

# Dealing with Spatial Heterogeneity in Pointwise to Gridded Data Comparisons

Amir H. Souri<sup>1\*</sup>, Kelly Chance<sup>1</sup>, Kang Sun<sup>2,3</sup>, Xiong Liu<sup>1</sup>, and Matthew S. Johnson<sup>4</sup>

<sup>1</sup>Atomic and Molecular Physics (AMP) Division, Harvard-Smithsonian Center for Astrophysics, Cambridge, MA, USA

<sup>2</sup>Department of Civil, Structural and Environmental Engineering, University at Buffalo, Buffalo, NY, USA

<sup>3</sup>Research and Education in Energy, Environment and Water Institute, University at Buffalo, Buffalo, NY, USA

<sup>4</sup>Earth Science Division, NASA Ames Research Center, Moffett Field, CA, USA

\*Corresponding author: [ahsouri@cfa.harvard.edu](mailto:ahsouri@cfa.harvard.edu)

## Abstract

Most studies on validation of satellite trace gas retrievals or atmospheric chemical transport models assume that pointwise measurements, which roughly represent the element of space, should compare well with satellite (model) pixels (grid box). This assumption implies that the field of interest must possess a high degree of spatial homogeneity within the pixels (grid box), which may not hold true for species with short atmospheric lifetimes or in the proximity of plumes. Results of this assumption often lead to a perception of a nonphysical discrepancy between data, resulting from different spatial scales, potentially making the comparisons prone to overinterpretation. Semivariogram is a mathematical expression of spatial variability in discrete data. Modeling the semivariogram behavior permits carrying out spatial optimal linear prediction of a random process field using kriging. Kriging can extract the spatial information (variance) pertaining to a specific scale, which in turn translating pointwise data to a gridded space with quantified uncertainty such that a grid-to-grid comparison can be made. Here, using both theoretical and real-world experiments, we demonstrate that this classical geostatistical approach can be well adapted to solving problems in evaluating model-predicted or satellite-derived atmospheric trace gases. This study suggests that satellite validation procedures using the present method must take kriging variance and satellite spatial response functions into account. We present the comparison of Ozone Monitoring Instrument (OMI) tropospheric NO<sub>2</sub> columns against 11 Pandora Spectrometer Instrument (PSI) systems during the DISCOVER-AQ campaign over Houston. The least-squares fit to the paired data shows a low slope ( $OMI=0.76 \times PSI + 1.18 \times 10^{15}$  molecules cm<sup>-2</sup>,  $r^2=0.67$ ) which is indicative of varying biases in OMI. This perceived slope, induced by the problem of spatial scale, disappears in the comparison of the convolved kriged PSI and OMI ( $0.96 \times PSI + 0.66 \times 10^{15}$  molecules cm<sup>-2</sup>,  $r^2=0.72$ ) illustrating that OMI possibly has a constant systematic bias over the area. To avoid gross errors in comparisons made between gridded data versus pointwise measurements, we argue that the concept of semivariogram (or spatial autocorrelation) should be taken into consideration, particularly if the field exhibits a strong degree of spatial heterogeneity at the scale of satellite and/or model footprints.

## 42 1. Introduction

43 Most of the literature on validation of satellite trace gas retrievals or atmospheric chemical  
44 transport models assume that geophysical quantities within a satellite pixel or a model grid box  
45 are spatially homogeneous. Nevertheless, it has long been recognized that this assumption can  
46 often be violated; spatially coarse atmospheric models or satellites are often not able to represent  
47 features, nor physical processes, transpiring at fine spatial scales. Janjic et al. [2016] used the term  
48 of *representation error* to describe this complication. They posit that this problem is a result of  
49 two combined factors: unresolved spatial scales and physiochemical processes. To elaborate on  
50 this definition, let us assume that an atmospheric model simulating CO<sub>2</sub> concentrations can  
51 represent the exact physiochemical processes but is fed with a constant CO<sub>2</sub> emission rate. This  
52 model obviously cannot resolve the spatial distribution of CO<sub>2</sub> concentration because we use an  
53 unresolved emission input. As another example, if we know the exact rates of CO<sub>2</sub> emissions but  
54 use a model unable to resolve atmospheric dynamics, the spatial distribution of CO<sub>2</sub> concentrations  
55 will be unrealistic due to unresolved physical processes.

56 Numerous scientific studies have reported on this matter. The simulations of short lifetime  
57 atmospheric compounds such as nitrogen dioxide (NO<sub>2</sub>), isoprene, formaldehyde (HCHO), and  
58 the hydroxyl radical (OH) have been found to be strongly sensitive to the model spatial resolution  
59 [Vinken et al., 2011; Valin et al., 2011; Yu et al., 2016; Pan et al., 2017]. Likewise, the performance  
60 of weather forecast models in resolving non-hydrostatic components heavily relies on both model  
61 resolution and parametrizations used. For example, when Kendon et al. [2014], Souri et al.  
62 [2020a], and Wang et al. [2017] defined a higher spatial resolution in conjunction with more  
63 elaborate model physics, they were able to more realistically simulate extreme or local weather  
64 phenomena such as convection and sea-land breeze circulation.

65 The spatial representation issue is not only limited to models. Satellite trace gas retrievals  
66 optimize the concentration of trace gases and/or atmospheric states to best match the observed  
67 radiance using an optimizer along with an atmospheric radiative transfer model. This procedure  
68 requires various inputs such as surface albedo, cloud and aerosol optical properties, and trace gas  
69 profiles, all of which come with different scales and representation errors. Moreover, the radiative  
70 transfer model by itself has different layers of complexity with regards to physics. A myriad of  
71 studies have reported that satellite-derived retrievals underrepresent spatial variability whenever  
72 the prognostic inputs used in the retrieval are spatially unresolved [e.g., Russell et al., 2011;  
73 Laughner et al., 2018; Souri et al., 2016; Goldberg et al., 2019; Zhao et al., 2020]. Additionally,  
74 the large footprint of some sensors relative to the scale of spatial variability of species inevitably  
75 leads to some degree of the representativity issues [e.g., Souri et al., 2020b, Tang et al., 2021; Judd  
76 et al., 2020]. It is because of this reason that several validation studies resorted to downscaling  
77 their relatively coarse satellite observations using high-resolution chemical transport models so  
78 that they could compare them to spatially finer datasets such as in-situ measurements [Kim et al.,  
79 2018; Choi et al., 2020]. Nonetheless, their results largely arise from modeling experiments which  
80 might be biased.

