 prior values of emission parameters are set to the *true* summertime source emission strength. Thus, unlike any of the other
248 methods studied in this work, the GP method integrates an important constraint on the emissions which implies that the
249 estimated values, hence the method's performance, are not entirely determined by the information contained within the
250 synthetic satellite observations alone. This limitation should be taken into account when applying this method to invert from
251 real satellite data emissions of sources whose amplitudes are barely known.

252 **2.1.4 Integrated mass enhancement (IME) method**

253 The IME method integrates the total mass enhancements of CO₂ or NO₂ above the background that can be associated with
254 detectable plumes. Then, following Frankenberg et al. (2016), the relationship between IMEs and emissions (Q) can be
255 approximated by a linear relationship defined by the residence times (τ) of the species within the plumes: $Q = \frac{1}{\tau}IME$. The
256 residence time can in turn be expressed as a characteristic plume length L divided an effective wind speed U_{eff} : $Q = \frac{1}{\tau}IME =$



257 $\frac{U_{eff}}{L}IME$. Varon et al. (2018), who applied the IME method with CH₄ observations, derived U_{eff} from 10 m wind speeds
258 using large eddy simulations (LES).

259 Here, the plume detection algorithm which identifies either CO₂ or NO₂ enhancements from the background is the same
260 as the one used by the CSF and GP methods, but the detected area of the plume over which the integration is performed is
261 dilated using a circular kernel in order to increase the number of integrated pixels (Hakkarainen et al., 2023b). Missing
262 values are filled using a normalised convolution and estimates are rejected when less than 75% of valid pixels are available
263 for the detected plume. The characteristic length L is computed from the centre-line of the plume as the arc length to the
264 most distant detected pixel minus 10 km, but at least 10 km. Moreover, the effective wind speed U_{eff} is extracted by using the
265 same vertically weighted average as the CSF method. If NO₂ observations are used in conjunction with CO₂ observations,
266 the integration area is established by the application of the plume detection algorithm with NO₂ data. Then, to estimate CO₂
267 emissions, the IME is calculated over this area with CO₂ observations. Finally, the uncertainty in the IME estimates is
268 computed by propagation of uncertainty from the single sounding precision of satellite data and an estimate of the
269 uncertainty in the wind speed.

270 2.1.5 Divergence method

271 The divergence method, initially introduced by Beirle et al. (2019, 2021), was used to estimate NO_x emissions based on
272 TROPOMI NO₂ observations. For this study, the method has been modified in order to estimate CO₂ emissions, as outlined
273 in Hakkarainen et al. (2022) where a detailed theoretical analysis of this approach can be found in the supplementary
274 material. The divergence method is based on the continuity equation at steady state (Jacob, 1999), where the divergence of a
275 vector field F (flux) is defined as the difference between emissions E and sinks S : $\nabla \cdot F = E - S$. Since CO₂ lifetime is
276 extremely long, the sink term can be neglected. However, before applying the divergence operator to XCO₂ images, the
277 atmospheric background needs to be removed in order to extract purely the XCO₂ enhancements. For this purpose, a median
278 filter is applied to the data and the resulting field is subtracted from the original data. Moreover, in order to improve the
279 accuracy of the estimates when CO₂ noise levels are high, data first undergo a denoising process using a 5×5 pixel mean
280 filter. The flux field F is then defined at each pixel as $F = (F_x, F_y) = (\Delta I \cdot U_{eff}, \Delta I \cdot V_{eff})$, where ΔI is the vertical column
281 density enhancement above background, and U_{eff} and V_{eff} are the eastward and northward winds, respectively, interpolated
282 at the location of the pixel and at the time of the satellite observations, and vertically averaged using the GNFR-A/SNAP-1
283 emission profile (Brunner et al., 2019).

284 Divergence maps are computed from flux fields by using a finite difference approximation and in order to clearly detect
285 point sources, the method needs to average the divergence fields over a long period. Here, divergence maps are averaged
286 over one year. For a specific source, the annual estimate of the emissions is then computed from the enhancement in the
287 averaged divergence field by using a peak fitting approach which fits the divergence map by a function including a Gaussian
288 and a linear term centred at the source (Beirle et al, 2021). Emissions, and more generally the parameters, of the peak



289 function are determined by an adaptive Markov chain Monte Carlo (MCMC) that also provides the uncertainties in the
290 estimates from the standard deviations of the sampled posterior distributions of the parameters.

291 **2.2. Synthetic satellite observations of CO₂ and NO₂**

292 In this study, synthetic satellite observations of CO₂ and NO₂ were generated from atmospheric simulations in order to
293 evaluate and compare the ability of the methods described in Sect. 2.1 for retrieving CO₂ or NO₂ emissions from point
294 sources or urban areas using satellite imagery akin to that provided by the upcoming CO2M mission. These simulated
295 satellite data are readable by the ddeq Python library and were produced as part of the SMARTCARB project and have been
296 extensively described and used in previous works (e.g. Brunner et al., 2019; Kuhlmann et al., 2019; 2020; 2021). They are
297 openly accessible from <https://doi.org/10.5281/zenodo.4048227> (Kuhlmann et al., 2020b).

298 Atmospheric concentrations of CO₂ and NO₂ were simulated by the COSMO-GHG atmospheric transport model (Jähn et
299 al., 2020) with a vertical resolution of 60 levels up to an altitude of 24 km and with a horizontal resolution of about 1 km × 1
300 km for a domain centred over the city of Berlin. The domain extends about 750 km in the east-west and 650 km in the south-
301 north direction. Simulations provided hourly outputs for nearly the entire year 2015. In order to generate realistic
302 simulations, initial and lateral boundary conditions for meteorological variables and tracers were extracted from products of
303 the European Centre for Medium-Range Weather Forecasts (ECMWF) and MeteoSwiss (Kuhlmann et al., 2019).
304 Furthermore, CO₂ emissions included both the anthropogenic and biospheric components which were interpolated onto the
305 COSMO grid at a temporal resolution of one hour: anthropogenic emissions were largely derived from the TNO/MACC-3
306 inventory (Kuenen et al., 2014) and biospheric fluxes were simulated with the Vegetation Photosynthesis and Respiration
307 Model (VPRM, Mahadevan et al., 2008). NO_x emissions were also derived from the TNO-MACC-3 inventory and
308 atmospheric simulations used a simplified NO_x chemistry with a fixed NO_x decay time of 4 hours. NO_x concentrations were
309 converted to NO₂ concentrations using an empirical equation for the evolution of NO₂ : NO_x ratios downwind of emission
310 sources (Düring et al., 2011).

311 To generate synthetic satellite observations similar to CO2M observations, the XCO₂ and NO₂ column densities derived
312 from the COSMO-GHG simulations were sampled at the resolution of 2 km × 2 km along 250 km wide satellite tracks
313 (Kuhlmann et al., 2019); these tracks were computed using an orbit simulator and correspond to a hypothetical constellation
314 of six CO2M satellites. In addition to XCO₂ and NO₂ column-average data, a cloud mask was generated from the total cloud
315 fraction computed by the COSMO-GHG model. For CO₂ data, all pixels with cloud fraction larger than 1% were removed as
316 CO₂ retrievals are strongly impacted by clouds (Taylor et al., 2016). For NO₂ data, less sensitive to clouds, a threshold of
317 30% on the cloud fraction was used to select valid pixels (e.g. Boersma et al., 2011). Figure 2 illustrates a COSMO-GHG
318 simulation of XCO₂ over the SMARTCARB domain, on which are represented synthetic XCO₂ data corresponding to a
319 CO2M satellite overpass.

320 For the purposes of this benchmarking study, we use the configuration of the SMARTCARB dataset where the CO2M
321 constellation consists of three satellites. By choosing this, we follow the recommendation of Kuhlmann et al. (2021) that a



322 constellation of at least three CO₂M satellites is necessary for a proper estimation of the annual emissions from weak
323 sources and in regions such as central Europe where cloud cover dramatically reduces the number of estimates. When
324 ignoring clouds, this constellation of three satellites leads to observing each local source within the SMARTCARB domain
325 once every other day; if we consider that a satellite image is usable if there are at least 50 data pixels next and downwind to
326 the source, then we can use about 3000 images to determine the emissions of the 16 local sources considered in this study.
327 But, if we consider the cloud cover, only 500 images remain usable.

328 The characteristics of the uncertainties in the synthetic CO₂M observations were computed using three different
329 uncertainty scenarios (low, medium, high). Simulated XCO₂ column densities were thus assigned random errors by
330 employing various levels of instrumental noise in the error parameterization formula. This formula, used for generating the
331 errors, takes into account the Solar Zenith Angle (SZA) and surface albedos (Buchwitz et al., 2013). The NO₂ column
332 densities were assumed to be characterised by random uncertainties of different constant values depending on the chosen
333 uncertainty scenario. These values are defined for clear sky conditions and increase in the presence of clouds; nearly
334 doubling for a cloud fraction of 30%. No systematic errors were prescribed for either XCO₂ or NO₂ column averaged data. In
335 this study, the characteristics of the random uncertainties prescribed to the synthetic data are chosen according to the
336 requirements of the CO₂M mission (Meijer et al., 2019). For XCO₂ retrievals, random errors are generated using the error
337 parameterization formula with a single sounding precision of 0.7 ppm for vegetation albedos and a SZA of 50°. For NO₂
338 retrievals, a single sounding precision in cloud-free conditions of 2×10^{15} molecules cm⁻² is prescribed.

339 2.3. Benchmarking scenarios

340 The relative performance of the different inversion methods to estimate CO₂ emissions are evaluated for the 15 strongest
341 point sources of the SMARTCARB domain and for the city of Berlin (Fig. 2 and Table 1 in Kuhlmann et al., 2021). These
342 16 sources covers a large emission range that extends from 3.7 MtCO₂.yr⁻¹ for the power plant located in Chvaletice (CZ) to
343 40.3 MtCO₂.yr⁻¹ for the power plant located in Jämschwalde (DE); these values being the annual mean emissions at the time
344 of the satellite overpass (10:30 UTC) used in the COSMO-GHG simulations. It is worth mentioning that the distribution of
345 the source emissions is skewed towards the lowest value as the median emission rate in the collection is around 9.6
346 MtCO₂.yr⁻¹ and 75% of the sources emit less than 14 MtCO₂.yr⁻¹.

347 In order to thoroughly evaluate the relative performance of the different methods and the sensitivity of these
348 performances to different factors, the benchmarking study is carried out according to several scenarios that share the same
349 features for the simulated data and for the source collection that have been described above. The most optimistic or ideal
350 scenario considers that inversions are performed with CO₂ and NO₂ cloud-free data using directly the winds from the
351 COSMO-GHG simulations (SMARTCARB winds). It is the ideal case because 1) with the inclusion of NO₂ data, the data
352 constraints on the estimates are stronger than when using CO₂ data only; 2) the absence of clouds maximises the number and
353 quality of the estimates, and 3) the winds are perfectly consistent with the data as they were used to simulate the XCO₂ and
354 NO₂ fields. The results derived from this benchmarking scenario should be seen as an upper limit of what the inversion



355 methods could achieve in terms of accuracy and number of estimates. The most realistic scenarios take cloud cover into
356 account and use winds extracted from the ERA5 wind product (Hersbach et al., 2020) that is independent from the inverted
357 data and whose resolution ($\sim 0.25^\circ$) is much coarser than that of the SMARTCARB winds ($\sim 0.01^\circ$). The results derived from
358 this benchmarking scenario should be seen as a lower limit for the method's performance.

359 The differences between the ERA5 and SMARTCARB wind products are significant at the 16 sources considered in this
360 study: the annual mean biases between these two wind products in 2015 range from 0.1 m.s^{-1} to 1.5 m.s^{-1} depending on the
361 source with an average value across the sources of 0.6 m.s^{-1} while RMSEs range from 1.1 m.s^{-1} to 2.1 m.s^{-1} depending on the
362 source with an average value across the sources of 1.5 m.s^{-1} (Fig. A3). The biases per source are systematically positive since
363 SMARTCARB tends to provide larger winds than ERA5. With such differences, comparing scenarios with the same
364 characteristics but using different wind products allows us to gain insight into the method's sensitivity to wind uncertainties.
365 Additional benchmarking scenarios were designed to test the sensitivity of the methods with respect to other factors,
366 including the consideration of cloud cover in satellite data and the use of NO_2 for plume detection and characterization. All
367 benchmarking scenarios are listed in Table 2.

368 2.4. Benchmarking metrics

369 For a given benchmarking scenario, the performances of the different inversion methods can be evaluated through the
370 number of single-image estimates that can be retrieved regarding the number of available satellite images: ~ 500 or ~ 3000
371 considering or not the cloud cover in the data. Performances can be assessed as well through the quality of the estimates; the
372 accuracies of the methods are then assessed by comparing the estimates retrieved from single satellite overpasses to the
373 corresponding *true* values that were used to generate the synthetic satellite data. More precisely, inversion results are
374 analysed in terms of distributions of the differences between the estimated and the true emissions of all the sources
375 considered in this study. We will refer to these differences in the following as *deviations*. More precisely, our analysis will
376 mostly focus on examining the distributions of the *relative* deviations, i.e. the differences between estimated and true
377 emissions divided by the true emissions, in order to fairly compare results across sources with significantly different
378 magnitudes (Sect. 2.3). Furthermore, to properly describe distributions that may be very different from Gaussian
379 distributions, box plots are used, in which the median values, the interquartile ranges (IQRs), the 10th and the 90th percentiles
380 of the distributions are represented.

381 The ability of the different inversion methods to estimate source emissions can also be analysed from the study of the
382 annual or monthly averages of the single-image estimates. Benchmarking results are then evaluated for each source in terms
383 of relative deviations of the annual/monthly estimates from the annual/monthly true emissions and, in terms of Root Mean
384 Square Errors (RMSE) in order to provide a global indicator for the accuracy of the annual/monthly estimates across all
385 sources.

386 In this study, the annual/monthly averages of the single-image estimates for a given source are computed using three
387 different methods which are 1) the arithmetic means of all the single-image estimates of the source emission that have been



388 generated from inverting one year/month of data, 2) the means of these estimates weighted by the inverse of their computed
389 variances (Sect. 2.1) and 3) the medians of these estimates. The annual/monthly inverse variance weighted means
390 incorporate the information provided by the methods on the quality of the estimates when averaging, whereas the
391 annual/monthly medians are statistical indicators that are more robust to outliers than the means. Moreover, since the Div
392 method is applied by temporally averaging satellite observations over the year, it produces only a single annual estimate for
393 each source; we will thus consider that the three types of annual/monthly estimates are all equal to this single estimate.

394 It is important to note that the annual and monthly estimates are affected by temporal sampling biases when inversion
395 methods use data filtered by cloud cover. Specifically, the presence of denser cloud cover during winter generally results in
396 over-representation of emission estimates during summer and hence could lead to an underestimation of annual estimates as
397 emissions are higher during winter due to increased electricity consumption. Although more advanced methods, such as
398 fitting periodic curves to capture seasonal cycles as demonstrated by Kuhlmann et al. (2021) could potentially enhance the
399 accuracy of estimates, they are not included in this study. However, these temporal sampling biases are integrated in the
400 results as the annual/monthly estimates are compared to the true annual/monthly emissions which are computed by
401 considering all the days of the year/months.

402 **3 Results on emission estimates based on individual images**

403 The following subsections present a comparative study of the CSF, GP, IME, and LCSF methods for estimating emissions
404 from single images. In the following, we will refer to these kinds of estimates as *single-image* estimates. Note that, as the
405 methods use different algorithms for plume detection and emission quantification, which include different rejection criteria
406 (Sect. 2.1), they produce different sets of estimates.

407 **3.1 Sensitivity to the emission strengths of the sources**

408 In the optimal scenario (cloud-free, SMARTCARB winds, CO₂ and NO₂ data), all methods tend to provide more accurate
409 estimates for strong sources than for weak sources, and this trend is particularly noticeable for the IME and CSF methods
410 (Fig. 3). The median values of the absolute relative deviations for weak sources (emissions ranging from 0 to 6.9 MtCO₂/yr
411 in the 1st row of Fig. 3) are 207% (IME method) and 54% (CSF method), respectively. In contrast, for strong sources
412 (emissions ranging from 15.6 to 53.2 MtCO₂/yr in the 4th row of Fig. 3), they are approximately 47% (IME) and 28% (CSF),
413 respectively. The inversion methods are also more prone to produce unrealistic values for weak sources as the distributions
414 are strongly skewed for this type of sources: the 95th percentile accuracy indicator is indeed 1128%, 584%, 172% and 178%
415 for the IME, CSF, GP and LCSF inversion models respectively (1st row in Fig. 3). For strong sources, this indicator is
416 significantly lower, decreasing to 200%, 108%, 90% and 76%, respectively (4th row in Fig. 3). Atmospheric signals
417 generated by strong sources are more distinct from the background than those from weak sources and as a result, the signal-
418 to-noise ratio in the XCO₂ and NO₂ images is better which helps to reduce uncertainties in the determination of their



419 emissions. For low-emitting sources, the performance of the inversion methods can be degraded by the limited number of
420 enhanced pixels that are detected in images with noise; this limitation makes the identification of plume centre-lines by the
421 CSF, IME and GP methods challenging (Sect. 2.1). This problem could have impacted the GP method, but its current
422 implementation incorporates prior knowledge filtering out estimates that fall outside the 25% to 400% range from the prior.
423 This filtering process is expected to improve the accuracy of the GP method, especially for weak sources.

424 Biases in the emission estimates may also depend on the strength of the source, as observed in the IME and CSF methods
425 which strongly overestimate the emissions of weak sources compared to strong sources. For weak sources, the median of the
426 deviation distributions for the IME and CSF models (blue bars, 1st row of Fig. 3) are +116% and +50%, respectively,
427 compared to +16% and +11% for strong sources (blue bars, 4th row of Fig. 3). This discrepancy is probably due to the plume
428 detection algorithm, which, for weak sources, may wrongly attribute enhancements from other sources in the vicinity of the
429 source of interest and thus artificially increase the amplitude of the detected emissions. Conversely, the LCSF approach
430 tends to underestimate the emissions of strong sources while slightly overestimating those of weak sources, with the median
431 of the deviation distribution being -26% (blue bar, 4th row of Fig. 3) and +12% (blue bar, 1st row of Fig. 3) respectively. The
432 underestimation of source emissions could be attributed to a tendency of the method to overestimate the amplitudes of the
433 background for non-isolated sources: contrary to the other methods, the LCSF method does not remove the influence of
434 neighbouring plumes when computing the background around a given source. Another explanation could lie in the fact that
435 this method uses 100-m winds as effective winds while, especially for strong emitting sources, these winds are lower than
436 the GNFR-A average winds used by the other methods.

437 **3.2 Impact of the use of NO₂ images for the detection of plumes**

438 The use of NO₂ data to identify and characterise plumes increases the number of estimates for all inversion methods
439 compared to CO₂-only inversions, as shown in Figure 4 (blue vs orange bars). The increase is significant for the IME and GP
440 methods (~93% and ~70%), moderate for the CSF method (~34%), and slight for the LCSF method (~4%). The IME, GP,
441 and CSF methods rely on a plume detection algorithm that is less reliable when using only CO₂ observations (Kuhlmann et
442 al. 2019). Of these three, the CSF method requires fewer pixels to detect and quantify plumes, resulting in a larger proportion
443 of still quantified plume cases than the IME and GP methods when having CO₂ data only. The detection of plumes by the
444 LCSF method is performed on data slices whose pixels are relatively close to sources and where XCO₂ enhancement signals
445 due to emissions are thus relatively strong; this may explain the only small benefit for this method of using joint CO₂ and
446 NO₂ images to better determine the shape of the plumes.