81 The validation of satellites or atmospheric models is widely done against pointwise  
82 measurements. Mathematically, a point is an element of space. Hence, it is not meaningful to  
83 associate a point with a spatial scale. If one compares a grid box to a point sample (i.e., apples to  
84 oranges), they are assuming that the point is the representative of the grid box. At this point, the  
85 fundamental question is: can the average of the spatial distribution of the underlying compound be  
86 represented by a single value measured at a subgrid location? This question was answered in  
87 Matheron [1963]. He advocated the notion of the semivariogram, a mathematical description of

88 the spatial variability, which finally led to the invention of kriging, the best unbiased linear  
 89 estimator of a random field. A kriging model can estimate a geophysical quantity in a common  
 90 grid. This is not exclusively special; a simple interpolation method such as the nearest neighbor  
 91 has the same purpose. The power of kriging lies in the fact that it takes the data-driven spatial  
 92 variability information into account and informs an error associated with the interpolated map.  
 93 This strength not only makes kriging a relatively superior model over simplified interpolation  
 94 methods, but also reflects the level of confidence pertaining to spatial heterogeneity dictated by  
 95 both data and the semivariogram model used through its variance [Chilès and Delfiner, 2009].

96 Different studies leveraged this classical geostatistical method to map the concentrations  
 97 of different atmospheric compounds at very high spatial resolutions [Tadić et al., 2017; Li et al.,  
 98 2019; Zhan et al., 2018; Wu et al., 2018]; To the best of our knowledge, Swall and Foley [2009]  
 99 is the only study that used kriging for a chemical transport model validation with respect to surface  
 100 ozone. They suggested that kriging estimation should be executed in grids rather than discrete  
 101 points. Kriging uses a semivariogram model in a continuous form. Optimizing the kriging grid size  
 102 (i.e., domain discretization) at which the estimation is performed is an essence to fully obtaining  
 103 the maximum spatial information from data. Another important caveat with Swall and Foley  
 104 [2009] is that averaging discrete estimates (points) to build grids is not applicable for remote  
 105 sensing data. Depending on the optics and the geometry, the spatial response function can  
 106 transform from an ideal box (simple average) to a sophisticated shape such as a super Gaussian  
 107 function (weighted average) [Sun et al., 2018]. Moreover, the footprint of satellites is not spatially  
 108 constant. We will address these complications in this study using both theoretical and real-world  
 109 experiments.

110 Our paper is organized with the following sections. Sections 2 is a thorough review of the  
 111 concept of the semivariogram and kriging. We then provide different theoretical cases, their  
 112 uncertainty, sensitivities with respect to difference tessellation, grid size, and the number of  
 113 samples. Section 3 proposes a framework for satellite (model) validation using sparse point  
 114 measurements and elaborates on the representation error using idealized experiments. Sections 4  
 115 introduces several real-world experiments.

## 116 **2. Semivariogram and Ordinary Kriging Estimator**

### 117 **2.1. Definition**

118 The semivariogram is a mathematical representation of the degree of spatial variability (or  
 119 similarity) in a function describing a regionalized geophysical quantity ( $f$ ), which is defined as  
 120 [Matheron, 1963]:

$$121 \gamma(h) = \frac{1}{2V} \iiint_V [f(x + \mathbf{h}) - f(x)]^2 dV \quad (1)$$

122 where  $x$  is a location in the geometric fields of  $V$ ,  $f(x)$  is the value of a quantity at the location of  $x$ ,  
 123 and  $\mathbf{h}$  is the vector of distance. If discrete samples are available rather than the continuous field,  
 the general formula can be simplified to the experimental semivariogram defined as:

$$124 \gamma(\mathbf{h}) = \frac{1}{2N(\mathbf{h})} \sum_{|x_i - x_j| - |\mathbf{h}| \leq \varepsilon} [Z(x_i) - Z(x_j)]^2 \quad (2)$$

125 where  $Z(x_i)$  (and  $Z(x_j)$ ) is discrete observations (or samples),  $N(\mathbf{h})$  is the number of paired  
 126 observations separated by the vector of  $\mathbf{h}$ .  $|\cdot|$  operator indicates the length of a vector. The condition  
 127 of  $|x_i - x_j| - |\mathbf{h}| \leq \varepsilon$  is to allow certain tolerance for differences in the length of the vector. For  
 128 simplicity, we only focus on an isotropic case meaning we rule out the directional (or angular)  
 dependency in  $\gamma(\mathbf{h})$ . Under this condition, the vector of  $\mathbf{h}$  becomes scalar ( $h = |\mathbf{h}|$ ).

129 If a reasonable number of samples is present, one can describe  $\gamma(h)$  through a regression  
 130 model (e.g., Gaussian or spherical shapes). The degree of freedom for this regression is:

$$dof = N - p \quad (3)$$

131 where  $p$  is the number of parameters defined in the model. For instance, to fit a Gaussian function  
 132 to the semivariogram with three parameters ( $p=3$ ), three paired ( $N=3$ ) observations are required at  
 133 minimum. Different regression models can be used to describe  $\gamma(h)$  depending on the  
 134 characteristic of the quantity of interest. In this study, we will use a stable Gaussian function:

$$\gamma(h) = a(1 - e^{-\left(\frac{h}{b}\right)^{c_0}}); c_0=1.5 \quad (4)$$

135 where  $a$  and  $b$  are fitting parameters. A non-linear least-squares algorithm based on Levenberg-  
 136 Marquardt method will be used to estimate the fitting parameters.

137 The kriging estimator predicts a value of interest over a defined domain using a  
 138 semivariogram model derived from samples [Chilès and Delfiner, 2009]. The kriging model is  
 139 defined as [Matheron, 1963]:

$$Z(x) = Y(x) + m(x) \quad (5)$$

140 where  $Y(x)$  is a zero-mean random function, and  $m(x)$  is a systematic drift. If we assume  
 141  $m(x) = a_0$ , the model is called ordinary kriging. Similar to an interpolation problem, the  
 142 estimation point ( $\hat{Z}$ ), is determined by linearly combining  $n$  number of samples with their weights  
 143 ( $\lambda_j$ ):

$$\hat{Z} = \sum_{j=1}^n \lambda_j Z(x_j) + \lambda_0 \quad (6)$$

144 where  $\hat{Z}$  is the estimation,  $\lambda_0$  is a constant weight,  $x_j$  is the location of samples, , and  $Z(x_j)$  is point  
 145 data (i.e., samples). The mean squared error of this estimation can be written as

$$E(\hat{Z} - Z_0)^2 = \text{Var}(\hat{Z} - Z_0) + \left[ \lambda_0 + \left( \sum_{j=1}^n \lambda_j - 1 \right) a_0 \right]^2 \quad (7)$$

146 Where  $Z_0$  is point observations ( $Z_0 = Z(x_j), j = 1, 2, \dots, n$ ), and  $a_0$  is the mean of  $Z$  which is  
 147 unknown. In order to estimate the weights, we are required to minimize Eq.7, but this cannot be  
 148 done without knowing the exact value of  $a_0$ . A solution is to assume  $\lambda_0 = 0$  and impose the  
 149 following condition:

$$\sum_{j=1}^n \lambda_j = 1 \quad (8)$$

150 This condition warrants  $E(\hat{Z} - Z_0)$  be zero and removes the need for the knowledge of  $a_0$ .  
 151 Therefore Eq.7 can be written as

$$E(\hat{Z} - Z_0)^2 = \text{Var}(\hat{Z} - Z_0) = \sum_{j_1=1}^n \sum_{j_2=1}^n \lambda_{j_1} \lambda_{j_2} \gamma_{j_1 j_2} - 2 \sum_{j_1=1}^n \lambda_{j_1} \gamma_{j_1 0} + \gamma_{00} \quad (9)$$