447 When using CO₂ and NO₂ data, the maximum number of estimates obtained from each inversion method varies
448 significantly: the IME method produces the smallest number of estimates, with 1661, while the LCSF method produces the
449 largest, with 2722. The GP and CSF methods, based on the same algorithm of plume detection as the IME method, produce
450 up to 1776 and 2012 estimates, respectively. These differences can be attributed to the differences in the number of detected
451 pixels below which the algorithm rejects plumes and, in the emission quantification algorithms used by the different



452 methods. In addition, the overall complexity of the IME, CSF and GP methods, which use a relatively large number of
453 rejection criteria likely explains why these three methods deliver much fewer estimates than the LCSF method. The relative
454 efficiency and robustness of the plume detection algorithm of the LCSF method is evidenced when using CO₂ data only to
455 determine emissions: the number and accuracy of estimates is hardly changed compared to the inversions performed with
456 CO₂ and NO₂ data; contrarily to the other methods whose algorithms are more sensitive to uncertainties in XCO₂ data and
457 which need NO₂ data to accurately fit a plume coordinate system to the data.

458 The inclusion of NO₂ data does not appear to significantly improve the overall performance of the GP and LCSF methods
459 in terms of accuracy of the CO₂ emission estimates (lower panel in Fig. 4). However, for the LCSF method, there is a notable
460 reduction in the 95th percentile of the relative absolute deviations from 175% without NO₂ to 115% with NO₂. For the CSF
461 method, the use of NO₂ data strongly improves its overall performance as the 3rd quartile and the median of the absolute
462 residuals are for example significantly decreased, from ~127% down to ~74% and from ~54% to ~36%, respectively. As the
463 CSF method rejects fewer estimates when using CO₂ data only than the GP method, its accuracy decreases because with a
464 more permissive filtering, it may include complex cases for which emissions are difficult to estimate. This may also explain
465 why the CSF estimates are less biased, with a significantly lower median relative deviation, in cases where inversions also
466 use NO₂ data (upper panel in Fig. 4).

467 In contrast, the precision of the IME method decreases when using NO₂ data, but this fact could be related to a numerical
468 artefact: the IME method performs much better for high-emitting sources than for low-emitting sources (see Sect. 3.1) and
469 the use of NO₂ data likely allows constraining small sources more efficiently than with CO₂ data only. Therefore, when
470 adding NO₂ data, the number of low-emitting sources which are estimated increases more than for the high-emitting sources
471 and then the overall performance degrades. This bias associated to the relative bad estimation of low-emitting sources is
472 confirmed when deviations are used to assess performance instead of relative deviations: the absolute deviations associated
473 to the IME estimates globally decrease with the use of NO₂ data with for example the median error decreasing from ~15 to
474 ~11.5 MtCO₂/yr.

475 3.3 Impact of the cloud cover

476 The impact of clouds is studied by comparing inversions with cloud-free images to inversions with cloud-filtered images
477 (Sect. 2.3). When disregarding cloudy pixels in the XCO₂ and column-averaged NO₂ data, the number of estimates from all
478 the methods is considerably reduced, with a decrease of 94%, 85%, 85% and 88% for the IME, CSF, GP and LCSF methods
479 respectively (Table 3). The number of estimates that can be provided for the cloud-filtered configuration with
480 SMARTCARB winds is at the maximum equal to 313 (LCSF) and decreases to 96 for the IME method which can provide
481 robust estimates for images free of clouds only as this method requires integrating enhancements over the full extent of
482 plumes. As sources are characterized by different cloud covers, the number of estimates per year and per source ranges from
483 1 to 12 (IME), from 6 to 28 (CSF), from 8 to 23 (GP) and from 15 to 26 (LCSF).



484 Furthermore, the use of cloud-filtered data not only affects the number of estimates but also impacts the performance of the
485 methods, although to a lesser extent. When comparing results obtained from the same images, cloud-free inversions produce
486 better outcomes than cloud-filtered inversions (Fig A2). This is because, in images partially masked by cloud cover, some
487 pixels containing useful information are likely removed, which can lead to less accurate determination of emissions.

488 **3.4 Impact of uncertainty in the wind**

489 As mentioned above, in order to assess the impact of potential uncertainties in the wind, a series of inversions is carried out
490 with a different wind product than the one used to generate the synthetic XCO₂ and NO₂ data. For this purpose, the
491 SMARTCARB winds are replaced by ERA5 winds and the differences between these two wind products are characterised at
492 the sites of this study by random and systematic components (Sect 2.3 and Fig. A3). Notably, ERA5 winds show
493 systematically lower values.

494 For all inversion methods, the global accuracies of the estimates, evaluated in terms of relative absolute deviations, are
495 only slightly reduced when using ERA5 winds instead of SMARTCARB winds (lower panel in Fig. 4, green vs red bars).
496 There are a few possible explanations for this: the temporal or spatial uncertainties in wind components are only a minor
497 source of uncertainty compared to other factors impacting the determination of the estimates by the different inversion
498 methods such as, for example, uncertainties in the XCO₂ and NO₂ columns densities (Sect. 2.2) or over-simplified
499 assumptions in plume detection or quantification algorithms. Kuhlmann et al. (2020, 2021) showed, for instance, that the
500 determination of the CO₂ background field could introduce significant uncertainties in the estimates. Furthermore, as
501 indicated by Reuter et al. (2019), one of the important benefits of satellite imagery is that uncertainties related to
502 meteorological variables likely average out when emission estimates are sampled along significant areas of plumes.

503 However, the fact that ERA5 wind values are systematically lower than those of SMARTCARB winds has an impact on
504 the median values of the relative deviations, i.e. on the biases in the estimates. While the accuracies in terms of relative
505 absolute deviations are slightly affected by using either wind product (bottom panel in Fig. 4, green vs red bars), biases can
506 be significantly increased, as in the cases of the GP and LCSF methods whose estimates are on average underestimated if
507 inversions use ERA5 winds instead of SMARTCARB winds. The lower amplitudes of the ERA5 winds explains also that the
508 results for the IME and CSF methods improve, especially for the 95th percentiles of the absolute deviation distributions
509 which respectively decrease from around 504% and 411% to 370% and 286% respectively. The systematic overestimation of
510 the estimates evidenced above for the CSF and the IME methods is therefore mitigated when using ERA5 winds (top panel
511 in Fig. 4).

512 As mentioned previously (Sect. 2.3), the benchmarking scenario for which inversions are performed with ERA5 winds
513 and data filtered for cloud cover, is the closest to real conditions of monitoring emissions from data images delivered by
514 satellites. For this scenario with CO₂ and NO₂ data, the GP and LCSF methods show the best performances in terms of
515 global accuracies with respectively IQRs of 25–62% and 17–55% for the distributions of the absolute relative deviations (red
516 boxes in Fig. 4). It is interesting to note that the overall accuracies of these methods are similar for this realistic scenario and



517 the ideal scenario where inversions are performed with cloud-free data and SMARTCARB winds. Contrarily, the number of
518 estimates strongly decreases when inversions are performed with cloud-filtered data such as, for example, from 2722 to 318
519 estimates for the LCSF method (see Table 3).

520 **4 Results on annual and monthly averages of the emissions**

521 **4.1 Annual estimates**

522 To evaluate how well an inversion method performs on an annual basis, we include all image estimates generated by the
523 method, regardless of their uncertainty. We calculate annual estimates for a given source using three methods, as described
524 in Sect. 2.4: 1) by taking the average of all available image estimates for the source over the entire year, 2) by taking the
525 weighted average of these image estimates based on their uncertainty, and 3) by taking the median value of these image
526 estimates. Because the Div method only provides one estimate per year, its annual estimates are the same, irrespective of the
527 calculation method used. In order to compare for a given source the three estimated annual values to the true emission, we
528 define this latter as the arithmetic mean of the true emissions values for the source over all 365 days of the year.

529 When annual estimates are calculated as arithmetic means or medians of individual image estimates, the GP and LCSF
530 methods generally outperform the other methods. Indeed, for cloud-free inversions with CO₂ and NO₂ data, the median
531 deviations for the annual arithmetic means (solid lines, 2nd column of Fig. 5) are 8% (GP), 14% (LCSF), 73% (IME), 35%
532 (CSF), and 64% (Div), and the median deviations for the annual medians (dotted lines, 2nd column of Fig. 5) are 14% (GP),
533 21% (LCSF), 54% (IME), 13% (CSF), and 64% (Div). However, if annual estimates are calculated as the means of image
534 estimates weighted by their uncertainty, the relative performance of the methods changes. In this case, the median deviations
535 for annual weighted means (dashed lines, 2nd column of Fig. 5) are 28% (GP), 48% (LCSF), 46% (IME), and 12% (CSF).
536 Thus, using weighted means to calculate annual estimates significantly improves, especially for low-emitting sources, the
537 performance of the IME and CSF methods while having a negative impact on the GP and LCSF methods. This finding
538 indicates the reliability of the uncertainties in the estimates produced by the IME and CSF methods compared to the other
539 methods and, if we use weighted means to compute annual estimates, the accuracies of the IME and CSF methods increase
540 significantly.

541 Figure 6 displays the inversion results for the annual estimates in a different but complementary way compared to Fig. 5:
542 the estimated annual emissions are represented with respect to the true ones which in particular allows illustrating whether
543 annual estimates are over- or under-estimated for a certain type of source and by a given inversion method. In order to
544 consider the best performance for each method according to what has been shown above, annual estimates represented in the
545 figure, and used for the analysis of the results made below, are arithmetic means of single-image estimates for the LCSF and
546 the GP methods, while they are weighted means for the IME and CSF methods. Furthermore, Fig. 6 illustrates more clearly
547 than Fig. 5 the fact that, when weighted averages are used as annual estimates, the latter methods produce annual estimates
548 whose precision is comparable for weak *and* strong sources while the global precision of estimates derived from single



549 images by these methods is significantly lower for weak sources (Fig. 3); averaging single-image estimates weighted by their
550 uncertainty thus strongly increases the performance of the IME and CSF methods at the annual scale for low-emitting
551 sources. However, even though the amplitudes of the relative deviations are similar between strong and weak sources, they
552 have opposite signs: annual estimates for strong sources are generally underestimated while annual estimates for weak
553 sources are generally overestimated.

554 Contrary to the results for the estimates retrieved from single images (Fig. 4), the CSF, GP and LCSF approaches show
555 similar performance, with a slight advantage for the GP method, when estimating annual emissions if we consider the
556 ensemble of the benchmarking scenarios. For example, in the case of inversions from cloud-filtered CO₂ and NO₂ data and,
557 with SMARTCARB/ERA5 winds, the relative RMSEs are 18/27% (CSF), 20/20% (GP) and 17/31% (LCSF). The analysis
558 of Fig. 3 shows that the LCSF method produces single-image estimates that are slightly more accurate but more biased than
559 that of the GP method. Thus, the compensation of errors when averaging single-image estimates over a year may be less
560 effective for the LCSF method than for the GP method leading to similar global accuracies for both methods. For instance,
561 the LCSF method has a greater tendency to underestimate high emissions (4th row of Fig. 3) which likely explain why,
562 contrarily to the GP method, it systematically underestimates the emissions of the strong emitting power plant located in
563 Jämschwalde, regardless of the inversion scenario (Fig. 6). With respect to its results for single-image estimates, the CSF
564 method has significantly better results at the annual scale when annual estimates are computed as weighted averages of
565 single-image estimates.

566 Even when annual estimates are computed for the IME method as weighted averages of the single-image estimates, this
567 method still show smaller accuracies compared to the CSF, GP and LCSF methods: the median values of the deviations for
568 the annual estimates are for example 39% (IME), 20% (CSF), 11% (GP) and 21% (LCSF) when considering the best scores
569 for the inversions performed with ERA5 winds and cloud-filtered data (4th column of Fig. 5). The relative performance of the
570 IME method is even worse when analysing the performance in terms of RMSE because, despite a weighting of estimates
571 according to their quality or uncertainty in the annual averages, this method produces for some sources annual estimates that
572 strongly deviate from the actual values, as in the cases of Boxberg or Schwarze Pumpe power plants (Fig. 6). Moreover, the
573 deviations of the Div method compared to that of the CSF, GP and LCSF methods are higher for most of sources except for
574 strong sources (true annual emissions > 15 MtCO₂/yr) when inversions are performed using cloud-filtered data and ERA5
575 winds (4th column of Fig. 5).

576 It is noteworthy that annual estimates for most inversion methods are comparable between inversions using data with or
577 without clouds (comparison between the 2nd and 3rd columns, Fig. 5), and surprisingly the deviations of the IME and Div
578 approaches are even smaller for inversions with cloud-filtered data. Despite significant differences in the number of image
579 estimates between those two (i.e., cloud-filtered and cloud-free) inversion configurations, annual estimates are *on average*
580 slightly affected when cloud cover is considered in the data, at least for the year and sources examined in this study.
581 However, even though the relatively small number of image estimates in the inversion configuration with clouds does not
582 hinder most methods from determining annual emissions of most sources, discrepancies can be high for some sources when



583 estimates do not sample correctly the entire year and thus introduce an important temporal bias. For example, the GP method
584 mostly estimates emissions during summer for the Jämschwalde power plant when it uses the cloud-filtered inversion setup,
585 explaining the strong underestimation of the annual emission of this source compared to the cloud-free case (top-left vs
586 bottom-left panel of Fig. 6); this explains additionally why the RMSE increases significantly for the GP method (from 13%
587 to 20% when inversions use SMARTCARB winds) when the cloud cover limits the number of single-image estimates. The
588 IME method is also impacted by this temporal bias when the number of estimates is too small to properly capture the
589 seasonal cycle of the emissions, as in the case of the Boxberg power plant. Moreover, whatever the benchmarking scenario,
590 most inversion methods produce annual estimates for all the sources studied in this work, with the notable exception of the
591 Div approach, which estimates annual emissions for only 10 out of 16 sources. This limitation, also present for cloud-free
592 data configurations, is related to the fact that some sources don't produce strong enough divergence peaks from which
593 annual estimates can be made by this method.

594 As for the results concerning single-image estimates, the use of ERA5 winds instead of SMARTCARB winds has on
595 average a very low impact on annual estimates delivered by the IME, CSF, GP and LCSF methods. For emissions estimated
596 from cloud-free CO₂ and NO₂ data, the median deviations when inversions use SMARTCARB winds are indeed 46% (IME),
597 12% (CSF), 8% (GP) and 14% (LCSF), and when inversions use ERA5 winds, they are equal to 46% (IME), 12% (CSF), 9%
598 (GP) and 12% (LCSF) as shown in the comparison between the 2nd and 4th columns of Fig. 5. On the other hand, the overall
599 accuracy of the Div method improves when inversions use ERA5 winds rather than SMARTCARB winds to estimate
600 emissions. In this case, annual estimates are less prone to overestimation due to the generally lower amplitude of ERA5
601 winds compared to SMARTCARB winds (Fig. A2). This also explains a stronger underestimation of the emissions of strong
602 sources by the LCSF method, resulting in a decrease in the accuracy of the annual estimates for this kind of sources when
603 this method uses ERA5 instead of SMARTCARB winds (left-bottom vs right-bottom panel of Fig. 6).

604 The overall precision of the annual estimates computed by the IME, CSF, GP and LCSF methods are, for all the
605 benchmarking scenarios, significantly higher than the overall precision of their single-image estimates. For example, when
606 inversions are performed with ERA5 winds and cloud-filtered data, which is the benchmarking scenario with the poorest
607 results, the median deviations of the annual estimates are 39%, 20%, 11% and 21% whereas the median deviations of the
608 single-image estimates are 73%, 35%, 46% and 37% for the IME, CSF, GP and LCSF methods. Despite the biases that can
609 hamper the image estimates, the compensation for errors when averaging across a year allow to generate annual estimates
610 that are more precise and this positive effect is amplified when error-weighted averages are used, as in the case of the IME
611 and CSF methods.

612 **4.2 Monthly estimates and seasonal cycle**

613 Monthly estimates can be computed using the same three methods as the annual estimates but, according to the results
614 analysed in the former section, we choose to estimate monthly emissions with the method leading to the best performance at
615 the annual scale: monthly estimates are thus calculated as the arithmetic means for the GP and LCSF methods and, as



616 weighted means for the CSF and IME methods. Then, considering the distributions of image estimates month by month
617 allows us to study how well inversion approaches capture the seasonal cycle of the true emissions. The analysis of Fig. 7
618 shows however that none of them are able to do this when the cloudy pixels are masked: the seasonal cycle of the actual
619 monthly emissions, i.e. maximal/minimal emissions for winter/summer months, is not reproduced by the inversion methods
620 whose estimates are characterised by an erratic monthly evolution leading to inconsistent seasonal cycles. Even though a
621 method correctly estimates annual emissions, some of its monthly estimates can be in important disagreement with the *true*
622 monthly emissions as it is the case for the CSF method on the Heyden source or for the LCSF method on the Dolna Odra
623 source (Fig. 7). Moreover, the methods generally fail to produce estimates for the winter months of the year due to the
624 temporal sparsity of data when the impact of the cloud cover is taken into account.

625 If the number of estimates is higher, i.e. when clouds are not considered in the data, seasonal cycles derived from
626 monthly estimates are in better agreement with that of the observations for most of inversion methods: the amplitude of the
627 seasonal cycle of the data can be well reproduced as it is the case for the Jänschwalde and Dolna Odra sources for example
628 (Fig. A4). But, the averaged values of the seasonal cycles of the monthly estimates, i.e. the annual estimates, can still be in
629 strong disagreement with that of the data even though the number of estimates is higher; this fact supports the presence of
630 systematic biases in the estimates that was evidenced for most of the methods in the analysis of the results for single-image
631 image estimates (Sect. 3.1).

632 **5 Discussion**

633 **5.1 Accuracy vs number of estimates**

634 For a given benchmarking scenario, the analysis conducted in Section 3 has evaluated the performance of the different
635 methods in inferring estimates from individual images by considering all the estimates provided by each method for this
636 scenario. In other terms, the analysis did not integrate any diagnostic regarding the quality of the estimates from these
637 methods. However, we demonstrated in Sect. 4.1 that computing annual means of estimates weighted by their uncertainties
638 can significantly improve the accuracy of the annual estimates when uncertainties are effectively characterised as in the case
639 of the IME and CSF methods. Therefore, a study of the performance of inversion methods for estimating single-image
640 estimates from synthetic XCO₂ images should as well integrate a characterization of the quality of its estimates. More
641 precisely, different performance indicators or error estimates can be derived from the application of the inversion methods
642 and such indicators can be used to identify and select the most reliable estimates. Nevertheless, there are no objective criteria
643 to impose a threshold on the quality of the estimates; higher quality thresholds come with smaller sets of estimates, and
644 optimal values depend on the inversion method. Indeed, not only do the different inversion methods calculate the
645 uncertainties in the estimates in different ways but also the computed uncertainties only reflect part of the total/actual
646 uncertainties, focusing on subsets of sources of uncertainties which differ across the different methods.