152 where  $\gamma_{j_1 j_2}$  is the spatial covariance between the point observations and  $\gamma_{j_1 0}$  is the spatial  
 153 covariance of between the observations and the estimation point. The spatial covariance is modeled  
 154 by a semivariogram. Using the method of Lagrange multiplier and considering the constraint on  
 155 the weights, Eq.9 can be minimized by solving the following problem [Chilès and Delfiner, 2009]:

$$\begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_n \\ \mu \end{pmatrix} = \begin{pmatrix} \gamma(x_1 - x_1) \cdots \gamma(x_1 - x_n) & 1 \\ \vdots & \vdots \\ \gamma(x_n - x_1) \cdots \gamma(x_n - x_n) & 1 \\ 1 & \cdots & 1 & 0 \end{pmatrix}^{-1} \begin{pmatrix} \gamma(x_1 - x_0) \\ \vdots \\ \gamma(x_n - x_0) \\ 1 \end{pmatrix} \quad (10)$$

156 where  $\mu$  is the Lagrange parameter and  $x_0$  is the location of estimation. The first term in the right  
 157 hand side of this equation shows the spatial variability described by the semivariogram model  
 158 among samples, whereas the second term indicates the modeled variability between samples and  
 159 the estimation point. The unknowns (the left hand side of the equation) have a unique solution if,  
 160 and only if, the semivariogram model is positive definite and the samples are unique [Chilès and  
 161 Delfiner, 2009]. The estimation error can be obtained by

$$\sigma^2 = E(\hat{Z} - Z_0)^2 = \sum_{j=1}^n \lambda_j \gamma_{j0} - \mu \quad (11)$$

162 This equation is an important component in the kriging estimator. Not only can we estimate  $Z(x)$   
 163 given a selection of data points, but also an uncertainty associated with such estimation can be  
 164 provided.

## 165 **2.2. Theoretical Cases**

### 166 *2.2.1. Sensitivity to spatial variability of the field*

167 The present section illustrates the application of ordinary kriging for several numerical  
 168 cases. Five idealized cases are simulated in a grid of  $100 \times 100$  pixels, namely, a constant field (C1),  
 169 a ramp starting from zero in the lower left to higher values in the upper right (C2), an intersection  
 170 with concentrated values in four corridors (C3), a Gaussian plume placed in the center (C4), and  
 171 multiple Gaussian plumes spread over the entire domain (C5). We randomly sample 200 data  
 172 points from each field as is, and successively create the semivariograms in 100 binned distances.  
 173 Except C1, which lacks a spatial variability thus  $\gamma(h) = 0$ , other semivariograms are fit with the  
 174 stable Gaussian function. Using the semivariogram model, we optimize Eq.10 to estimate  $\hat{Z}(x)$   
 175 for each pixel (i.e.,  $100 \times 100$ ) with the estimation errors based on Eq.11. Figure 1 depicts the truth  
 176 field ( $Z(x)$ ), semivariograms made from the samples, estimated values ( $\hat{Z}(x)$ ), difference of  $Z(x)$   
 177 and  $\hat{Z}(x)$ , and error associated with the estimation.

178 As for C1, the uniformity results in a constant semivariogram leading the estimation to be  
 179 identical to the truth. This estimation signifies the unbiased characteristic of ordinary kriging. C1  
 180 is never met in reality, however, it is possible to assume some degree of uniformity among data  
 181 restrained to background values; a typical example of this can be seen in the spatial distribution of  
 182 a number of trace gases in pristine environments such as  $\text{NO}_2$  [e.g., Wang et al., 2020] and  $\text{HCHO}$   
 183 [Wolfe et al., 2019]. Under this condition, any data point within the field (i.e., the satellite  
 184 footprint) can be assumed to be representative of the spatial variability in truth.

185 Concerning C2, the semivariogram shows a linear shape meaning data points at larger  
 186 distances exhibit larger differences. Generally geophysical samples are uncorrelated at large  
 187 distances, thereby one expects the semivariogram to increase more slowly as the distance gets  
 188 further. The steady increase in  $\gamma(h)$  is indicative of a systematic drift in the data invalidating the  
 189 assumption of  $m(x) = a_0$ . In many applications, a simple polynomial can explain  $m(x)$  and  
 190 subsequently be subtracted from the data points. An example of this problem is tackled by Onn  
 191 and Zebker [2006]; it concerns the spatial variability of water vapor columns measured by GPS  
 192 signals. Onn and Zebker [2006] observed a strong relationship between the water vapor columns  
 193 and GPS altitudes resulting from the vertical distribution of water vapor in the atmosphere.  
 194 Because of this complication, a physical drift model describing the vertical dependency was fit

195 and removed from the measurements so that they could focus on the horizontal fluctuations. In  
196 terms of C2, one can effortlessly reproduce  $Z(x)$  by fitting a three-dimensional plane to barely  
197 three samples, indicating that the semivariogram is of little use.

198 C3 is an example of an extremely inhomogeneous field manifested in the stabilized  
199 semivariogram at a value of  $\gamma \sim 500$ , called the sill [Chilès and Delfiner, 2009], indicating  
200 insignificant information (variance) from the samples beyond this distance ( $\sim 20$ ), called the range.  
201 Range is defined as the separation distance at which the total variance in data is extracted. The  
202 smaller the range is, the more heterogeneous the samples will be. While the estimated field roughly  
203 captures the shape of the intersections, it is spatially distorted at places with relatively sparse data  
204 points. The kriging model error is essentially a measure of the density of information. It converges  
205 to zero in the samples location and diverges to large values in gaps.

206 C4 is a close example of a point source emitter with faint winds and turbulence. The  
207 semivariogram exhibits a bell shape. As samples get further from the source, the variance diverges,  
208 stabilizes, and then sharply decreases. This is essentially because many data points with low  
209 values, apart from each other, have negligible differences. This tendency is recognized as the hole  
210 effect which is characterized for high values to be systemically surrounded by low values (and  
211 vice versa). It is possible to mask this effect by fitting a semivariogram model stabilizing at certain  
212 sill (like the one in Figure 1). Nonetheless, if the semivariogram shows periodic holes, the fitted  
213 model should be modified to a periodic cosine model [Pyrcz and Deutsch, 2003].

214 The last case, C5, shows a less severe case of the hole effect previously observed in C4.  
215 This is due to the presence of more structured patterns in different parts of the domain. The range  
216 is roughly twice as large as the previous case (C4) denoting that there is more information  
217 (variance) among the samples at larger distances. A number of experiments using this particular  
218 case will be discussed in the following subsections.