647 For a given inversion method, we attempt an effective quality indicator (QI) which would allow selecting estimates in a
648 manner that the global accuracy of the method increases when the QI increases, and which would provide indications on the
649 actual/total errors. We assume that the uncertainties in the estimates derived by the methods provide the best basis we can
650 get from the algorithms described in Sect. 2.1 for the derivation of such an indicator. In principle, since dealing with sources
651 of quantitatively different amplitudes (see Sect. 2.3) we should derive the QI in terms of *relative* uncertainties. And, if we
652 define the QI as a threshold selecting the estimates whose relative uncertainties are below it, we should select the most
653 reliable estimates regardless of the strength of the source they are associated with. However, this would be true if the
654 methods perform independently with respect to the amplitudes of the emissions and this is not the case for most methods as
655 illustrated in Sect 3.1. The CSF and IME methods for example strongly overestimate low-emitting sources compared to
656 high-emitting sources which implies that the relative uncertainties of weak sources are underestimated by these methods
657 (Fig. 3). Therefore, if the threshold value of relative uncertainty was decreased, we would tend to select more bad than good
658 estimates and the overall performance would decrease. Therefore, for these methods, we prefer to select estimates with
659 respect to their uncertainties, and not to their *relative* uncertainties, which will mitigate the impact of the bias in the
660 estimation of low-emitting sources.

661 In any case, determining whether a QI should be based on absolute or relative uncertainties depends on whether the
662 overall performance of the method improves when estimates with decreasing absolute or relative uncertainties are chosen.
663 Preliminary tests (not shown here) have established that the overall accuracy of the IME and CSF methods increases when
664 the *absolute* uncertainty below which estimates are selected is decreased. For the GP and LCSF methods, this behaviour is
665 obtained when *relative* uncertainties are used to discriminate estimates. Consistently, for all methods, the increase of
666 performance is then associated with a reduction in the number of estimates and, in order to get a significant number of high-
667 quality estimates, the value of uncertainty corresponding to the maximal accuracy of the method is arbitrarily set to the 10th
668 percentile of the distribution of the absolute/relative uncertainties. Then, by varying its QI between this value and the
669 maximal uncertainty of its estimates, each method can be thus associated to a range of accuracies with their respective
670 number of estimates for a specific benchmarking scenario (e.g. cloud-filtered or cloud-free). In other words, inversion results
671 can be represented by curves of accuracy *vs* number of estimates, which gives for each inversion method a complete
672 overview of its performance in terms of accuracy and number of estimates.

673 To assess the inherent performance of the methods without considering the impact of the cloud cover or of the
674 uncertainty in the winds, inversion results are analysed for the inversion configuration using XCO₂ and NO₂ cloud-free data
675 and SMARTCARB winds, *i.e.* the same winds used to generate the synthetic XCO₂ and NO₂ observations. Figure 8
676 illustrates that the overall accuracies of the CSF and IME methods are highly dependent on the selection of their estimates,
677 and are therefore strongly correlated with their number of estimates. For instance, the IME and CSF methods exhibit large
678 increases in the 3rd quartiles of their deviation distribution when the QIs of their estimates decrease: from 81% to 231%
679 (IME) and from 43% to 75% (CSF) respectively. For these methods, the selection of estimates based on their quality
680 indicators appears to be effective, as the 3rd quartiles and 95th percentiles, which indicate the proportion of poor estimates,



681 significantly decrease with increasing quality index, *i.e.* with decreasing number of estimates. Therefore, the IME and CSF
682 methods are very likely to produce reliable uncertainty estimates in the individual emission estimates and the definition and
683 derivation of their QI reflect the level of accuracy of their estimates.

684 The LCSF and GP methods display a slight correlation between most of their accuracy indicators and the number of
685 estimates. For instance, the 3rd quartiles of the distributions of relative absolute deviations remain relatively stable, varying
686 only from 46% to 56% and from 51% to 59% for the LCSF and GP methods respectively, over their entire range of number
687 of estimates. For these methods, the tradeoff between precision and number of estimates is not a critical issue and retrieving
688 an important number of estimates does not imply a significant deterioration in accuracy. On the other hand, this also
689 indicates that the current quality indicators for the GP and LCSF methods do not reflect the total/actual uncertainties in their
690 estimates.

691 As the methods present different sensitivities of the accuracy to the number of estimates, the relative performances of the
692 methods in terms of accuracy change according to the number of estimates. In other terms, as is the case for the LCSF and
693 CSF methods in Fig. 8, one method may outperform another method depending on the number of estimates we consider.
694 Indeed, below 1000 estimates, the CSF method is characterised by a better precision than the LCSF method for all the
695 statistical indicators and in particular for the 95th percentile of the deviation distribution. The best performance of the CSF
696 methods in terms of precision is then reached for ~400 estimates where the median of the deviations is ~25% compared to
697 ~29% for the LCSF method. But, if the number of estimates increases beyond 1000, the LCSF method starts outperforming
698 the CSF method with respect to the 95th percentile and when estimates are not filtered by their QI (right ends of the curves of
699 Fig. 8), it totally outperforms the CSF method not only in terms of precision but also in terms of number of estimates: if all
700 estimates are considered, the LCSF/CSF method generates 2722/2028 estimates whose deviations from the truth are
701 characterised by an IQR of 17%–56%/17%–75%. Furthermore, the LCSF method discards outliers much more efficiently
702 than the CSF method insofar as the 95th percentile of the deviation distribution is much lower for the former (118%) than for
703 the latter method (341%).

704 Selecting one method over another involves making a trade-off between precision and the number of estimates obtained.
705 Taking the example from Fig. 8, if the primary objective of an application is to obtain as many estimates as possible, the
706 LCSF method would be the preferred choice, as it can provide 2722 estimates with an IQR of the deviations ranging from
707 17% to 56%. On the contrary, if the main priority is to obtain estimates with the highest precision, the CSF method would be
708 more suitable, providing approximately 400 estimates with an IQR of the deviations ranging from 11% to 45%. The trade-off
709 between accuracy and number of estimates in the choice of method is even more accentuated in the case where inversions
710 are made with ERA5, as the use of this wind product increases the accuracy of the CSF method through bias compensation
711 (Sect. 3.4): in this case, using the CSF method, a maximum precision can be obtained, with an IQR equal to 11%–42%, for
712 650 estimates. If, on the other hand, the LCSF method is used, a maximum number of estimates, 2670, can be obtained with
713 an IQR of 18%–55% (Fig. A5).



714 The difficulty in achieving the best possible precision for a given method lies in determining an appropriate QI for their
715 estimates. Here, we adopted a relatively simple approach by defining high-quality estimates as those with relative or absolute
716 errors below the 10th percentile of the distribution relative to all the uncertainties of the estimates. However, as seen in the
717 curves of Fig. 8, highest precision may not be achieved at this value but at a higher one as in the examples of the IME and
718 CSF method. This is because misleading estimates, such as those resulting from the overlap of plumes from two sources, can
719 be characterised by very small uncertainties but at the same time by important deviations from the truth, and their impact on
720 the results becomes significant when the number of estimates gets relatively small. More generally, the QIs defined in this
721 study reflect the actual uncertainties in the estimates more or less well and the definition of a more reliable QI that ensures
722 increased accuracy with higher values of the indexes and deliver the maximum achievable precisions for all of the methods
723 is beyond the scope of this study, as it likely requires extensive studies in order to provide a common and an accurate
724 characterization of the total uncertainties in the estimates for all the inversion methods. Finally, we will note that all the
725 qualitative insights stated above about the relationships between accuracy and number of estimates are also valid when
726 considering inversions using cloud-filtered data and ERA5 winds (Fig. A6).

727 **5.3 Single methods vs ensemble approaches**

728 In this study, we create ensemble approaches by averaging the single-image estimates – for the same source and from the
729 same individual image – produced by different inversion methods. The aim is to obtain more robust and reliable predictions
730 if individual biases and errors associated with each approach compensate each other. We want thus to analyse whether an
731 ensemble method, although more expensive from a computational point of view, would perform quantitatively better than a
732 single method among CSF, GP and LCSF; these methods clearly outperforming the IME method in terms of accuracy and
733 number of estimates.

734 Four sets of ensemble approaches are considered: the first one integrates the CSF, GP and LCSF inversion methods, and
735 the remaining three ensemble approaches integrate pairs of methods (CSF & GP, CSF & LCSF and GP & LCSF). Moreover,
736 in order to assess the impact of the QIs of the different inversion methods on the performance of the ensemble methods,
737 results are analysed by considering 1) all the estimates and 2) only the best estimates produced by each method. As results
738 are assessed for the inversions using ERA5 winds and cloud-filtered data which provide a relatively small number of
739 estimates, we consider the best estimates as the estimates whose relative/absolute errors are below the 25th percentile of their
740 respective error distribution.

741 The ensemble approaches do not provide clear improvements in terms of estimate accuracy over the individual methods
742 from which they are derived (Fig. 9), with the exception of the important number of outliers produced by the CSF method
743 when estimates are not filtered: the 95th percentile of the deviation distribution is equal to 286% for the CSF method only,
744 while it decreases to 160% for the ensemble approach gathering the CSF, GP and LCSF methods. On the other hand, the
745 skewness of the CSF distribution of deviations lead to an increase of the 95th percentile of the deviations of the ensemble
746 approaches compared to the 95th percentiles of the LCSF and GP methods. Otherwise, the IQR of the deviations are similar



747 for all the ensemble and individual approaches and roughly ranges from 15% to 65% when estimates are not filtered and
748 from 15% to 60% when best estimates are selected. Therefore, errors and biases in the estimates produced by a given method
749 are generally not compensated by the estimates of other inversion methods which suggest that in general, for the same
750 images and sources, the estimates produced by other inversion methods may also present larger errors or similar biases.

751 The great benefit of using ensemble approaches lies in the significant increase in the number of estimates, which is a
752 crucial issue in the real world when the amount of satellite data is strongly limited by the cloud cover. The ensemble
753 approach gathering the CSF, GP and LCSF methods can supply a maximum of 412 estimates over the year analysed in this
754 study, representing a 30% increase compared to the LCSF method which is the individual method that supplies the most
755 estimates (318). This result indicates that the CSF, GP and LCSF methods can provide estimates from different images, i.e. if
756 one method does not provide an estimate from a given image, another method from the ensemble may, conversely, provide
757 one (Fig. A7). This allows the ensemble method to produce a maximum number of estimates (412) that is close to the
758 number of usable satellite images (~500). When only best estimates are considered, the ensemble approach generates more
759 than twice as many values compared to the LCSF method (195 vs 80) whereas the other ensemble approaches (CSF & GP,
760 CSF & LCSF and GP & LCSF) only provide about 140 estimates.

761 While combining the estimates generated by the CSF, GP and LCSF methods seems to be the optimal choice for an
762 ensemble approach providing the largest number of predictions, the computational cost of using these methods together may
763 not outweigh the benefits in terms of number of estimates compared to using a single method. For example, in the most
764 realistic scenario of inversions conducted with cloud-filtered data and ERA5 winds, the computational time required for the
765 CSF-GP-LCSF ensemble method is more than three times that of the LCSF method alone (see Sect. 2.1) whereas the overall
766 precision of the LCSF method is better and the increase in the number of estimates is only 30% when using the ensemble
767 approach. Therefore, if the performance of computer systems remains an important factor to take into account, one would
768 prefer to use the LCSF method, which is the fastest method of this study, instead of using an ensemble approach.

769 In order to investigate the benefit of using ensemble approaches for the estimation of annual emissions, we use the same
770 three individual methods that produce much better results than the IME and Div methods (see Sect. 4.1), but we consider
771 different definitions of the annual estimates depending on the inversion method: annual estimates are arithmetic means of
772 image estimates for the LCSF and the GP methods whereas they are weighted means for the CSF method. This choice
773 corresponds to the best performance at the annual scale that has been found in this study for each method (Sect. 4.1.)
774 Besides, no selection of the estimates was performed to compute the annual estimates although the quality of the estimates is
775 integrated within the annual estimates of the CSF method which are averages weighted by the errors in the estimates. Among
776 the ensemble methods considered here, only the approach gathering the CSF and GP methods yields better results than the
777 best individual method composing it for most of benchmarking scenarios (Fig. A8). For example, when inversions are
778 performed with cloud-filtered data and SMARTCARB winds, the CSF, GP and their ensemble approach are characterised by
779 relative RMSE equal to 18%, 20% and 16%, respectively. The benefit of using ensemble methods for estimating annual
780 estimates is thus questionable, especially considering that the gain in accuracy, if any, is very small compared to the



781 individual methods which, depending on the inversion scenario, produce the more accurate annual estimates. This is due to
782 the fact that the inversion methods generate annual estimates that are generally biased in the same way: emissions of strong
783 sources are generally underestimated while emissions of weak sources are generally overestimated (see median values in
784 Fig. 6).

785 **6 Conclusions**

786 In this paper, we tested and benchmarked several lightweight data-driven inversion methods for estimating local (city and
787 power plant) emissions from XCO₂ and NO₂ satellite images. The five methods that have been studied are the Integrated
788 Mass Enhancement (IME), the Cross-Sectional Flux (CSF), the Gaussian Plume (GP), the Light Cross-Sectional Flux
789 (LCSF) and the Divergence (Div); the last method generating only annual estimates. In a domain centred over the city of
790 Berlin, which extends about 750 km in the east-west and 650 km in the south-north direction, inversions were performed
791 with almost one year of synthetic SMARTCARB XCO₂ and tropospheric column NO₂ satellite observations with similar
792 characteristics as the upcoming CO2M mission. The ability of the inversion methods to estimate emissions has been assessed
793 by comparing the deviations of estimates from the corresponding “true” values used in the simulations, for 16 sources
794 including the city of Berlin and 15 power plants. To get a complete overview of performance, several benchmarking
795 scenarios were considered in order to analyse the benefit of using auxiliary NO₂ data or the impacts of the cloud cover in the
796 data or of uncertainties in the wind data.

797 In terms of quantifying emissions from single satellite images, the implementations of the CSF, GP and LCSF methods
798 used in this study outperform that of the IME method. Furthermore, we have demonstrated that the performance in terms of
799 accuracy and number of estimates varies, to a greater or a lesser extent depending on the method, with the selection of the
800 estimates based on their relative or absolute uncertainty. The overall accuracies of the IME and CSF methods are
801 significantly enhanced when a strict screening for high quality estimates is applied but at the cost of an important decrease in
802 the number of estimates. The GP and LCSF methods, on the other hand, perform more robustly showing only a variation in
803 their global precisions with increasing quality screening. This behaviour points out the need for these methods of a better
804 characterization of the uncertainties in the estimates. When estimates are filtered, the CSF method yields the best results in
805 terms of accuracy while, when estimates are not filtered, the LCSF method provides the highest number of estimations with
806 a slight decrease in accuracy. Overall, the CSF, GP and LCSF methods show similar accuracies for all the benchmarking
807 scenarios and when the less reliable estimates of the CSF method are removed: most of IQRs of the absolute deviations
808 range from 15% to 60% with an average median around 35%. Moreover, for the most realistic benchmarking scenario, i.e.
809 for the inversions using cloud-filtered NO₂ & CO₂ data and ERA5 winds, the IME, CSF, GP and LCSF methods generate on
810 average 6 (IME), 18 (CSF), 17 (GP) and 20 (LCSF) estimates per source and per year with great differences between sources
811 (See Sect. 3.3), which is equivalent to a maximum number of estimates equal to 96 (IME), 295 (CSF), 274 (GP) and 318
812 (LCSF) for all 16 sources. These figures are significantly lower than the number of usable images (~500) that can provide a



813 hypothetical constellation of 3 satellites as analysed here; this suggests that methodological improvements could increase the
814 number of estimates.

815 The accuracy of the CSF and IME methods was found to depend on the strength of the sources with important errors
816 when determining low emissions; the GP and LCSF methods, in contrast, show similar performances across different ranges
817 of emissions. Moreover, the advantage of using co-located NO₂ signal for plume detection and quantification appeared to be
818 clear for the CSF, IME and GP methods, for which the number of single-image estimates significantly increased, while it
819 was rather weak for the LCSF method. When a cloud cover mask was taken into account in the data, the number of estimates
820 significantly decreased for all the inversion methods with an average reduction of 85%; the global precision however hardly
821 decreased and even improved for the IME method. For all the inversion methods, the sensitivities of the results to wind
822 uncertainties were surprisingly found to be insignificant when replacing the SMARTCARB winds (used in the simulation)
823 by ERA5 reanalysis winds. Finally, if we do not take computational cost into account, the interest in using ensemble
824 approaches instead of a single method lies mainly in an increased number of single-image estimates as the availability of
825 estimates from the different methods complements each other.

826 Part of the effectiveness of the implementations of the cross-sectional flux method may come from the generation of
827 multiple estimates of cross-sectional fluxes along plumes and the subsequent averaging in order to get an unique emission
828 estimate. Probably, errors in the satellite data or in the simplifying assumptions of the cross-sectional approaches partly
829 cancel out when averaging. The CSF implementation uses a complex algorithm of plume detection which makes it possible
830 to use the total detectable plume, probably leading to more accurate estimates than for the LCSF implementation, which only
831 uses observations near the source. However, the plume detection and the computation of the curved centerline can fail for
832 weak sources (i.e. short plumes) at the cost of having a large number of outliers. On the contrary, the LCSF implementation
833 uses a simpler but more robust algorithm that uses the wind vector to estimate the location of the plume, which likely
834 explains why this method generates more estimates, and without the need of NO₂ data, compared to the CSF implementation.
835 However, efforts should be made to correct the systematic underestimation of strong emissions by the LCSF implementation.
836 A way forward can be merging the CSF and LCSF method into a single algorithm that takes the advantages of both
837 approaches.

838 When compared to other methods, the relative ability of the GP method in estimating emissions probably relies on the
839 use of a Gaussian function whose optimization determines the emissions while taking into account the entire structure of the
840 plumes, and calculating effective winds that are consistent with that of the plumes. However, this optimization and thus the
841 performance of the GP method highly depend on the first-guessed values to be assigned to its parameters (not shown). And,
842 in this study, the first-guessed values of the emissions are the summer average emissions for each source; this could be a
843 strong constraint on the estimated values and could lead to an overestimation of the GP performance in this benchmarking
844 study. Finally, the GP method is computationally expensive due to the heavy plume detection algorithm and to the multi-
845 parameter optimization required for the Gaussian fitting of the plumes (Table 1).



846 The IME method also integrates information retrieved from the entire structure of the plumes but, contrarily to the GP
847 method, it does not use this information when computing effective winds. Therefore, these winds may be inconsistent with
848 the characteristic lengths of plumes used by the IME method to estimate CO₂ emissions (Sect. 2.1.4) and this could explain
849 the relatively poor performance of the IME method in this study. Varon et al. (2018) probably found that the IME method
850 was adapted to estimate CH₄ emissions from high-resolution plumes because they inferred a relationship between the
851 effective winds and the characteristic lengths through LES simulations. Another drawback of the IME method is that it is
852 very sensitive to missing data as it needs an entire coverage of the plume area by data to efficiently integrate the total mass
853 enhancement. Other single-image methods (GP, CSF and LCSF) are less sensitive to missing data as they fit functions to the
854 data and can handle data gaps; this explains why these methods provide a much larger number of estimates when the impact
855 of cloud cover on the data is considered (see Sect. 3.3).