### 219 *2.2.2. Sensitivity to the number of samples*

220 It is often essential to optimize the number of samples used for kriging. The kriging  
221 estimator somewhat recognizes its own capability at capturing the spatial variability through  
222 Eq. 11. Thus, if the target is spatially too complex and/or the samples are too limited, the estimator  
223 essentially informs that  $\hat{Z}(x_0)$  is unreliable through large variance. However, there is a caveat;  
224  $Y(x)$  must be a Gaussian random model with a zero mean so that kriging can capture the statistical  
225 distribution of  $\hat{Z}$  given the data points. Except this case, the kriging variance can either be  
226 underestimated or overestimated depending on the level of skewness of the statistical distribution  
227 of  $Y(x)$  [Armstrong, 1994]. Figure 2 shows the kriging estimation for C5 using 5, 25, 50, 100, and  
228 500 random samples in the entire field. Immediately apparent is a better description of the  
229 semivariogram when larger number of samples are used, which in turn, results in a better  
230 estimation of  $Z(x)$ . The optimum number of samples to reproduce  $Z(x)$  depends on the  
231 requirement for the relative error ( $\sigma/Z(x)$ ) being met at a given location.

### 232 *2.2.3. Sensitivity to the tessellation of samples*

233 A common application of kriging is to optimize the tessellation of data points for a fixed  
234 number of samples to achieve a desired precision. In real-world practices, the objective of such  
235 optimization is very purpose-specific, for example, one might prefer a spatial model representing  
236 a certain plume in the entire domain. Different ways for data selection exist [e.g., Rennen, 2008],  
237 but for simplicity, we focus on four categories: purely random, stratified random, a uniform grid,  
238 and an optimized tessellation. Figure 3 demonstrates the estimation of C5 using 25 samples chosen  
239 based on those four procedures.

240 Concerning the random selection, the lack of samples over two minor plumes cause the  
241 estimation to deviate largely from the truth. While a random selection may seem to be practical  
242 because it is independent of the underlying spatial variability, it can suffer from under sampling  
243 issues, thus being inefficient. As a remedy, it might be advantageous to group the domain into  
244 similar zones and randomly sample from each, which is commonly known as stratified random  
245 selection. We classify the domain into four zones by running the k-mean algorithm on the  
246 magnitudes of  $Z(x)$  (not shown) and randomly sample six to seven points from each one (total 25).  
247 We achieve a better agreement between the estimated field and the truth because we exploited  
248 some prior knowledge (here the contrast between low and high values).

249 As for the uniform grid, we notice that there are fewer data points in the semivariogram  
250 stemming from redundant distances which is indicative of correlated information. Nonetheless, if  
251 the desired tessellation is neutral with regard to location meaning that all parts of the domain is  
252 equal of scientific interest, the uniform grid is the most optimal design for the prediction of  $Z(x)$   
253 under an ideally isotropic case. A mathematical proof for this claim can be found in Chilès and  
254 Delfiner [2009].

255 To execute the last experiment, we select 25 random samples for 1000 times and find the  
256 optimal estimation by finding the minimum sum of  $|\hat{Z}(x_0) - Z(x)|$ . It is worth mentioning that  
257 the optimized tessellation is essentially a local minimum based on 1000 kriging attempts. The  
258 optimized location of samples seems to more clustered over areas with large spatial gradients. Not  
259 too surprisingly, we observe the smallest discrepancy between the estimation and the truth.

260 A lingering concern over the application of these numerical experiments is that the truth is  
261 assumed to be known. The truth is never known, by this means we may never exactly know how  
262 well or poorly the kriging estimator is performing. However, it is highly unlikely for some prior  
263 understandings or expectations of the truth to be absent. If this is the case, which is rare, a uniform  
264 grid should be intuitively preferred to deliver the local estimations of average values in uniform  
265 blocks. In contrast, if the prior knowledge is articulated by previous site visits, model predictions,  
266 theoretical experiments, pseudo-observations, or other relevant data, the tessellation needs to be  
267 optimized.

268 It is important to recognize that the uncertainties associated with the prior knowledge  
269 directly affects the level of confidence in the final answer. Accordingly, the prior knowledge error  
270 should ultimately be propagated to the kriging variance. The determination of the prior error is  
271 often done pragmatically. For example, if the goal is to design the location of thermometer sites to  
272 capture surface temperature during heat waves using a yearly averaged map of surface  
273 temperature, it would be wise to specify a large error with this specific prior information to play  
274 down the proposed design. This is primarily because the averaged map underrepresents such an  
275 atypical case. A possible extension of this example would be to use a weather forecast model with  
276 quantified errors capable of capturing retrospective heat waves. Although a reasonable forecast in  
277 the past does not necessarily guarantee a reasonable one in the future, it is rational to assume for  
278 the uncertainty with a new tessellation design using the weather model forecast to be lower than  
279 that of using the averaged map.

280 A general roadmap for the data tessellation design is shown in Figure 4. As proven in Chilès  
281 and Delfiner [2009], if the field is purely isotropic, the uniform grid is the most intuitive sensible  
282 choice when the prior information on the spatial variability is lacking. When the prior knowledge  
283 with quantified errors is available, an optimum tessellation can be achieved by running a large  
284 number of kriging models with suitable  $\gamma(h)$  and picking the one yielding the minimum difference  
285 between the prior knowledge and the estimation. The choice of the cost function (here L1 norm)

286 is purpose-specific. For example, if the reconstruction of a major plume was the goal, using a  
287 weighted cost function, geared towards capturing the shape of plume, would be more appropriate.

#### 288 *2.2.4. Sensitivity to the grid size*

289 A kriging model can estimate a geophysical quantity at a desired location considering the  
290 data-driven spatial variability information. Since the kriging model is practically in a continuous  
291 form, the desired locations can be anywhere within the field of  $V$ . A question is whether or not it  
292 is necessary to map the data onto a very fine grid. There is a trade-off between the computational  
293 cost and the accuracy of the interpolated map. The range of the underlying semivariogram helps  
294 in finding the optimal solution. The greater the range (i.e., a more homogeneous field), the less  
295 important to map the data in a finer grid.

296 Figure 5a depicts an experiment comparing the estimates of C2 at different grid sizes with  
297 the truth. The departure of the estimate from the truth is rather negligible for several coarse grids  
298 (e.g.,  $10 \times 10$ ). The homogeneous field, manifested by the large range (Figure 1), allows for a  
299 reasonable estimation of  $Z(x)$  at coarse resolutions with inexpensive computational costs. Figure  
300 5b shows the same experiment but on C5 with the optimized tessellation. As opposed to the  
301 previous experiment, the estimate substantially diverges from the truth when increasing the grid  
302 size, suggesting that a finer resolution should be used for fields with smaller ranges (i.e.,  
303 heterogeneous fields).

304 The complexity of directly using the range for choosing the optimal grid size arises from  
305 the fact that the level of spatial homogeneity can vary within the domain. In fact, the range is  
306 derived from a semivariogram model representing a crude estimate of varying ranges occurring at  
307 various scales. It is intuitively clear that depending on the degree of heterogeneity, which is  
308 spatiotemporally variable, the grid size needs to be adaptively adjusted [Bryan, 1999]. For the sake  
309 of simplicity, but at a higher computational cost, we adopt a numerical solution which is to first  
310 simulate on a coarse grid, then on a finer one until the difference with respect to the previous grid  
311 size across all pixels reaches to an acceptable value ( $< 1\%$ ). We name this output ( $1 \times 1$ ) with the  
312 optimized tessellation for C5 as C5opt.

### 313 **3. Comparison of points to satellite pixels**

#### 314 *3.1. Synching the scales between the gridded field and satellite pixels*

315 To minimize the complications of different spatial scales between two gridded data, we  
316 first need to upscale the finer resolution data to match the coarse ones. In case of numerical  
317 chemical transport or weather forecast models, the size of the grid box is definitive. Likewise, a  
318 satellite footprint, mainly dictated by the sensor design, the geometry, and signal-to-noise  
319 requirements [Platt et al., 2021], is known. However, the grid size of the kriging estimation is a  
320 variable subject to optimization which has been discussed previously.

321 When we compare the grid size of the kriging estimate to that of a satellite (or a model),  
322 three situations arise: First, the kriging spatial resolution is coarser than the satellite, a condition  
323 occurring when either the field is homogeneous or the field is under sampled. In situations where  
324 the field is homogeneous ( $\gamma(h) \cong 0$ ), it is safe to directly compare the data points to the satellite  
325 measurements without having to use kriging. If the under sampling is the case (see Figure 2 with  
326 5 samples), it is sensible to first investigate if the field is homogeneous within the satellite footprint  
327 using different data (if any). If the homogeneity is met, we either can compare two datasets without  
328 kriging or to match the size of kriging grid cell with the satellite footprint and statistically involve  
329 the kriging variance in the comparison (discussed later); nonetheless, the kriging estimate beyond  
330 the location of samples must be used with extra caution because their variance very quickly  
331 departures from zero to extremely large numbers (see Figure 1). Thus, there is a compromise



332 between increasing the number of paired samples between two datasets and enhancing the level of  
 333 confidence in statistics. If independent observations suggest that there might be large heterogeneity  
 334 within a satellite footprint, it is strongly advised against quantitatively comparing the points to the  
 335 satellite observations. Second, the number of samples is fewer than three observations in the field  
 336 so it is in principal impossible to build a semivariogram. Validating a satellite under this condition  
 337 is prone to misinterpretation because the spatial heterogeneity cannot be modeled. Nonetheless, if  
 338 one presumes a good degree of homogeneity within the sensor footprint (such as very high-  
 339 resolution remote sensing airborne data), the direct comparison of point measurements might be  
 340 possible. Third, the satellite footprint is coarser than the kriging estimate. Under this condition, we  
 341 upscale the kriging map to match the spatial resolution of the satellite using

$$\hat{Z}_c = \hat{Z}_f * S = \int \hat{Z}_f(x)S(x - y)dy \quad (12)$$

342 where  $S$  is the spatial response function,  $\hat{Z}_c$  is the coarse kriging field,  $\langle * \rangle$  is the convolution  
 343 operator,  $y$  is shift, and  $\hat{Z}_f$  is the fine field. In discrete form we can rewrite Eq.12 in

$$\hat{Z}_c[i, j] = \sum_m \sum_n \hat{Z}_f[i - m, j - n] S[m, n] \quad (13)$$

344 where  $m$  and  $n$  are the dimension of the response function. The mathematical formulation of  
 345  $S[m, n]$  for a number of satellites can be represented by two-dimensional super Gaussian functions  
 346 as discussed in Sun et al. [2018]. Atmospheric models have a uniform response to the simulated  
 347 values within a grid box, therefore  $S[m, n] = \frac{1}{m \times n} J_{m,n}$ , where  $J$  is the matrix of ones. In the same  
 348 way, the kriging variance should be convolved through

$$\sigma_c^2[i, j] = \sum_m \sum_n \sigma_f^2[i - m, j - n] S^2[m, n] \quad (14)$$

349 where a superscript of 2 denotes squaring, and  $\sigma_c^2$  and  $\sigma_f^2$  are the kriging variance in the coarse  
 350 and the fine grids, respectively.

351 To demonstrate the upscaling procedure, we use C5opt ( $1 \times 1$ ) and upscale it at six grid sizes  
 352 ( $m, m$ ) of  $5 \times 5$ ,  $10 \times 10$ ,  $15 \times 15$ ,  $20 \times 20$ ,  $25 \times 25$ , and  $30 \times 30$ . For simplicity, we consider  $S = \frac{1}{m^2} J_{m,m}$ ;  
 353 this spatial response function results in averaging the values in the grid boxes. Figure 6 shows the  
 354 resultant map overplotted with the samples along with the error estimation. Two tendencies from  
 355 this experiment can be identified: First, the discrepancy of the point data and  $\hat{Z}$  is becoming more  
 356 noticeable as the grid size grows; this directly speaks to the notion of the spatial representativeness;  
 357 large grid boxes are less representative of sub-grid values. Second, the gradients of the field along  
 358 with the estimation error become smoother primarily due to convolving the field with the spatial  
 359 response function, which acts as a low pass filter.

360 We further directly compare  $\hat{Z}$  to the samples (i.e., observations) shown in Figure 7. We  
 361 see an excellent comparison between  $\hat{Z}$  at  $1 \times 1$  resolution with the observations underscoring the  
 362 unbiasedness characteristic of the kriging estimator. Conversely, the upscaled field gradually  
 363 diverges from the observations. This divergence is *the problem of scale*.

### 364 **3.2. Point to pixel vs pixel to pixel**

365 To elaborate on the problem of scale, we design an idealized experiment theoretically  
 366 validating pseudo satellite observations against some pseudo point measurements. The pseudo  
 367 satellite observations are created by upscaling the C5 truth ( $Z$ ) to  $30 \times 30$  grid footprint considering

368  $S = \frac{1}{m^2} J_{m,m}$ , meaning that the satellite is observing the truth but in a different scale (Figure S1).  
 369 The pseudo point measurements are the ones used for C5opt. Figure 8a shows the direct  
 370 comparison of the satellite pixel with the point observations. By ignoring the fundamental fact that  
 371 these two datasets are inherently different in nature, displaying the same geophysical quantity by  
 372 at different scales, we observe a perceived discrepancy ( $r^2=0.64$ ). The comparison suggests a  
 373 wrong conclusion that the satellite observations are biased-low. This discrepancy is unrelated to  
 374 any observational or physical errors, rendering any physical interpretation of the comparison  
 375 biased due to spatial-scale differences in the data sets. Figure 8b depicts the comparison of each  
 376 grid box of the upscaled kriging estimate ( $30 \times 30$ ) with that of the satellite. This direct comparison  
 377 shows a strong degree of agreement ( $r^2=0.98$ ), shaking off the erroneous idea of directly comparing  
 378 point to gridded data when the field exhibits substantial spatial heterogeneity.