856 For estimating annual emissions, the CSF, GP and LCSF methods outperform the Div and IME methods when annual
857 estimates are computed as error-weighted means of single-image estimates for the CSF method and as arithmetic means of
858 these estimates for the GP and LCSF methods. Across the different benchmarking scenarios, the GP method shows better
859 precisions in its annual estimates because its single-image estimates have similar absolute deviations from the truth but are
860 less affected by biases compared to the CSF and LCSF methods (see Fig. 3). However, despite biases, errors in the single-
861 image estimates provided by the CSF, GP and LCSF methods likely compensate when averaging and these methods also
862 generate annual estimates with a better precision than for their single-image estimates. In the most realistic benchmarking
863 scenario – where inversions use cloud-filtered XCO₂ & NO₂ data and ERA5 winds and where performances are the lowest
864 compared to other scenarios – the relative RMSE for the annual emissions of the 16 sources is 20% (GP), 27% (CSF), 31%
865 (LCSF), 55% (IME) and 79% (Div). The relatively weak performance of the Div method could be explained by the fact that
866 this method was originally developed for the estimation of NO_x emissions and the fields of this chemical species are
867 generally characterised by stronger divergence peaks than for CO₂ fields. The performances of ensemble approaches
868 gathering several inversion methods in terms of annual estimations is not better, and in some cases even worse, than the
869 individual methods. Finally, none of the methods were able to correctly reproduce the monthly seasonal cycle of the
870 emissions when data underwent a cloud-filtering, i.e. when data were not available for some months, which points out the
871 need for an extensive temporal coverage of the observations when aiming to capture the monthly variability in emissions.

872 In addition to the technical improvements that could be made on the algorithms of the methods, further developments
873 could extend this study such as the integration of new data streams for estimating CO₂ emissions such as satellite data of
874 other co-emitted gases than NO₂, e.g. CO data provided by the TROPOMI instrument. A companion paper (Hakkarainen et
875 al., 2023c) analyses the ability of the inversion methods in determining NO_x emissions, from synthetic and TROPOMI NO₂
876 satellite data for the Matimba and Medupi power plants in South-Africa. The NO₂ synthetic data are extracted from the high-
877 resolution MicroHH Large Eddy Simulations (LES) (Van Heerwaarden et al., 2017) and used in particular to study the
878 nitrogen dioxide to nitrogen oxide scaling factors that are required for satellite-based estimations of NO_x emissions.
879 Moreover, the capacity of the inversion methods to estimate city emissions has been analysed in this study on the single



880 example of the city of Berlin and, as most of the methods have provided correct estimates for its emissions, it would be
881 interesting to expand this study to other cities and other local sources. Finally, this benchmarking study has not integrated the
882 new and promising type of inversion methods that are the methods derived from deep learning techniques (e.g. Lary et al.,
883 2016). After a potentially complex training phase, deep-learning methods could quickly process large amounts of data and
884 provide estimations with similar or better accuracy than the methods studied here (Dumont le Brazidec et al., 2023). They
885 could also complement these methods by allowing a fine differentiation of the plumes compared to the background with
886 advanced image segmentation techniques.

887 The aim of this study is to contribute to the development of the CO₂ Monitoring and Verification Support system that
888 will use the upcoming CO2M satellite data. And, although this benchmarking study has been performed with synthetic
889 observations, the methods studied here can be easily adapted to the analysis of real satellite observations and to deal with
890 sources of unknown location as demonstrated in Hakkarainen et al. (2023c).

891
892 *Code and data availability.* The code repository of the python package *ddeq* is available on Gitlab.com:
893 <https://gitlab.com/empa503/remote-sensing/ddeq>. The SMARTCARB dataset is available on Zenodo:
894 <https://doi.org/10.5281/zenodo.4048227>.

895
896 *Author contributions.* DS made the diagnostics and led the analysis for the intercomparison of the results from the different
897 inversion methods. All co-authors contributed to the decisions for the configuration, diagnostics and analysis of the
898 intercomparison. DS wrote the manuscript with inputs from all co-authors. DS, GB and FC carried out the analysis specific
899 to the LCSF method. JH, II, HL, JN and LA carried out the analysis specific to the Div method. GK developed the original
900 *ddeq* library that has been used as a basis for the application of the different methods. GK provided the SMARTCARB
901 dataset used to test the different methods. GK carried out the analysis specific to the IME method. EK carried out the
902 analysis specific to the CSF and GP inversion methods. The project was coordinated by JT, DB and GB.

903
904 *Competing Interests.* Some authors are members of the editorial board of Atmospheric Measurement Techniques. The
905 authors have no other competing interests to declare.

906
907 *Acknowledgements.* Most of the work performed in this paper was done in the framework of EU H2020 project CoCO2
908 (grant No. 958927). The FMI team would also like to thank the Research Council of Finland project 353082. All authors
909 would also like to thank the ICOS Carbon Portal for providing access to their JupyterLab servers, which were used for code
910 development and data sharing.



911 **References**

- 912 Beirle, S., Borger, C., Dörner, S., Li, A., Hu, Z., Liu, F., et al. (2019). Pinpointing nitrogen oxide emissions from space.
913 Science Advances 5. doi:10.1126/sciadv.aax9800
- 914 Beirle, S., Borger, C., Dörner, S., Eskes, H., Kumar, V., de Laat, A., et al. (2021). Catalog of NO_x emissions from point
915 sources as derived from the divergence of the NO₂ flux for TROPOMI. Earth System Science Data 13, 2995–3012.
916 doi:10.5194/essd-13-2995-2021
- 917 Boersma, K. F., Eskes, H. J., Dirksen, R. J., van der A, R. J., Veefkind, J. P., Stammes, P., Huijnen, V., Kleipool, Q. L.,
918 Sneep, M., Claas, J., Leitão, J., Richter, A., Zhou, Y., and Brunner, D.: An improved tropospheric NO₂ column retrieval
919 algorithm for the Ozone Monitoring Instrument, Atmos. Meas. Tech., 4, 1905–1928, [https://doi.org/10.5194/amt-4-1905-](https://doi.org/10.5194/amt-4-1905-2011)
920 2011, 2011.
- 921 Bovensmann, H., Buchwitz, M., Burrows, J. P., Reuter, M., Krings, T., Gerilowski, K., et al. (2010). A Remote Sensing
922 Technique for Global Monitoring of Power Plant CO₂ Emissions from Space and Related Applications. Atmos. Meas. Tech.
923 3, 781–811. doi:10.5194/amt-3-781-2010
- 924 Broquet, G., Bréon, F.-M., Renault, E., Buchwitz, M., Reuter, M., Bovensmann, H., et al. (2018). The Potential of Satellite
925 Spectro-Imagery for Monitoring CO₂ Emissions from Large Cities. Atmos. Meas. Tech. 11, 681–708. doi:10.5194/amt-11-
926 681-2018
- 927 Brunner, D., Kuhlmann, G., Marshall, J., Clément, V., Fuhrer, O., Broquet, G., Löscher, A., and Meijer, Y.: Accounting for
928 the vertical distribution of emissions in atmospheric CO₂ simulations, Atmos. Chem. Phys., 19, 4541–4559,
929 <https://doi.org/10.5194/acp-19-4541-2019>, 2019.
- 930 Brunner, D., Kuhlmann, G., Henne, S., Koene, E., Kern, B., Wolff, S., ... & Fix, A. (2023). Evaluation of simulated CO₂
931 power plant plumes from six high-resolution atmospheric transport models. *Atmospheric Chemistry and Physics*, 23(4),
932 2699-2728.
- 933 Buchwitz, M., Reuter, M., Bovensmann, H., Pillai, D., Heymann, J., Schneising, O., et al. (2013). Carbon Monitoring
934 Satellite (CarbonSat): Assessment of Atmospheric CO₂ and CH₄ Retrieval Errors by Error Parameterization. Atmos. Meas.
935 Tech. 6, 3477–3500. doi:10.5194/amt-6-3477-2013.
- 936 Chevallier, F., Feng, L., Bösch, H., Palmer, P. I., and Rayner, P. J.: On the impact of transport model errors for the
937 estimation of CO₂ surface fluxes from GOSAT observations, Geophys. Res. Lett., 37,
938 21, <https://doi.org/10.1029/2010GL044652>, 2010.
- 939 Chevallier, F., Zheng, B., Broquet, G., Ciais, P., Liu, Z., Davis, S. J., et al. (2020). Local anomalies in the column-averaged
940 dry air mole fractions of carbon dioxide across the globe during the first months of the coronavirus recession. *Geophysical
941 Research Letters*, 47, e2020GL090244. <https://doi.org/10.1029/2020gl090244>



- 942 Chevallier, F., Broquet, G., Zheng, B., Ciais, P., & Eldering, A. (2022). Large CO₂ emitters as seen from satellite:
943 Comparison to a gridded global emission inventory. *Geophysical Research Letters*, 49, e2021GL097540.
944 <https://doi.org/10.1029/2021GL097540>
- 945 Ciais, P., Crisp, D., v. d. Gon, H., Engelen, R., Heimann, M., Janssens-Maenhout, G., Rayner, P., and Scholze, M.: Towards
946 a European Operational Observing System to Monitor Fossil CO₂ emissions – Final Report from the expert group,
947 Copernicus climate Change Service, Report, European Commission, Brussels, 2015.
- 948 Crisp, D., Pollock, H. R., Rosenberg, R., Chapsky, L., Lee, R. A. M., Oyafuso, F. A., et al. (2017). The on-orbit performance
949 of the Orbiting Carbon Observatory-2 (OCO-2) instrument and its radiometrically calibrated products. *Atmos. Meas.Tech.*
950 10, 59–81. doi:10.5194/amt-10-59-2017
- 951 Dumont Le Brazidec, J., Vanderbecken, P., Farchi, A., Broquet, G., Kuhlmann, G., & Bocquet, M. (2023). Deep learning
952 applied to CO₂ power plant emissions quantification using simulated satellite images. *Geoscientific Model Development*
953 *Discussions*, 2023, 1-30.
- 954 Düring, I., Bächlin, W., Ketzler, M., Baum, A., Friedrich, U., and Wurzler, S. (2011). A New Simplified NO/NO₂ Conversion
955 Model under Consideration of Direct NO₂-Emissions. *metz* 20, 67–73. doi:10.1127/0941-2948/2011/0491
- 956 Ehret, T., De Truchis, A., Mazzolini, M., Morel, J. M., D’aspremont, A., Lauvaux, T., ... & Facciolo, G. (2022). Global
957 tracking and quantification of oil and gas methane emissions from recurrent sentinel-2 imagery. *Environmental science &*
958 *technology*, 56(14), 10517-10529.
- 959 Frankenberg, C., Thorpe, A. K., Thompson, D. R., Hulley, G., Kort, E. A., Vance, N., Borchardt, J., Krings, T., Gerilowski,
960 K., Sweeney, C., and Conley, S.: Airborne methane remote measurements reveal heavy-tail flux distribution in Four Corners
961 region, *P. Natl. Acad. Sci. USA*, 113, 9734–9739, <https://doi.org/10.1073/pnas.1605617113>, 2016.
- 962 Hakkarainen, J., Ialongo, I., and Tamminen, J. (2016). Direct space-based observations of anthropogenic CO₂ emission areas
963 from OCO-2. *Geophysical Research Letters* 43, 11,400–11,406. doi:10.1002/2016GL070885
- 964 Hakkarainen, J., Ialongo, I., Koene, E., Szeląg, M., Tamminen, J., Kuhlmann, G., and Brunner, D. (2022). Analyzing local
965 carbon dioxide and nitrogen oxide emissions from space using the divergence method: An application to the synthetic
966 SMARTCARB dataset. *Frontiers in Remote Sensing* 3. doi:10.3389/frsen.2022.878731.
- 967 Hakkarainen, J., Ialongo, I., Oda, T., Szeląg, M. E., O’Dell, C. W., Eldering, A., and Crisp, D. (2023a). Building a bridge:
968 Characterizing major anthropogenic point sources in the South African Highveld region using OCO-3 carbon dioxide
969 Snapshot Area Maps and Sentinel-5P/TROPOMI nitrogen dioxide columns. *Environmental Research Letters*, 18(3),
970 doi:10.1088/1748-9326/acb837.
- 971 Hakkarainen, J., Tamminen, J., Nurmela, J., Lindqvist, H., Santaren, D., Broquet, G., Chevallier, F., Koene, E., Kuhlmann,
972 G. and Brunner, D. (2023b). Benchmarking of plume detection and quantification methods. Technical Report. FMI. URL:
973 <https://www.coco2-project.eu/node/366>. CoCO2: Prototype system for a Copernicus CO₂ service.
- 974 Hakkarainen, J., Kuhlmann, G., Koene, E., Santaren, D., Meier, S., Krol, M.C., van Stratum, B.J.H, Ialongo, I., Chevallier,
975 F., Tamminen, J., Brunner, D., Broquet, G. (2023c). Analyzing nitrogen dioxide to nitrogen oxide scaling factors for



976 computationally light satellite-based emission estimation methods: a case study of Matimba/Medupi power stations in South
977 Africa, Technical Note in review for Atmospheric Environment: X.

978 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., et al. (2020). The ERA5 global
979 reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 1, 51. <https://doi.org/10.1002/qj.3803>

980 Houweling, S., Aben, I., Breon, F.-M., Chevallier, F., Deutscher, N., Engelen, R., Gerbig, C., Griffith, D., Hungershofer,
981 K., Macatangay, R., Marshall, J., Notholt, J., Peters, W., and Serrar, S.: The importance of transport model uncertainties for
982 the estimation of CO₂ sources and sinks using satellite measurements, *Atmos. Chem. Phys.*, 10, 9981–
983 9992, <https://doi.org/10.5194/acp-10-9981-2010>, 2010.

984 Jacob, D. J. (1999). *Introduction to Atmospheric Chemistry* (Princeton University Press).

985 Jacob, D. J., Varon, D. J., Cusworth, D. H., Dennison, P. E., Frankenberg, C., Gautam, R., ... & Duren, R. M. (2022).
986 Quantifying methane emissions from the global scale down to point sources using satellite observations of atmospheric
987 methane. *Atmospheric Chemistry and Physics*, 22(14), 9617-9646.

988 Jähn, M., Kuhlmann, G., Mu, Q., Haussaire, J. M., Ochsner, D., Osterried, K., ... & Brunner, D. (2020). An online emission
989 module for atmospheric chemistry transport models: implementation in COSMO-GHG v5. 6a and COSMO-ART v5. 1-
990 3.1. *Geoscientific Model Development*, 13(5), 2379-2392.

991 Jacob, D. J. (1999). *Introduction to Atmospheric Chemistry* (Princeton University Press).

992 Jacob, D. J., Varon, D. J., Cusworth, D. H., Dennison, P. E., Frankenberg, C., Gautam, R., ... & Duren, R. M. (2022).
993 Quantifying methane emissions from the global scale down to point sources using satellite observations of atmospheric
994 methane. *Atmospheric Chemistry and Physics*, 22(14), 9617-9646.

995 Janssens-Maenhout, G., Pinty, B., Dowell, M., Zunker, H., Andersson, E., Balsamo, G., et al. (2020). Toward an Operational
996 Anthropogenic CO₂ Emissions Monitoring and Verification Support Capacity. *Bull. Am. Meteorol. Soc.* 101, E1439–E1451.
997 doi:10.1175/BAMS-D-19-0017.1

998 Kasahara, M., Kachi, M., Inaoka, K., Fujii, H., Kubota, T., Shimada, R., & Kojima, Y. (2020, September). Overview and
999 current status of GOSAT-GW mission and AMSR3 instrument. In *Sensors, Systems, and Next-Generation Satellites XXIV*
1000 (Vol. 11530, p. 1153007). SPIE.

1001 Koene, E., Brunner, D. and Kuhlmann, G. (2021). Documentation of plume detection and quantification methods. Tech. rep.,
1002 Empa. CoCO₂: Prototype system for a Copernicus CO₂ service. <https://coco2-project.eu/node/329>

1003 Koene, E., Brunner, D., 2023. Assessment of plume model performance. Technical Report. Empa. URL: [https://www.coco2-](https://www.coco2-project.eu/node/357)
1004 [project.eu/node/357](https://www.coco2-project.eu/node/357). CoCO₂: Prototype system for a Copernicus CO₂ service.

1005 Kort, E. A., Frankenberg, C., Miller, C. E., and Oda, T.: Space-based observations of megacity carbon dioxide, *Geophys.*
1006 *Res. Lett.*, 39, L17806, <https://doi.org/10.1029/2012gl052738>, 2012.

1007 Kuenen, J. J. P., Visschedijk, A. J. H., Jozwicka, M., and Denier van der Gon, H. A. C.: TNO-MACC_II emission inventory;
1008 a multi-year (2003–2009) consistent high-resolution European emission inventory for air quality modelling, *Atmos. Chem.*
1009 *Phys.*, 14, 10963–10976, <https://doi.org/10.5194/acp-14-10963-2014>, 2014.



- 1010 Kuhlmann, G., Broquet, G., Marshall, J., Clément, V., Löscher, A., Meijer, Y., et al. (2019). Detectability of CO₂ emission
1011 plumes of cities and power plants with the Copernicus Anthropogenic CO₂ Monitoring (CO₂M) mission. *Atmospheric*
1012 *Measurement Techniques* 12, 6695–6719. doi:10.5194/amt-12-6695-2019.
- 1013 Kuhlmann, G., Brunner, D., Broquet, G., and Meijer, Y. (2020). Quantifying CO₂ emissions of a city with the Copernicus
1014 Anthropogenic CO₂ Monitoring satellite mission. *Atmospheric Measurement Techniques* 13, 6733–6754. doi:10.5194/amt-
1015 13-6733-2020.
- 1016 Kuhlmann, G., Clément, V., Marshall, J., Fuhrer, O., Broquet, G., Schnadt-Poberaj, C., et al. (2020b). Synthetic XCO₂, CO
1017 and NO₂ Observations for the CO₂M and Sentinel-5 Satellites. doi:10.5281/zenodo.4048228
- 1018 Kuhlmann, G., Henne, S., Meijer, Y., and Brunner, D. (2021). Quantifying CO₂ Emissions of Power Plants With CO₂ and
1019 NO₂ Imaging Satellites. *Frontiers in Remote Sensing* 2, 14. doi:10.3389/frsen. 2021.
- 1020 Kuhlmann, G., Koene, E., Meier, S., Brunner, D., Santaren, D., Broquet, G., Chevallier, F., Hakkarainen, J., Nurmela, J.,
1021 Amoros, L., and Tamminen, J.. The ddeq Python library for point source quantification from remote sensing images
1022 (Version 1.0). *To be submitted to GMD as a model description paper*.
- 1023 Landgraf, J., Rusli, S., Cooney, R., Veeffkind, P., Vemmix, T., de Groot, Z., Bell, A., Day, J., Leemhuis, A., and Sierk, B.:
1024 The TANGO mission: A satellite tandem to measure major sources of anthropogenic greenhouse gas emissions, EGU
1025 General Assembly 2020, Online, 4–8 May 2020, EGU2020-19643, <https://doi.org/10.5194/egusphere-egu2020-19643>, 2020.
- 1026 Lary, D. J., Alavi, A. H., Gandomi, A. H., & Walker, A. L. (2016). Machine learning in geosciences and remote
1027 sensing. *Geoscience Frontiers*, 7(1), 3-10.
- 1028 Mahadevan, P., Wofsy, S. C., Matross, D. M., Xiao, X., Dunn, A. L., Lin, J. C., ... & Gottlieb, E. W. (2008). A
1029 satellite-based biosphere parameterization for net ecosystem CO₂ exchange: Vegetation Photosynthesis and Respiration
1030 Model (VPRM). *Global Biogeochemical Cycles*, 22(2).
- 1031 Meijer, Y., Boesch, H., Bombelli, A., Brunner, D., Buchwitz, M., Ciais, P., et al. (2019). Copernicus CO₂ monitoring
1032 mission Requirements document (MRD). Netherlands, Europe: European Space Agency, Earth and Mission Science
1033 Division.
- 1034 Nassar, R., Hill, T. G., McLinden, C. A., Wunch, D., Jones, D. B. A., and Crisp, D. (2017). Quantifying CO₂ emissions from
1035 individual power plants from space. *Geophys. Res. Lett.* 44, 10045-10053. doi:10.1002/2017GL074702
- 1036 Nassar R, Moeini O, Mastrogiacomo J-P, O'Dell CW, Nelson RR, Kiel M, Chatterjee A, Eldering A and Crisp D (2022),
1037 Tracking CO₂ emission reductions from space: A case study at Europe's largest fossil fuel power plant. *Front. Remote Sens.*
1038 3:1028240. doi: 10.3389/frsen.2022.1028240
- 1039 Pascal, V., Buil, C., Loesel, J., Tauziede, L., Jouglet, D., & Buisson, F. (2017, November). An improved microcarb
1040 dispersive instrumental concept for the measurement of greenhouse gases concentration in the atmosphere. In *International*
1041 *Conference on Space Optics—ICSO 2014* (Vol. 10563, pp. 1028-1036). SPIE.



- 1042 Pillai, D., Buchwitz, M., Gerbig, C., Koch, T., Reuter, M., Bovensmann, H., et al. (2016). Tracking City CO₂ Emissions
1043 from Space Using a High-Resolution Inverse Modelling Approach: a Case Study for Berlin, Germany. *Atmos. Chem. Phys.*
1044 16, 9591–9610. doi:10.5194/acp-16-9591-2016
- 1045 Pinty, B., Janssens-Maenhout, G., Dowell, M., Zunker, H., Brunhes, T., Ciais, P., Dee, D., Denier van der Gon, H. A. C.,
1046 Dolman, H., Drinkwater, M., Engelen, R., Heimann, M., Holmlund, K., Husband, R., Kentarchos, A., Meyer, A., Palmer, P.,
1047 and Scholze, M.: An operational anthropogenic CO₂ emissions monitoring and verification support capacity. Baseline
1048 requirements, model components and functional architecture, EUR28736 EN, European Commission Joint Research Centre,
1049 Ispra, Italy, <https://doi.org/10.2760/08644>, 2017.
- 1050 Reuter, M., Buchwitz, M., Schneising, O., Krautwurst, S., O'Dell, C. W., Richter, A., et al. (2019). Towards Monitoring
1051 Localized CO₂ Emissions from Space: collocated Regional CO₂ and NO₂ Enhancements Observed by the OCO-2 and S5P
1052 Satellites. *Atmos. Chem. Phys.* 19, 9371–9383. doi:10.5194/acp-19-9371-2019
- 1053 Santaren, D., Broquet, G., Bréon, F.-M., Chevallier, F., Siméoni, D., Zheng, B., and Ciais, P.: A local- to national-scale
1054 inverse modeling system to assess the potential of spaceborne CO₂ measurements for the monitoring of anthropogenic
1055 emissions, *Atmos. Meas. Tech.*, 14, 403–433, <https://doi.org/10.5194/amt-14-403-2021>, 2021.
- 1056 Sierk, B., Bézy, J.-L., Löscher, A., and Meijer, Y. (2019). The European CO₂ Monitoring Mission: Observing
1057 Anthropogenic Greenhouse Gas Emissions from Space 11180. Proceedings, International Conference on Space Optics—
1058 ICSO 2018. 12 July 2019. Chania, Greece. 111800M. doi:10.1117/12.2535941
- 1059 Singer, A.M., Branham, M., Hutchins, M.G., Welker, J., Woodard, D. L., Badurek, C. A., et al. (2014). The role of CO₂
1060 emissions from large point sources in emissions totals, responsibility and policy. *Environ. Sci. Policy* 44, 190–200.
1061 doi:10.1016/j.envsci.2014.08.001
- 1062 Taylor, T. E., O'Dell, C. W., Frankenberg, C., Partain, P. T., Cronk, H. Q., Savtchenko, A., Nelson, R. R., Rosenthal, E. J.,
1063 Chang, A. Y., Fisher, B., Osterman, G. B., Pollock, R. H., Crisp, D., Eldering, A., and Gunson, M. R.: Orbiting Carbon
1064 Observatory-2 (OCO-2) cloud screening algorithms: validation against collocated MODIS and CALIOP data, *Atmos. Meas.*
1065 *Tech.*, 9, 973–989, <https://doi.org/10.5194/amt-9-973-2016>, 2016.
- 1066 Van Heerwaarden, C. C., Van Stratum, B. J., Heus, T., Gibbs, J. A., Fedorovich, E., & Mellado, J. P. (2017). MicroHH 1.0:
1067 A computational fluid dynamics code for direct numerical simulation and large-eddy simulation of atmospheric boundary
1068 layer flows. *Geoscientific Model Development*, 10(8), 3145–3165.
- 1069 Varon, D. J., Jacob, D. J., McKeever, J., Jervis, D., Durak, B. O. A., Xia, Y., et al. (2018). Quantifying methane point
1070 sources from fine-scale satellite observations of atmospheric methane plumes. *Atmospheric Measurement Techniques* 11,
1071 5673–5686. doi:10.5194/amt-11-5673-2018.
- 1072 Wang, Y., Broquet, G., Bréon, F.-M., Lespinas, F., Buchwitz, M., Reuter, M., et al. (2020). PMIF v1.0: Assessing the
1073 Potential of Satellite Observations to Constrain CO₂ Emissions from Large Cities and point Sources over the globe Using
1074 Synthetic Data. *Geosci. Model. Dev.* 13, 5813–5831. doi:10.5194/gmd-13-5813-2020



1075 Worden, J. R., Doran, G., Kulawik, S., Eldering, A., Crisp, D., Frankenberg, C., O'Dell, C., and Bowman, K.: Evaluation
1076 and attribution of OCO-2 XCO₂ uncertainties, *Atmos. Meas. Tech.*, 10, 2759–2771, [https://doi.org/10.5194/amt-10-2759-](https://doi.org/10.5194/amt-10-2759-2017)
1077 2017, 2017.

1078 Ye, X., Lauvaux, T., Kort, E., Oda, T., Feng, S., Lin, J., Yang, E., & Wu, D. (2020). Constraining Fossil Fuel CO₂ Emissions
1079 From Urban Area Using OCO-2 Observations of Total Column CO₂. *Journal of Geophysical Research: Atmospheres*, 1-29.

1080 Zheng, T., Nassar, R., and Baxter, M. (2019). Estimating power plant CO₂ emission using OCO-2 XCO₂ and high resolution
1081 WRF-Chem simulations. *Environ. Res. Lett.* 14, 085001. doi:10.1088/1748-9326/ab25ae

1082 Zheng, B., Chevallier, F., Ciais, P., Broquet, G., Wang, Y., Lian, J., et al. (2020). Observing Carbon Dioxide Emissions over
1083 China's Cities and Industrial Areas with the Orbiting Carbon Observatory-2. *Atmos. Chem. Phys.* 20, 8501–8510.
1084 doi:10.5194/acp-20-8501-2020

1085

1086

1087

1088

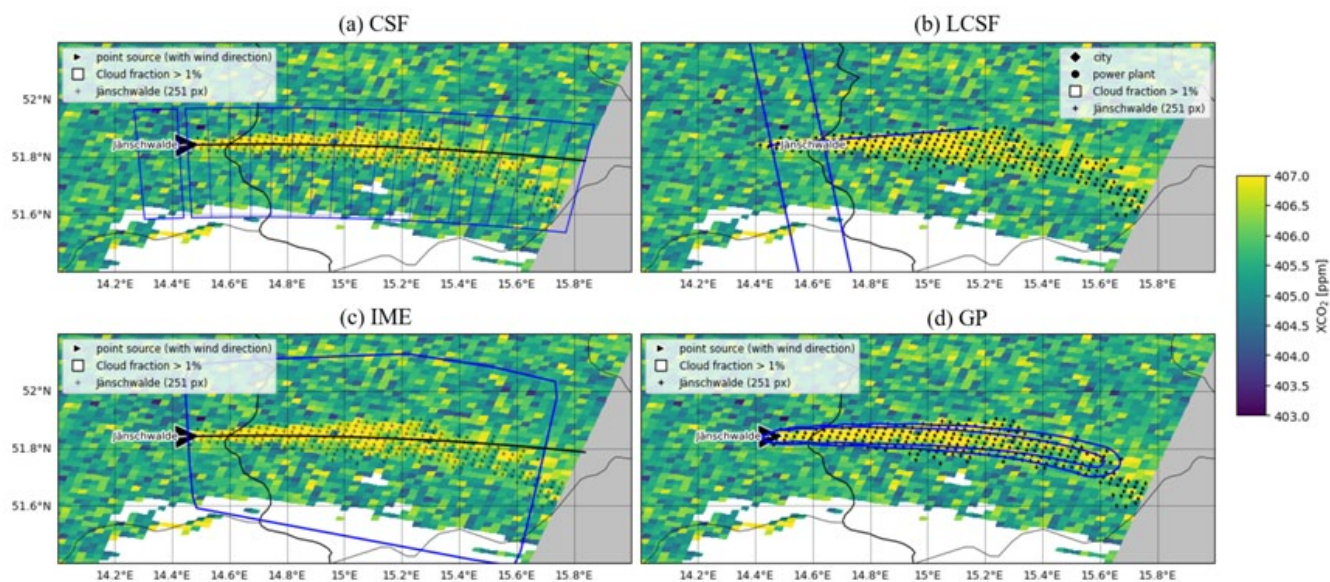
1089

1090

1091

1092

1093



1094

1095

1096

1097

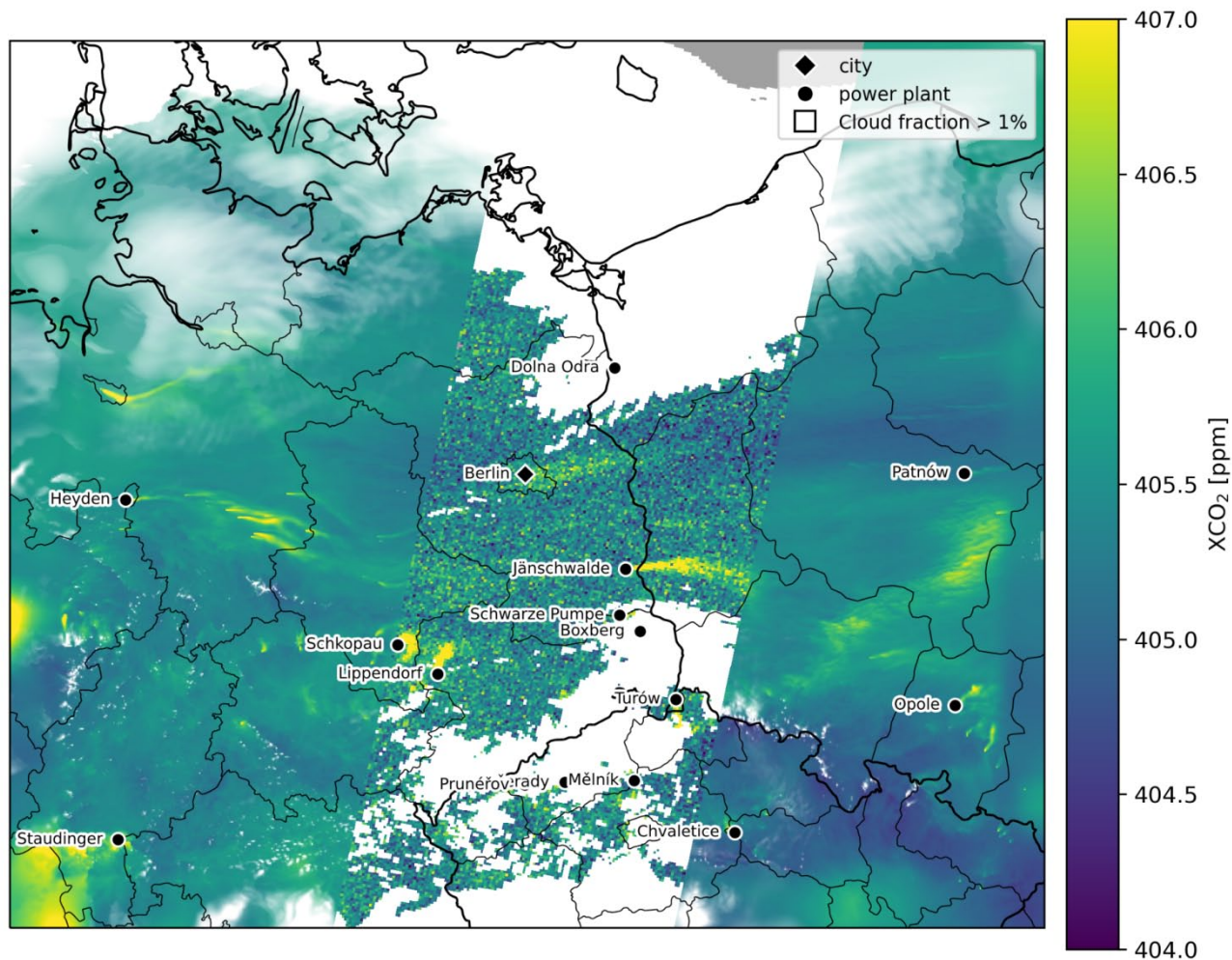
1098

1099

1100

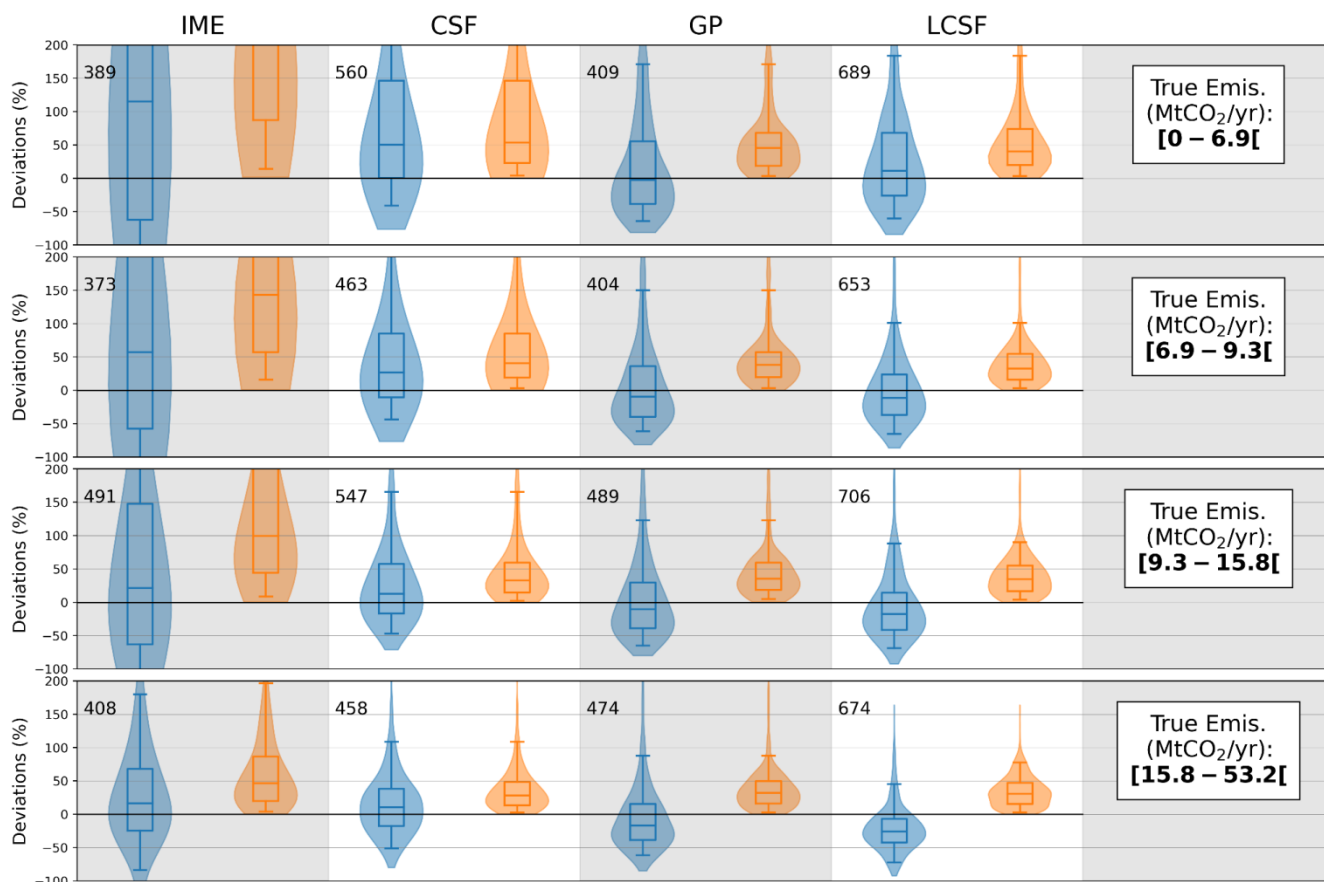
Figure 1: Illustration of different inversion methods for a plume produced by the Jämschwalde power plant on April 23rd, 2015. For all figures, pixels with dots are the selected enhancements representing the plume a) CSF method: the blue boxes depict the areas where the Gaussian fits of the plume cross-sections are made and the black line the centre-line of the plume. b) LCSF method: the blue lines represent the domain where the Gaussian fits of the plume cross-sections are made and the black line the along-wind direction at the source. c) IME method: the blue curve represents the domain on which mass enhancements are integrated. d) GP method: Blue curves depict contour lines of the 2-dimensional Gaussian curve that fits the plume.