379 Yet, the comparison misses an important point: the kriging estimate is considered error-  
 380 free. We attempt to incorporate the kriging variance through a Monte Carlo linear regression  
 381 method. Here, the goal is to find an optimal linear fit ( $y = ax + b + \varepsilon$ ) such that  $\chi^2 =$   
 382  $\sum \frac{[y-f(x_i, a, b)]^2}{\sigma_y^2 + a^2 \sigma_x^2}$  is minimized.  $\sigma_y^2$  and  $\sigma_x^2$  are the variances of  $y$  (here the satellite) and  $x$  (the kriging  
 383 variance), respectively. We set the errors of  $y$  to zero, and randomly perturb the errors of  $x$  based  
 384 on a normal distribution with zero mean and a standard deviation equal to that of kriging estimate  
 385 15,000 times. The average of optimized  $a$  and  $b$  coefficients derived from each fit are then  
 386 estimated and their deviation at 95% confidence interval assuming a Gaussian distribution is  
 387 determined. Figure 8b,c show the linear fit with and without considering the kriging error estimate.  
 388 The linear fit without involving the kriging error gives a strong impression that it is nearly perfect,  
 389 following closely to the paired observations. This is essentially explainable by the primary goal of  
 390  $\chi^2$  which is to minimize the L2 norm of residuals ( $y - f(x_i, a, b)$ ), portraying a very optimistic  
 391 picture of the satellite validation. The linear fit considering the kriging errors is different. The  
 392 uncertainties associated with  $a$  and  $b$  are larger since  $x$  is variable (shown in horizontal error bars).  
 393 The optimal fit gravitates towards the points with smaller standard deviations as they impose a  
 394 larger weight. The confidence in the linear fit at higher values is lower due to their errors being  
 395 large. This fit is a more realistic portrayal of the satellite validation.

396 Figure 9 summarizes the general roadmap for satellite (and model) validations against point  
 397 measurements. To fit the semivariogram with at least two parameters, we are required to have  
 398 three samples at minimum. Therefore, it is implausible to derive the spatial information from the  
 399 point data where sampling is extremely sparse (<3 samples within the field). The only case of  
 400 directly comparing point and satellite pixels is when the field within satellite footprint or the field  
 401 in general is rather homogeneous confirmed by independent data/models. Having more samples  
 402 allows to acquire some information on the spatial heterogeneity. The information carried by the  
 403 data is considered more and more robust with increasing the number of samples. Subsequently,  
 404 the kriging map along with its variance derived from a reasonable semivariogram at an optimized  
 405 grid resolution should be convolved with the satellite response function so that we can conduct an  
 406 apples-to-apples comparison. A real-world example on the satellite validation will be shown later.

## 407 **4. Real-world experiments**

### 408 **4.1. Spatial distribution of NO<sub>2</sub>**

409 We begin with focusing on tropospheric NO<sub>2</sub> columns observed by TROPOMI sensor  
 410 [Copernicus Sentinel data processed by ESA and Koninklijk Nederlands Meteorologisch Instituut  
 411 (KNMI), 2019; Boersma et al., 2018] at ~13:30 LST. We choose NO<sub>2</sub> primarily due to its spatial  
 412 heterogeneity [e.g., Souri et al., 2018; Nowlan et al., 2016, 2018; Valin et al., 2011; Judd et al.,

2020]. We oversample good quality pixels ( $qa\_flag > 0.75$ ) through a physical-based gridding approach [Sun et al., 2018] over Texas at  $3 \times 3 \text{ km}^2$  resolution in four seasons in 2019. We extract samples by uniformly selecting the  $\text{NO}_2$  columns in the center of each  $30 \times 30 \text{ km}^2$  block. The semivariogram along with its model are calculated, and then we kriging the samples. Figure 10 shows the  $\text{NO}_2$  columns map for four different seasons, the semivariogram, the kriging estimates, and the differences between the estimate and the field. High levels of  $\text{NO}_2$  are confined to cities indicating the sources being predominantly anthropogenic. Wintertime  $\text{NO}_2$  columns are larger than summertime mainly due to meteorological conditions and the OH cycle, the major sink of  $\text{NO}_2$ . All semivariograms exhibit the hole effect. This is because of high values of  $\text{NO}_2$  being systematically surrounded by low values. Regardless of the season, we fit the stable Gaussian to variances at distances smaller than  $2.5^\circ$  ( $\sim 275 \text{ km}^2$ ). The  $b_0$  parameter explaining the length scale is found to be 0.94, 0.88, 0.71, and 0.83 degree for DJF, MAM, JJA, and SON, respectively. These numbers strongly coincide with the seasonal lifetime of  $\text{NO}_2$  [Shah et al., 2020]; wintertime  $\text{NO}_2$  columns are spatially more uniform around the sources thus in relative sense, they are more homogeneous (spatially correlated) than those in warmer seasons. On the other hand, the shorter  $\text{NO}_x$  lifetime in summer results in a steeper gradient of  $\text{NO}_2$  concentrations. This tendency should not be generalized because transport and various  $\text{NO}_x$  sources including biomass burning, soil emissions, and lightning and can have large spatiotemporal variability resulting in different length scales in different times of a year. The differences between the kriging estimate and the field show some spatial structures indicating that  $\text{NO}_2$  is greatly heterogenous.

#### 4.2. Optimized tessellation over Houston

The preceding TROPOMI data enabled us to optimize a tessellation of ground-based point spectrometers over Houston. Our goal here is to propose an optimized network for winter 2021 given our knowledge on the spatial distribution of  $\text{NO}_2$  columns in winter 2019 measured by TROPOMI. The assumption of using a retrospective  $\text{NO}_2$  field for informing a hypothetical future campaign is not entirely unrealistic. If we have a consistent number of pixels from TROPOMI between two years, it is unlikely for the spatial variance of  $\text{NO}_2$  to be substantially different for the same season. We follow the framework proposed in Sect. 2.2.3 involving randomly selecting samples from the field (for 50000 iteration), and calculating kriging estimates for a given number of spectrometers. We then chose the optimum tessellation based on the minimum sum of  $|\hat{Z}(x_0) - Z(x)|$ .

Figure 11 shows the optimized tessellation given 5, 10, 15, and 20 spectrometers over Houston. The Houston plume is better represented with more samples being used. All cases share the same feature; the optimized samples are clustered in the proximity or within the plume. This tendency is clearly intuitive. We are required to place the spectrometers in locations where a substantial gradient (variance) in the field is expected. The difference between kriging estimate and the TROPOMI observations using 20 samples does not substantially differ in comparison to the one using 15 samples. Therefore, to keep the cost low, a preferable strategy is to keep the number of spectrometers as low as possible while achieving a reasonable accuracy. Based on the presented results, the optimized tessellation using 15 samples is preferred among others because it achieves roughly the same accuracy as the one with 20 samples.

#### 4.3. Validating OMI tropospheric $\text{NO}_2$ columns during DISCOVER-AQ 2013 campaign using Pandora

In order to understand ozone pollution [e.g., Mazzuca et al., 2016; Pan et al., 2017; Pan et al., 2015], characterize anthropogenic emissions [Souri et al., 2016, 2018], and validate satellite data [Choi et al., 2020], an intensive air quality campaign was made in September 2013 over

459 Houston (DISCOVER-AQ). The campaign encompassed a large suite of Pandora spectrometer  
460 instrument (PSI) (11 stations) measuring total NO<sub>2</sub> columns with a high precision ( $2.7 \times 10^{14}$   
461 molecules cm<sup>-2</sup>) and a moderate nominal accuracy ( $2.7 \times 10^{15}$  molecules cm<sup>-2</sup>) under the clear-sky  
462 condition [Herman et al., 2009]. We remove the observations with an error of >0.05 DU,  
463 contaminated by clouds, and averaged them over the month of September at 13:30 LST ( $\pm 30$   
464 mins). We attempt to validate OMI tropospheric NO<sub>2</sub> columns version 3.0 [Bucsela et al., 2013]  
465 refined in Souri et al. [2016] with the 4-km model profiles. The OMI sensor resolution varies from  
466  $13 \times 34$  km<sup>2</sup> at nadir to  $\sim 40 \times 160$  km<sup>2</sup> at the edge of the scan line. Biased pixels were removed based  
467 on cloud fraction > 0.2, terrain reflectivity > 0.3, and main (xtrack) quality flags =0. Following  
468 Sun et al. [2018], we oversample high quality pixels in the month of September 2013 over Houston  
469 at  $0.2 \times 0.2^\circ$  resolution. To remove the stratospheric contributions from PSI measurements, we  
470 subtract OMI stratospheric NO<sub>2</sub> ( $2.8 \pm 0.16 \times 10^{15}$  molecules cm<sup>-2</sup>) from the total columns over the  
471 area. Figure 12 shows the monthly-averaged tropospheric NO<sub>2</sub> columns measured by OMI  
472 overplotted by 11 PSIs. The elevated NO<sub>2</sub> levels (up to  $\sim 6 \times 10^{15}$  molecules cm<sup>-2</sup>) are seen over the  
473 center of Houston.

474 We then follow the validation framework shown in Figure 9 in which the number of point  
475 measurements and the level of heterogeneity are the main factors in deciding if we should directly  
476 compare them to the satellite pixels. Figure 13 shows the monthly-averaged PSI measurements  
477 along with the semivariogram and resulting kriging estimate at an optimized resolution ( $\sim 2$  km<sup>2</sup> =  
478 13800 data over the entire region) and errors. The distribution of semivariogram suggests that there  
479 is a strong degree of spatial heterogeneity, necessitating the use of kriging. We fit a stable Gaussian  
480 to the semivariogram resulting in  $2.23 \times (1 - e^{-\frac{h}{0.19}})^{1.5}$ . The spatial information (variance) levels  
481 off at  $0.19^\circ$  ( $\sim 21$  km) with a maximum variance equal to  $2.23$  molecules<sup>2</sup> cm<sup>-4</sup>. The measurements  
482 beyond this range (21 km) have a minimal weight due to this length scale. It is because of this  
483 reason that we see the kriging estimate converges to a fixed value at places being further than this  
484 range. The kriging errors of those grid boxes are constantly large (40% relative error). The  
485 optimum grid size for kriging is found to be 2 km<sup>2</sup> (<1% difference across all grid boxes).  
486 Subsequently, we use the super Gaussian spatial response function described in Sun et al. [2018]  
487 to convolve both the kriging estimate and error within (see Figure S2). Figure 14 shows the  
488 differences between the kriging estimate and error before and after convolution. The response  
489 function (OMI pixel) tends to be on average coarser than 2 km<sup>2</sup> resulting in smoothing of both the  
490 kriging estimate and error.

491 We ultimately conduct two different sets of comparison: directly comparing PSI to OMI  
492 pixels, and comparing convolved kriged PSI to OMI. It is worth noting that PSI measurements are  
493 monthly-averaged; similarly OMI data are oversampled in a monthly basis. In terms of the PSI,  
494 we only account for grid boxes whose kriging error is below  $1.2 \times 10^{15}$  molecules cm<sup>-2</sup> (1193  
495 samples, 8% of total kriging grid boxes). As for the grid-to-grid comparison, the kriging variance  
496 is considered in the linear polynomial fitted to the data through the Monte Carlo of chi-square with  
497 5,000 iterations. The variability with the OMI stratospheric NO<sub>2</sub> columns ( $0.16 \times 10^{15}$  molecules  
498 cm<sup>-2</sup>) is added to the PSI error for both analyses. The left and right panels of Figure 15 show the  
499 comparisons. As for the direct comparison of actual points (PSI) to pixels (OMI), the PSI  
500 measurements indicate a deviation of the slope ( $r^2=0.66$ ) from the unity line. This suggests that  
501 there is an unresolved magnitude-dependent systematic error. The grid-to-grid comparison not  
502 only offers a clearer picture of the distribution of data points, but also it hints at the offset being  
503 rather constant ( $0.66 \pm 0.18 \times 10^{15}$  molecules cm<sup>-2</sup>;  $r^2=0.72$ ). We also observe that the statistics  
504 between the satellite and the benchmark are moderately improved. This comparison in general

505 provides an important implication: the varying offsets in a plume shape environment (high to low  
506 values) are not necessarily due to variable offsets in the satellite retrieval, as the kriging estimate  
507 suggests that those varying offsets in point-to-pixel comparison, manifested in slope = 0.76, are a  
508 result of varying spatial scales.

### 509 **Summary**

510 There needs to be increased attention to the spatial representativity in the validation of  
511 satellite (model) against pointwise measurements. A point is the element of space, whereas satellite  
512 (model) pixels (grid box) are (at best) the product of the integration of infinitesimal points and a  
513 normalized spatial response function. If the spatial response function is assumed to be an ideal  
514 box, the resulting grid box will represent the average. Essentially, no justifiable theory exists to  
515 accept that the averaged value of a population should absolutely match with a sample, unless all  
516 samples are identical (i.e., a spatially homogeneous field). This glaring fact is often overlooked in  
517 the atmospheric science community. At a conceptual level, we are required to translate pointwise  
518 data to the grid format (i.e., rasterization). This can be done by modeling the spatial autocorrelation  
519 (or semivariogram) extracted from the spatial variance (information) among measured sample  
520 points. Assuming that the underlying field is a random function with an unknown mean, the best  
521 linear unbiased predictions of the field can be achieved by kriging using the modeled  
522 semivariograms.