1101



1102
1103
1104
1105
1106

Figure 2. Simulations of XCO₂ on 23 April 2015 over the SMARTCARB domain. Synthetic XCO₂ observations over a 250 km wide swath are represented in the centre of the figure for a low noise scenario. Missing XCO₂ observations due to a cloud fraction larger than 1% are shown in white. The 16 emission sources considered in this study are highlighted along with their names



1107

1108

1109

1110

1111

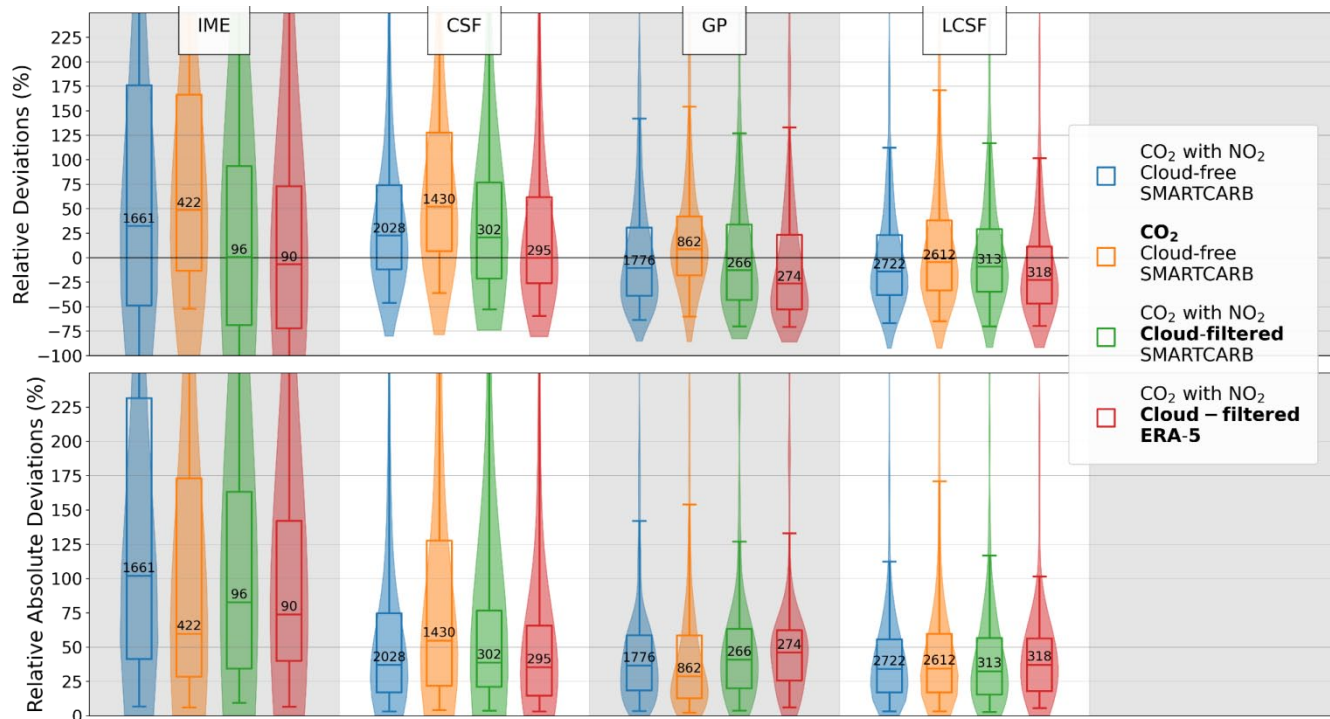
1112

1113

Figure 3. Performance when estimating CO₂ emissions from individual images of the different single-image inversion methods (columns) across different ranges of true emissions (rows) using SMARTCARB winds and cloud-free CO₂ and NO₂ data. The distributions of relative deviations (in blue) and relative absolute deviations (in orange) are illustrated using violin plots. The inter-quartiles are represented by the boxes, while the whiskers indicate the 5th and 95th percentiles, and medians are the lines inside the boxes. The numbers alongside boxes show the numbers of estimates corresponding to true emissions ranges and inversion methods.

1114

1115



1116

1117

1118

1119

1120

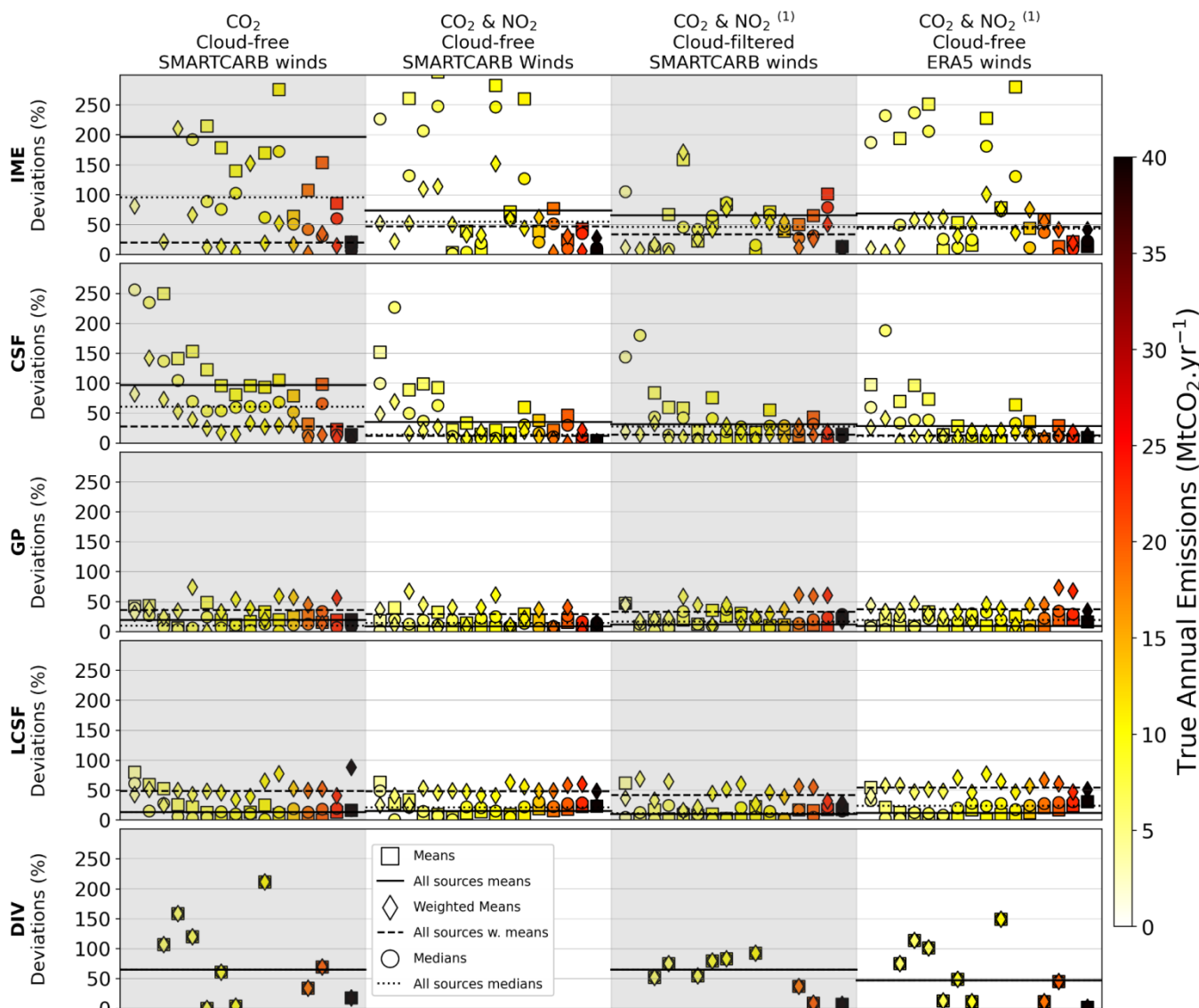
1121

1122

1123

1124

Figure 4. Performances of the inversion methods when estimating emissions from single images for different benchmarking scenarios: cloud-free CO₂ and NO₂ data with SMARTCARB winds (in blue), cloud-free CO₂ data only with SMARTCARB winds (in orange), cloud-filtered CO₂ and NO₂ data with SMARTCARB winds (in green), cloud-filtered CO₂ and NO₂ data with ERA5 winds (in red). Bold texts in the legend indicate the elements of benchmarking scenarios that differ from those in the ideal benchmarking scenario. Distributions of the relative deviations (top panel) and relative absolute deviations (bottom panel) are illustrated using violin plots. Boxes are the inter-quartiles of the distributions, the whiskers are the 5th and 95th percentiles, and the lines within boxes are the medians.



1125

1126

1127

1128

1129

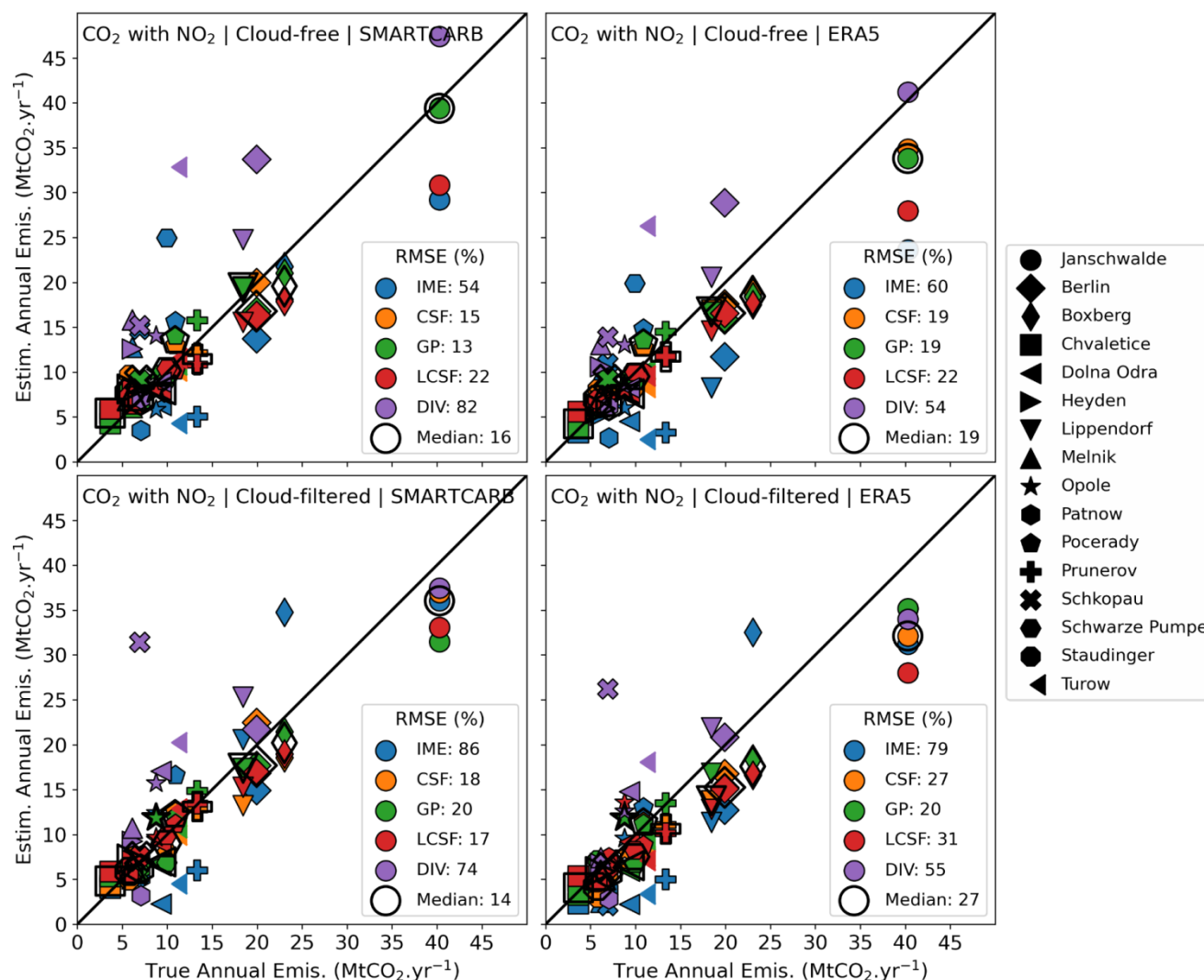
1130

1131

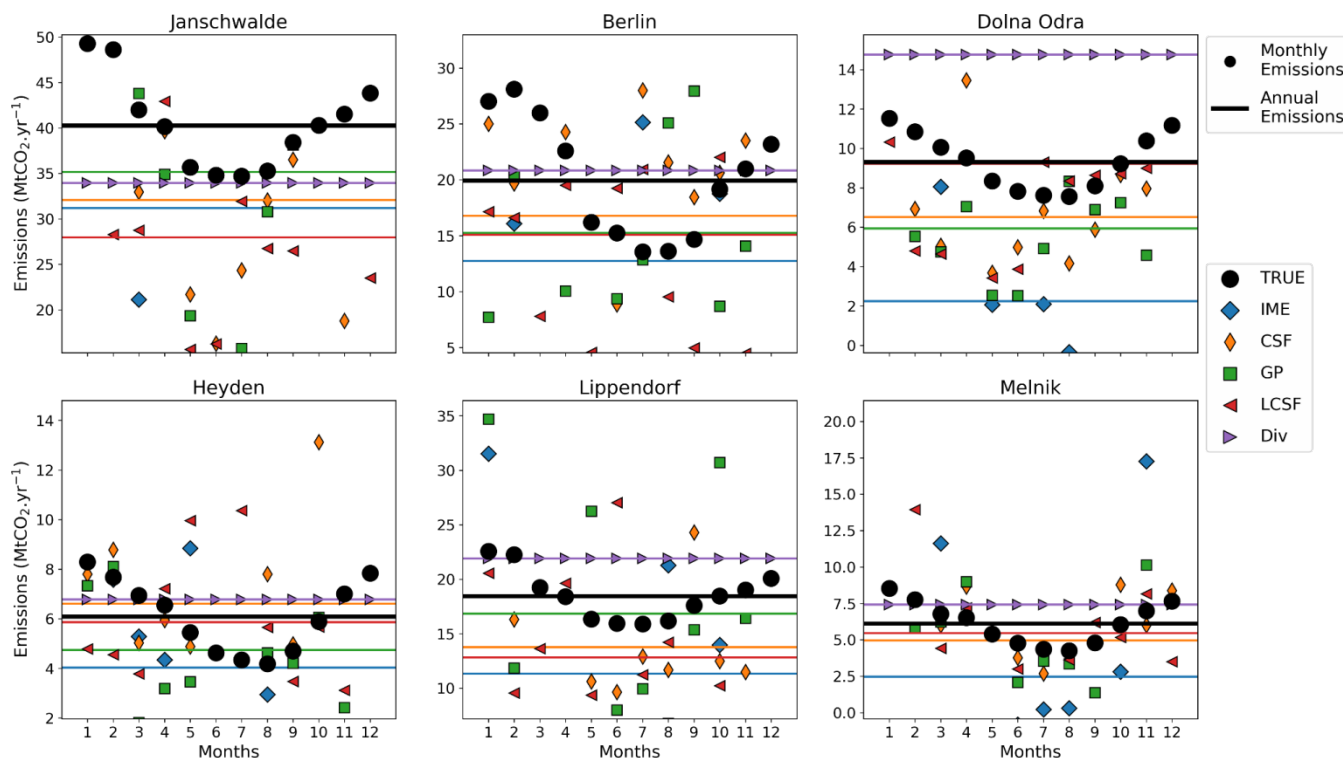
1132

Figure 5. Performance of the inversion methods for annual estimates. The markers represent for a given source the relative absolute deviations from the true annual emissions of the arithmetic means (squares), the weighted means (diamonds) and the medians (circles) of the estimates over a year. The lines represent the median values of the annual estimates over the entire set of sources. The inversions are performed using CO₂ cloud-free data and SMARTCARB winds (1st column), using CO₂ and NO₂ cloud-free data and with SMARTCARB winds (2nd column), using CO₂ and NO₂ cloud-filtered data and SMARTCARB winds (3rd column), and using CO₂ and NO₂ cloud-free data and with ERA5 winds (4th column). (1) For the Divergence methods, the inversions of the 3rd and 4th columns are performed using CO₂ data only.

1133



1134
 1135 **Figure 6. Estimated vs true annual emissions for 4 inversion scenarios (titles of the panels). For the IME and CSF methods, annual**
 1136 **estimates are weighted means of the single-image estimates while they are arithmetic means for the GP, LCSF and Divs methods.**
 1137 **Each marker represents a given emission source and each color a given inversion method. The unfilled markers represent the**
 1138 **median values of all the estimates for each source. The divergence inversion method uses CO₂ data for all the inversion scenarios.**
 1139 **The plain line represents the 1:1 line. The bottom-right legends display for each inversion method the relative RMSE which is the**
 1140 **RMSE between estimated and true annual emissions divided by the median of true annual emissions of all sources (~9.6 MtCO₂.yr⁻¹**
 1141 **¹).**
 1142
 1143



1144

1145

1146

1147

1148

1149

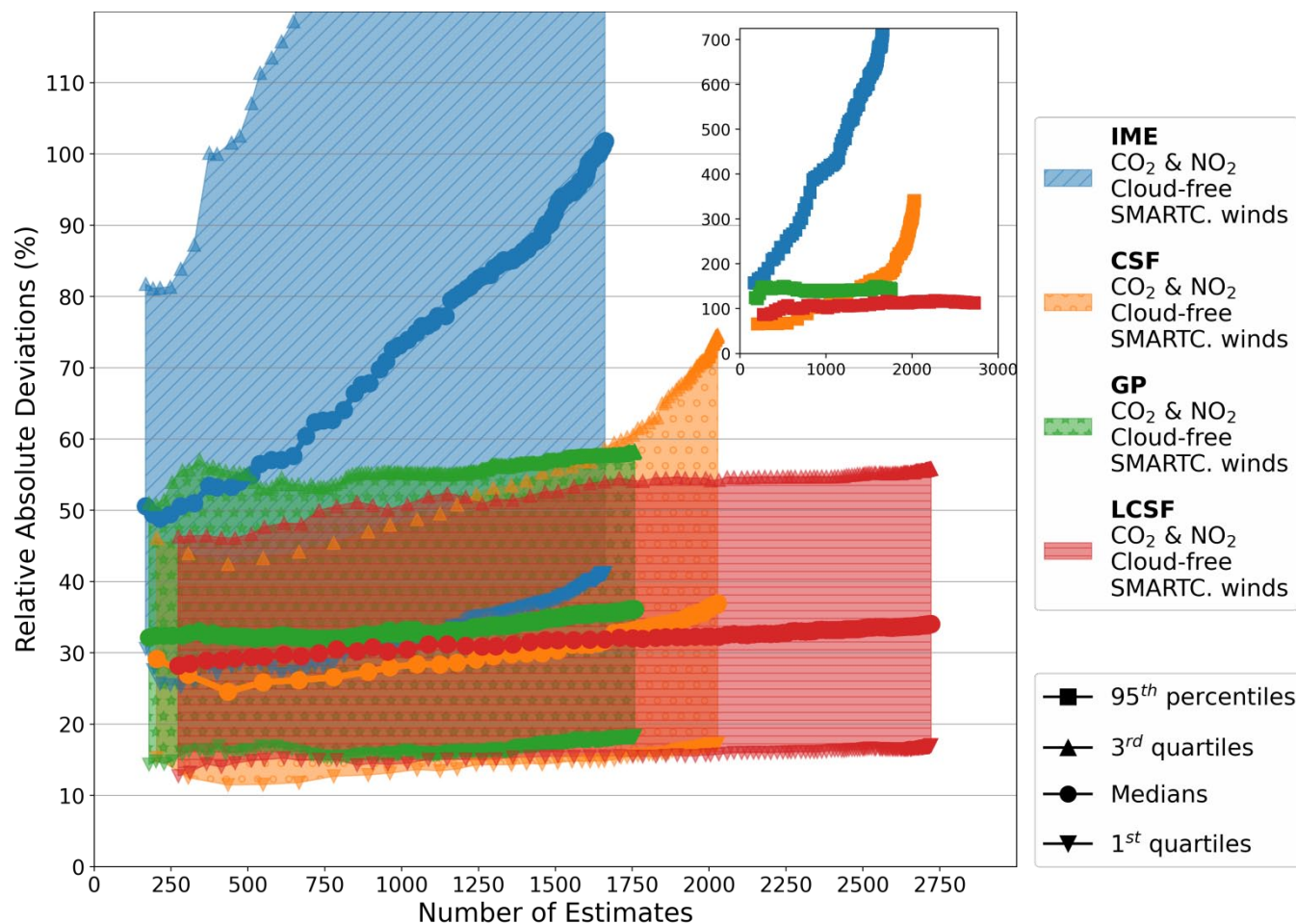
1150

1151

1152

1153

Figure 7. Annual and monthly estimates of the true and estimated emissions for different sources and for different inversion methods. Each panel is associated with a given source. Plain lines and markers represent annual averages and monthly averages respectively. Colors and markers are associated with different inversion methods (true emissions are represented by black circles). Annual and monthly estimates for the IME and CSF methods are weighted means of image estimates. Annual and monthly estimates for the GP and LCSF are means of image estimates while for the divergence method, we use the annual estimate also for monthly estimates. All inversion methods use CO₂ and NO₂ cloud-filtered data (CO₂ data only for the Div method) with ERA5 winds.

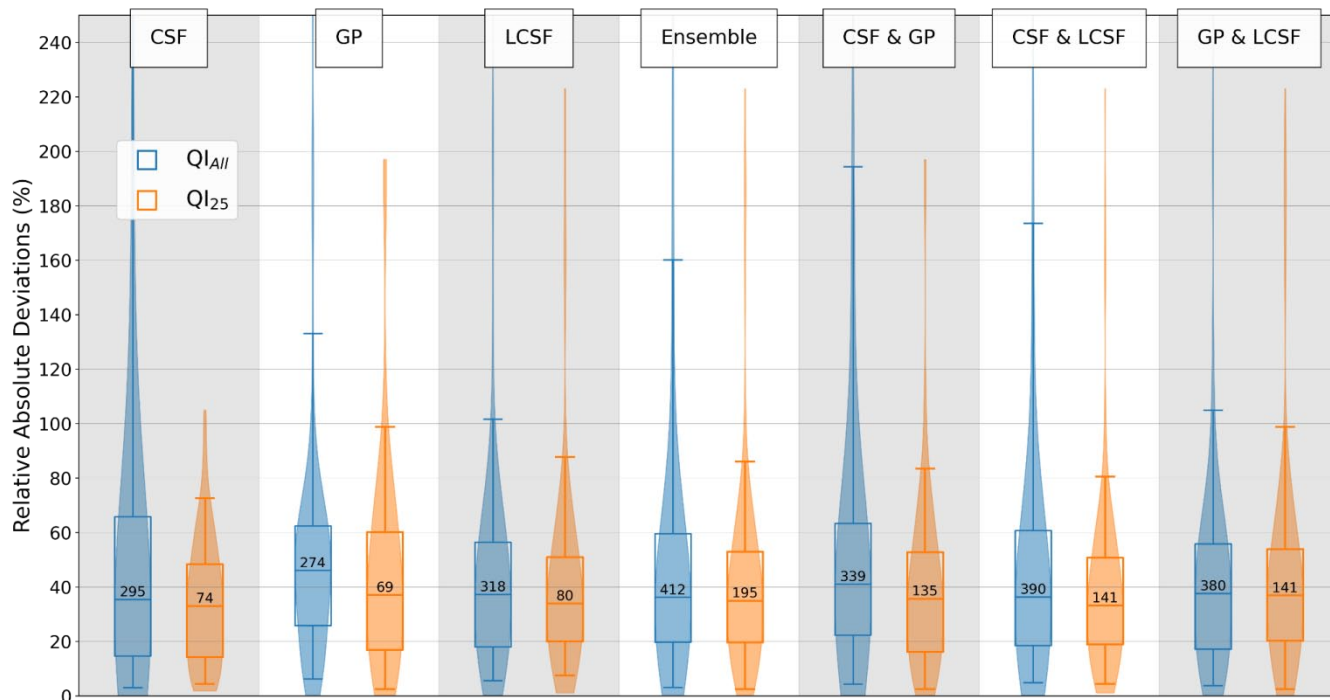


1154

1155 **Figure 8.** Accuracy of inversions vs number of single-image estimates. The inversion methods shown here use CO₂ and NO₂ cloud-
 1156 free data and SMARTCARB winds. The filled areas represent the inter-quartiles of the distributions of the relative absolute
 1157 deviations depending on the number of estimates. The 90th percentiles of the distributions are represented in the inset. Points
 1158 belonging to a same curve are associated to different QIs and from left to right along curves, points are associated with a
 1159 decreasing QI; the points at the left and right ends of the curves are associated with the maximal and minimal QIs respectively.

1160

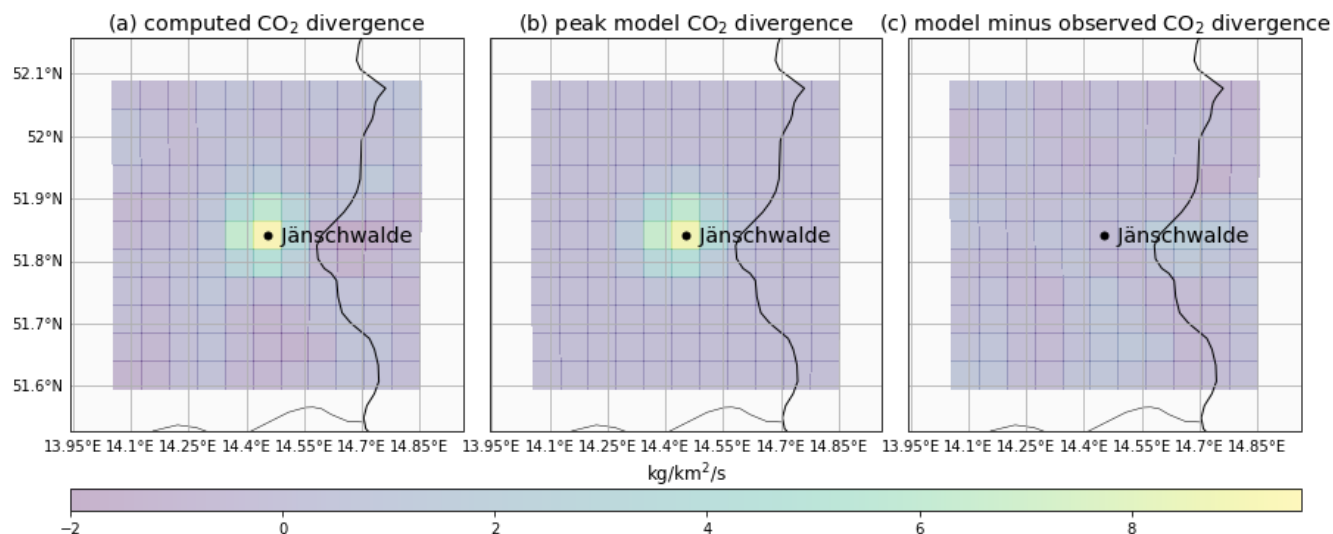
1161



1162

1163 **Figure 9:** Performance of the inversion methods and ensemble approaches for estimating the emissions with cloud-filtered CO₂ &
 1164 NO₂ data and with ERA5 winds. The distributions of the relative absolute deviations for all the inversion results (in blue) and for
 1165 the best estimates (in orange) provided by each method (see text) are illustrated using violin plots. Boxes represent the inter-
 1166 quartiles of the distributions, the whiskers the 5th and 95th percentiles, and the lines within boxes the medians.

1167

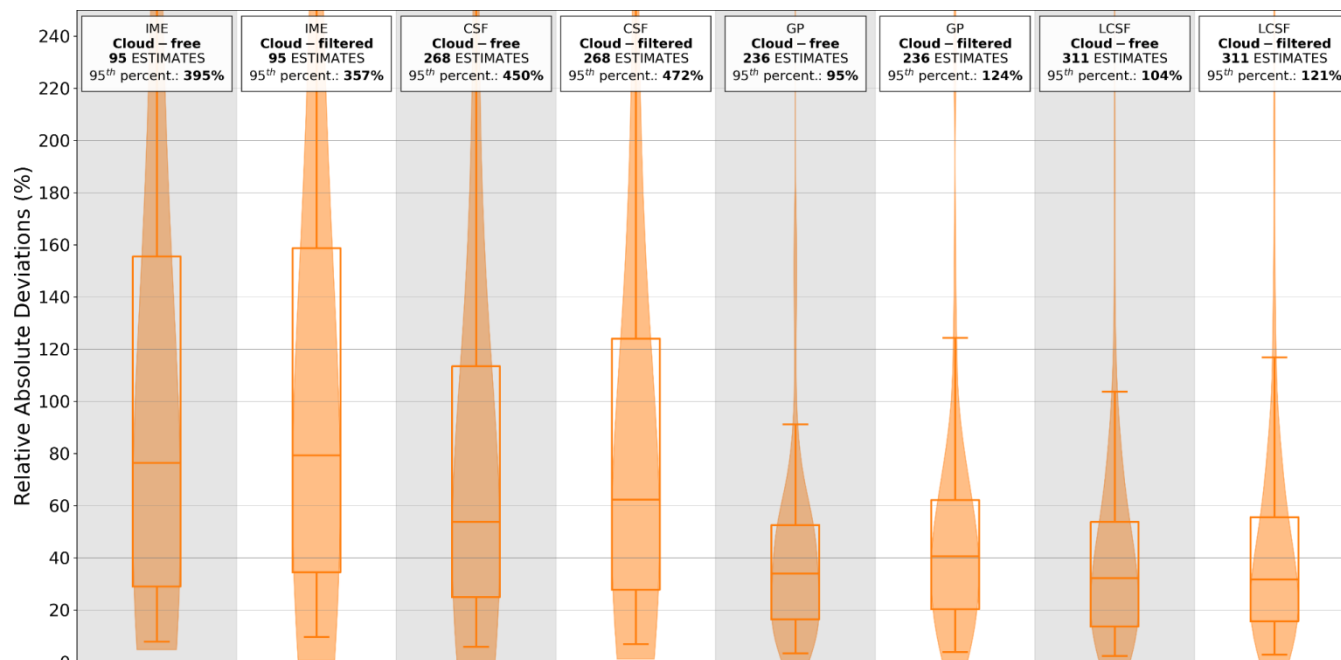


1168

1169 **Figure A1:** Illustration of the divergence method for the Jämschwalde power station in 2015 based on the synthetic SMARTCARB
 1170 dataset (see text). The figures represent the annual fields of the computed CO₂ divergence (a), the modeled CO₂ divergence (b) and
 1171 the difference of both quantities (c). Of note that as sink terms are considered negligible for CO₂, divergence fields are considered
 1172 equal to the emission fields for CO₂.



1173



1174

1175

1176

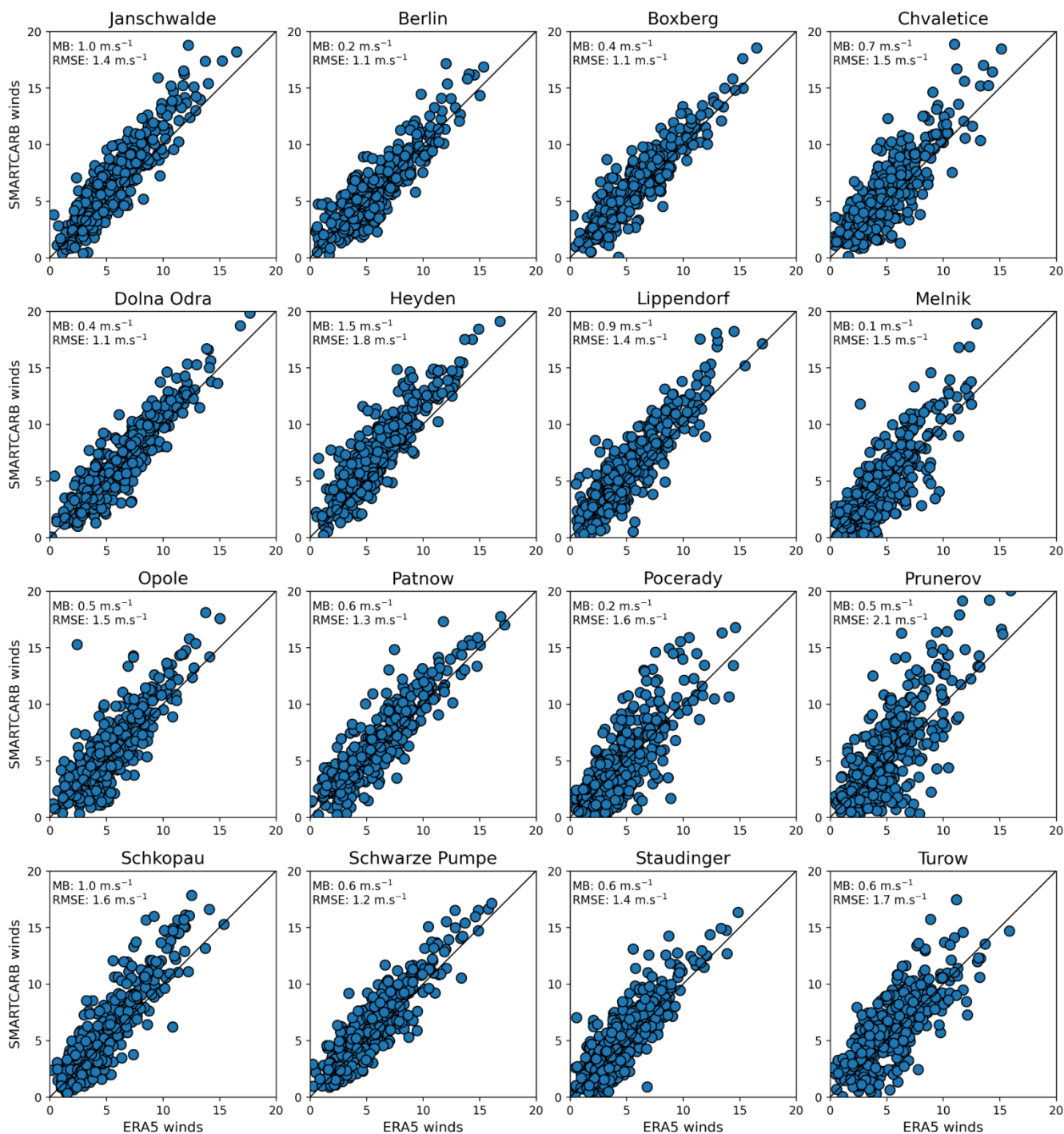
1177

1178

Figure A2: Performance of the inversion methods when using data with or without clouds for the emissions estimated from the same images. The inversion methods use CO₂ and NO₂ data and SMARTCARB winds. The boxes represent the inter-quartiles of the distributions of the absolute relative deviations, the whiskers the 5th and 95th percentiles, and the lines within boxes the medians.



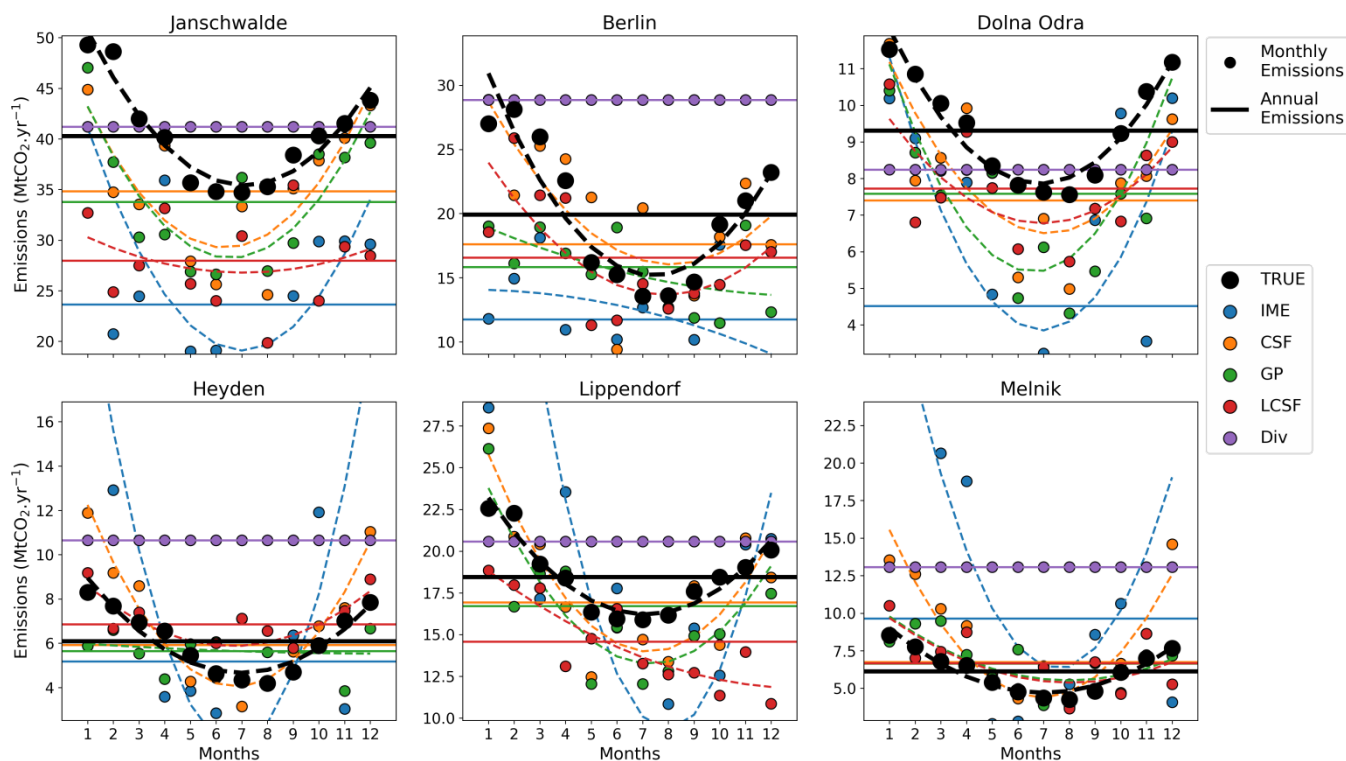
1179



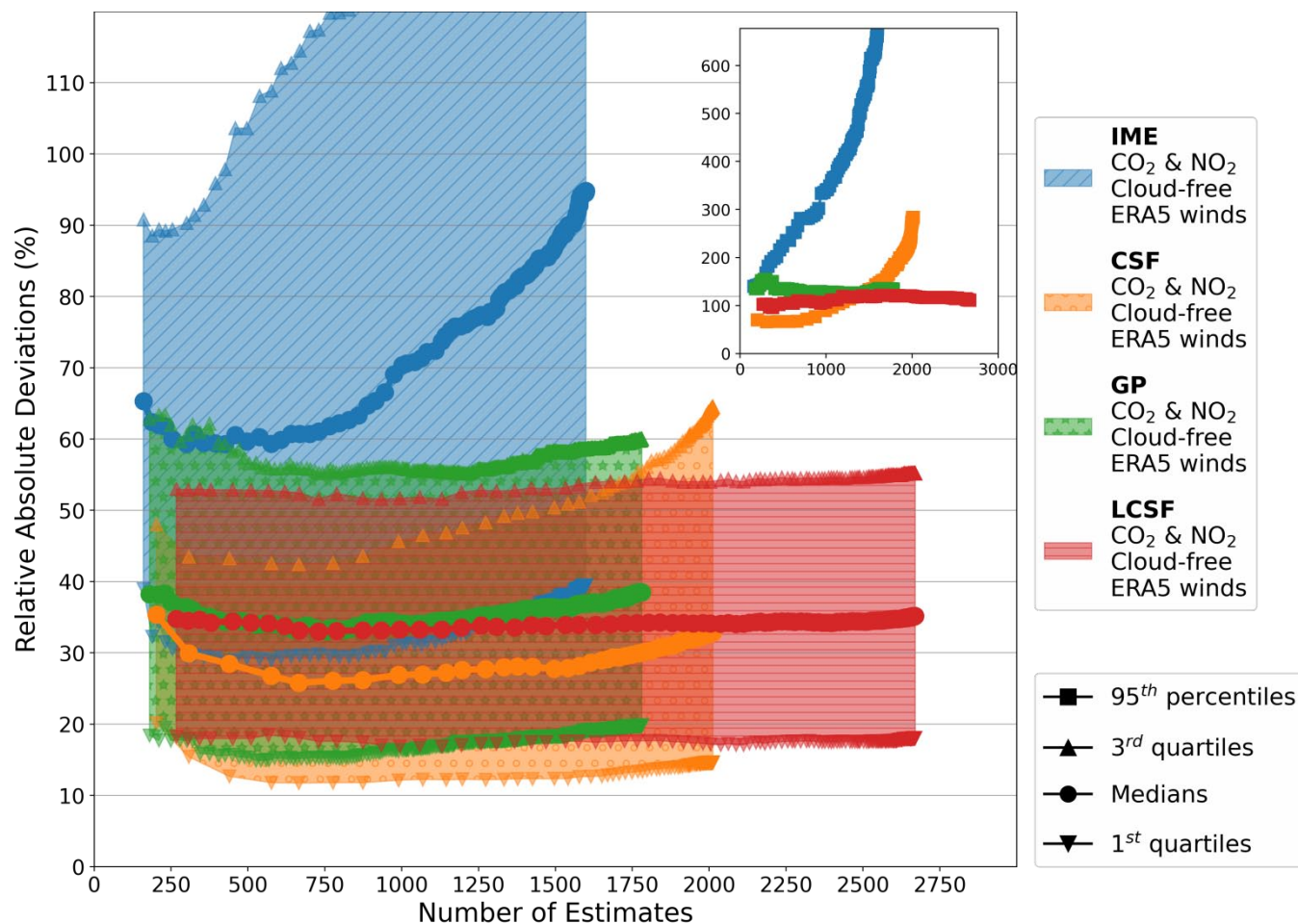
1180



1181 **Figure A3: Norms of the ERA5 winds vs norms of the SMARTCARB winds at the sources considered in this study and for all the**
 1182 **days of 2015. Black lines represent the 1:1 agreement line. Mean biases of the SMARTCARB norms minus the ERA5 norms and**
 1183 **RMSEs are noted at the top left of the figures.**



1184 **Figure A4: Annual and monthly estimates of the true and estimated emissions for different sources and for different inversion**
 1185 **methods. Each panel is associated with a given source. Plain lines and markers represent annual averages and monthly averages**
 1186 **respectively. Dashed lines represent the fits by a 2nd order polynomial of the monthly estimates. Colours are associated with**
 1187 **different inversion methods (true emissions are in black). Annual and monthly estimates for the IME and CSF methods are**
 1188 **weighted means of image estimates. Annual and monthly estimates for the GP and LCSF are means of image estimates while for**
 1189 **the divergence method, we use the annual estimate also for monthly estimates. All inversion methods use CO₂ and NO₂ cloud-free**
 1190 **data (CO₂ data only for the Divs methods) with ERA5 winds.**
 1191



1192

1193

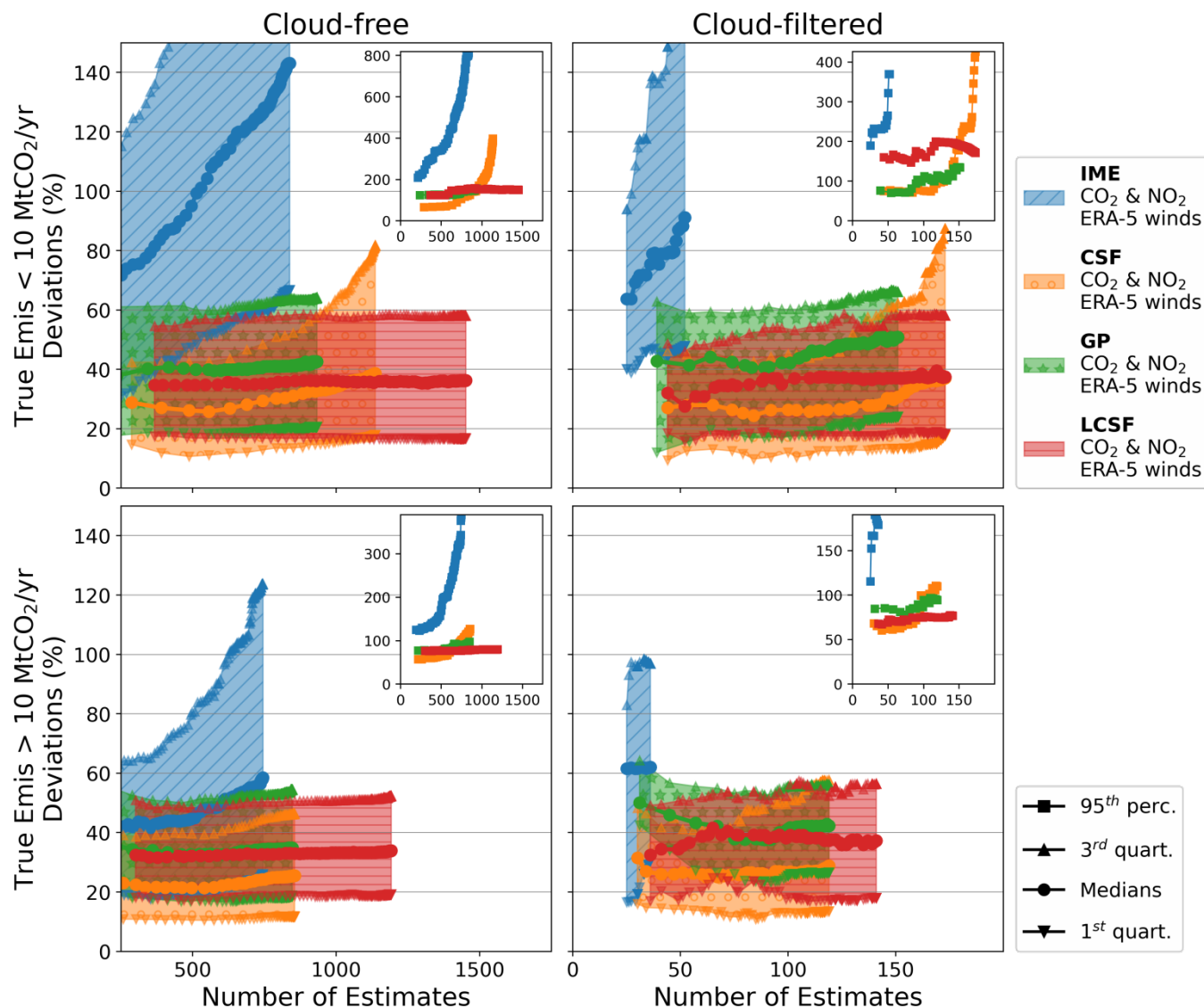
1194

1195

1196

1197

Figure A5. Accuracy of inversions vs number of instant estimates. The inversion methods shown here use CO₂ and NO₂ cloud-free data and ERA5 winds. The filled areas represent the inter-quartiles of the distributions of the relative absolute deviations depending on the number of estimates. The 90th percentiles of the distributions are represented in the inset. Points belonging to a same curve are associated to different QIs and from left to right along curves, points are associated with a decreasing QI; the points at the left and right ends of the curves are associated with the maximal and minimal QIs respectively.



1198

1199

1200

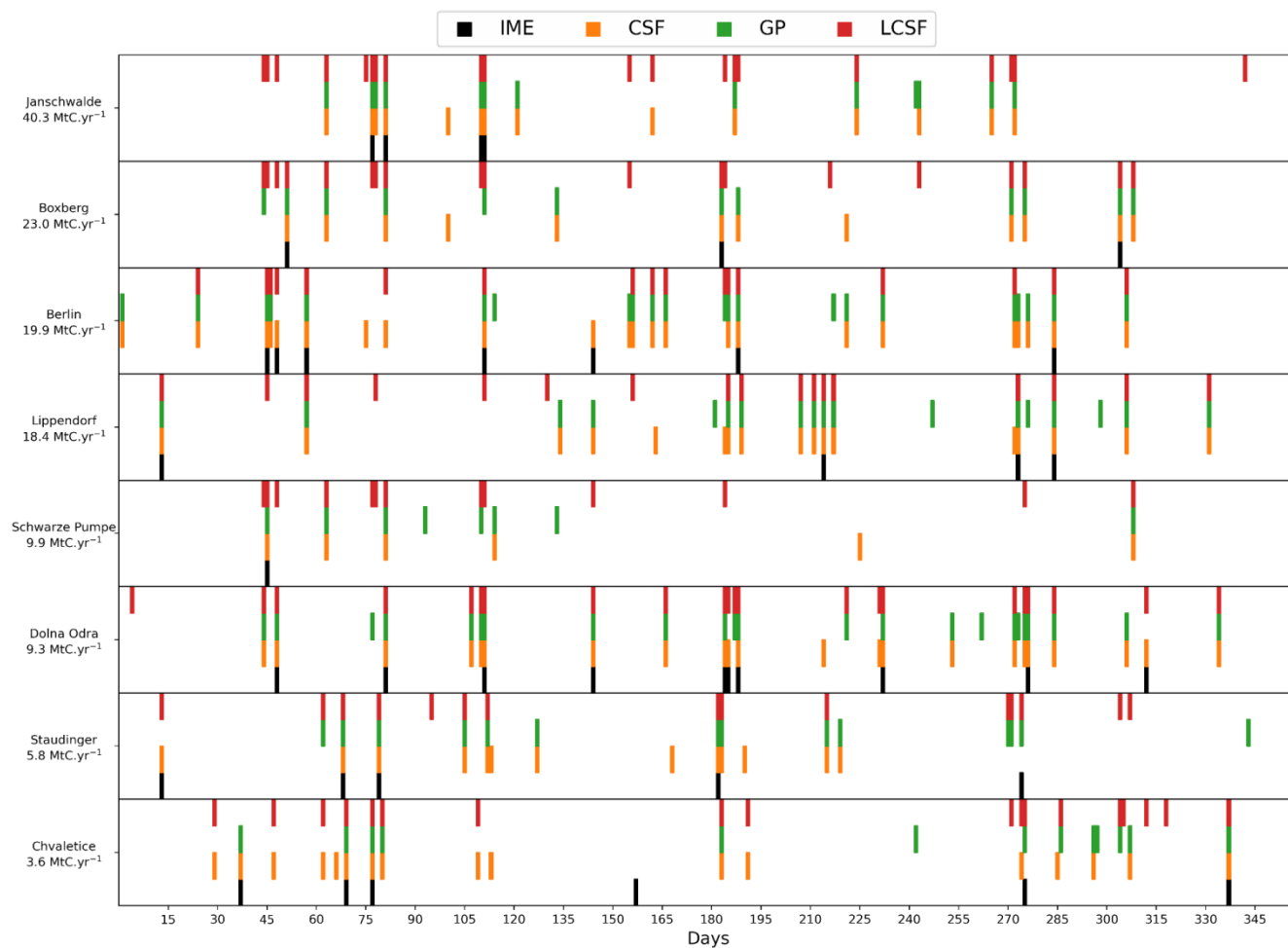
1201

1202

1203

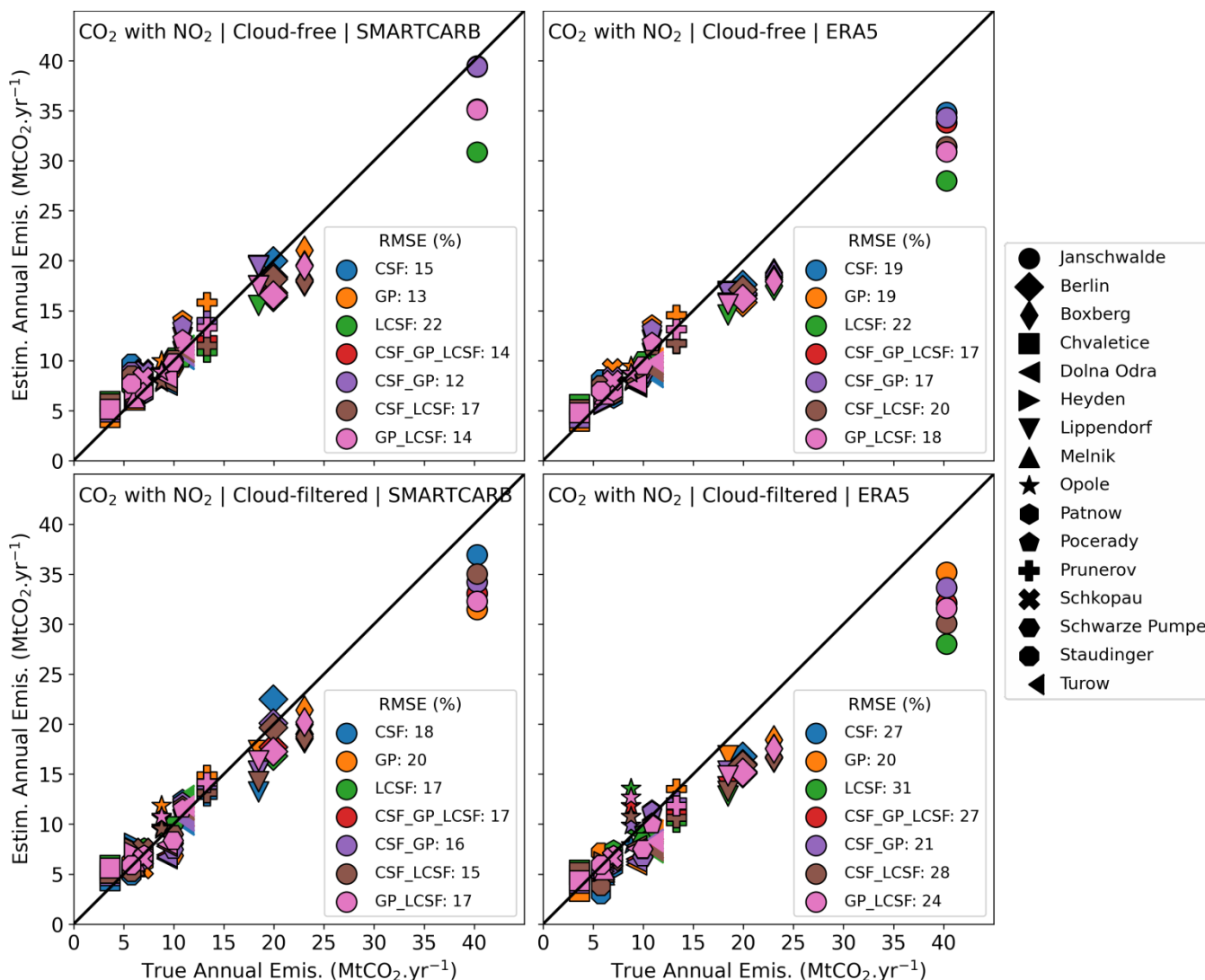
1204

Figure A6: Accuracy of inversions vs number of instant estimates. The inversion methods shown here use CO₂ and NO₂ data, ERA5 winds and for cloud-free (1st column) and cloud-filtered data (2nd column). Results are shown for the cases where true emissions of sources are below (1st row) and above (2nd row) 10 MtCO₂ yr⁻¹. The filled areas represent the inter-quartiles of the distributions of the relative absolute deviations depending on the number of estimates. The 90th percentiles of the distributions are represented in the insets. Each point belonging to a same curve is associated with a different QI and from left to right along a same curve; points are associated with a decreasing QI.



1205
1206
1207
1208

Figure A7: Days of 2015 (x-axis) for which the IME, CSF, GP and LCSF methods produce estimates for the emissions of eight sources (y-axis). For a given day, the availability of an estimate from a given inversion method is illustrated by a color bar (for color explanation, see legend of the figure). Inversions use CO₂ and NO₂ cloud-filtered data and ERA5 winds.



1209

1210 **Figure A8:** Estimated vs true annual emissions for 4 inversion scenarios (titles of the panels). Results are displayed for the CSF,
 1211 GP, LCSF and ensemble methods that gather 2 or 3 of these individual methods. For the CSF method, annual estimates are
 1212 weighted means of the instant estimates while they are arithmetic means for the GP and LCSF methods. Each marker represents a
 1213 given emission source and each color a given inversion method. The divergence inversion method uses CO₂ data only for all the
 1214 inversion scenarios. The plain line represents the 1:1 line. The bottom-right legends display for each inversion method the relative
 1215 RMSE which is the RMSE between estimated and true annual emissions divided by the median of true annual emissions of all
 1216 sources (~9.6 MtCO₂.yr⁻¹).

1217

1218

1219



Method	Time frame	Potential for joint use of NO ₂ to detect plumes	Computational cost (1)
Integrated Mass Enhancement (IME)	Single-Image estimates	Yes	Medium: ~20 mn
Cross-Sectional Flux (CSF)	Single-Image estimates	Yes	Medium: ~25 mn
Gaussian Plume (GP)	Single-Image estimates	Yes	High: ~110 mn
Light Cross-Sectional Flux (LCSF)	Single-Image estimates	Yes	Low: ~10 mn
Divergence (Div)	Averaged estimates from ensemble of images	No	Medium: ~23 mn

1220 **Table 1: Summary of characteristics of the benchmarked methods. (1) Computation time was estimated by inverting one month of**
 1221 **CO₂ and NO₂ cloud-free SMARTCARB data on the same server using the ddeq package (Kuhlmann et al., 2023)**

1222

Benchmark Scenario	Wind dataset	Cloud fraction thresholds	Joint use of NO ₂ and CO ₂
Scenario 1	SMARTCARB	100% (no clouds)	Yes
Scenario 2	SMARTCARB	1% for CO ₂ , 30% for NO ₂	No
Scenario 3	SMARTCARB	100% (no clouds)	No
Scenario 4	SMARTCARB	1% for CO ₂ , 30% for NO ₂	Yes
Scenario 5	ERA5	100% (no clouds)	Yes
Scenario 6	ERA5	1% for CO ₂ , 30% for NO ₂	No
Scenario 7	ERA5	100% (no clouds)	No
Scenario 8	ERA5	1% for CO ₂ , 30% for NO ₂	Yes

1223 **Table 2: List of the different benchmarking scenarios: from the most optimistic (scenario 1) which considers inversions with cloud-**
 1224 **free data and SMARTCARB winds to the most realistic (Scenario 8) with cloud-filtered data and with ERA5 winds.**

1225

1226



1227

Inversion method	Cloud-free data	Cloud-filtered data
IME	1661	96
CSF	2028	302
GP	1776	266
LCSF	2722	313

1228 **Table 3. Number of estimates for each inversion method when data with or without clouds are used. Inversions are**
1229 **performed with CO₂ and NO₂ data and, with SMARTCARB winds.**

1230