523 In this study, we discussed methods for the kriging estimation of several idealized cases.  
524 Several key tendencies were observed through this experiment: first, the range corresponded to the  
525 degree of spatial heterogeneity; a larger range indicated the less presence of heterogeneity. Second,  
526 the kriging variance explaining the density of information quickly diverged from zero to large  
527 values when the field exhibited large spatial heterogeneity. This tendency mandates increasing the  
528 number of samples (observations) for those cases. Third, while the semivariogram models were  
529 constructed from discrete pair of samples, they are mathematically in a continuous form. It is  
530 because of this reason that we determined the optimal spatial resolution of the kriging estimate by  
531 incrementally making the grids finer and finer until a desired precision (=1%) was met.

532 The present study applied kriging to achieve an optimum tessellation given a certain  
533 number of samples such that the difference between our prior knowledge of the field, articulated  
534 by previous observations, models or theory, and the estimation is minimal. Usually there is  
535 uncertainty about the prior knowledge that should be propagated to the final estimates. The  
536 optimum tessellation for a range of idealized and real-world data consistently voted for placing  
537 more samples in areas where the gradients in the measurements were significant such as those  
538 close to point emitters.

539 This study also revisited the spatial representativity issue; it limits the realistic  
540 determination of biases associated with satellites (models). In one experiment, we convolved the  
541 kriging estimate for a multi-plume field with a box filter but various sizes. The perfect agreement  
542 ( $r=1.0$ ) between the samples (point) and kriging output (pixel) seen at a high spatial resolution  
543 gradually vanished with coarsening of the resolution of grid boxes ( $r=0.8$ ). We also directly  
544 compared samples (point) with pseudo satellite observations (showing the truth) with a coarse  
545 spatial resolution which led to a flawed conclusion about the satellite being biased-low. We  
546 modeled the semivariogram of those samples, estimated the field using kriging, and convolved  
547 with the pseudo-satellite spatial response function. The direct comparison of this output with that  
548 of the satellite showed a completely different story suggesting that the data were rather free of any  
549 bias. A serious caveat with using a spatial model (here kriging) is that it consists of errors: the  
550 estimations being further from samples are less certain. It is widely known that discounting the

551 measurement/model errors in true straight-line relationship between data can introduce artifacts.  
552 To consider the kriging variance in the comparisons we employed a Monte Carlo method on chi-  
553 square optimization which ultimately allowed us to not only provide a set of solutions within the  
554 range of the uncertainty of the kriging model, but also to assign smaller weights on gross estimates.

555 We further validated monthly-averaged Ozone Monitoring Instrument (OMI) tropospheric  
556 NO<sub>2</sub> columns using 11 Pandora Spectrometer Instrument (PSI) observations over Houston during  
557 NASA's DISCOVER-AQ campaign. A pixel-to-point comparison between two dataset suggested  
558 varying biases in OMI manifested in a slope far from the identity line. By contrast, the kriging  
559 estimate from the PSI measurements, convolved with the OMI spatial response function, resulted  
560 in an inter-comparison slope close to the unity line. This suggested that there was only a constant  
561 systematic bias ( $0.66 \pm 0.18 \times 10^{15}$  molecules cm<sup>-2</sup>) associated with the OMI observations which  
562 does not vary with tropospheric NO<sub>2</sub> column magnitudes.

563 The central tenants of satellite and model validation are pointwise measurements. Our  
564 experiments paved the way for a clear roadmap explaining how to transform these pointwise  
565 datasets to a comparable spatial scale relative to satellite (model) footprints. It is no longer  
566 necessary to ignore *the problem of scale*. The validation against point measurements can be  
567 carefully conducted in the following steps:

- 568
- 569 i. Construct the experimental semivariogram if the number of point measurements  
570 allows (usually  $\geq 3$  within the field; the field can vary depending on the length  
571 scale of the compound).
- 572 ii. Drop the quantitative assessment if the number of point measurements are  
573 insufficient to gain spatial variance and the prior knowledge suggests a high  
574 likelihood of spatial heterogeneity within the field.
- 575 iii. Choose an appropriate function to model the semivariogram.
- 576 iv. Estimate the field with kriging (or any other spatial estimator capable of digesting  
577 the semivariogram) and calculate the variance.
- 578 v. Estimate the optimum grid resolution of the estimate.
- 579 vi. Convolve the kriging estimate and its variance with the satellite (model) spatial  
580 response function (which is sensor specific).
- 581 vii. Conduct the direct comparison of the convolved kriged output and the satellite  
582 (model) considering their errors through a Monte Carlo (or a weighted least-squares  
583 method).
- 584

585 Recent advances in satellite trace gas retrievals and atmospheric models have helped  
586 extend our understanding of atmospheric chemistry but an important task before us in improving  
587 our knowledge on atmospheric composition is to embrace the semivariogram (or spatial auto-  
588 correlation) notion when it comes to validating satellites/models using pointwise measurements,  
589 so that we can have more robust quantitative applications of the data and models.

## 590 **Acknowledgement**

591 Amir Souri and Matthew Johnson were funded for this work through NASA's Aura Science Team  
592 (grant number: 80NSSC21K1333). Kang Sun acknowledges support by NASA's Atmospheric  
593 Composition: Modeling and Analysis (ACMAP) program (grant number: 80NSSC19K09). We  
594 thank many scientists whose concerns motivated us to tackle the presented problem. In particular,  
595 we thank Chris Chan Miller, Ron Cohen, Jeffrey Geddes, Gonzalo González Abad, Christian  
596 Hogrefe, Lukas Valin, and Huiqun (Helen) Wang.

597 **Author contributions**

598 AHS designed the research, executed the experiments, analyzed the data, made all figures, and  
599 wrote the paper. KS implemented the oversampling method, provided the spatial response  
600 functions, and oversampled TROPOMI data. KC, XL, and MSJ helped with the conceptualization  
601 of the study and the interpretation of the results. All authors contributed to discussions and edited  
602 the paper.  
603

604 **References**

- 605
- 606 Armstrong, M.: Is Research in Mining Geostats as Dead as a Dodo?, in: *Geostatistics for the*  
607 *Next Century: An International Forum in Honour of Michel David's Contribution to*  
608 *Geostatistics*, Montreal, 1993, edited by: Dimitrakopoulos, R., Springer Netherlands,  
609 Dordrecht, 303–312, [https://doi.org/10.1007/978-94-011-0824-9\\_34](https://doi.org/10.1007/978-94-011-0824-9_34), 1994.
- 610 Boersma, K. F., Eskes, H. J., Richter, A., De Smedt, I., Lorente, A., Beirle, S., van Geffen, J. H.  
611 G. M., Zara, M., Peters, E., Van Roozendaal, M., Wagner, T., Maasakkers, J. D., van der  
612 A, R. J., Nightingale, J., De Rudder, A., Irie, H., Pinardi, G., Lambert, J.-C., and  
613 Compernelle, S. C.: Improving algorithms and uncertainty estimates for satellite NO<sub>2</sub>  
614 retrievals: results from the quality assurance for the essential climate variables  
615 (QA4ECV) project, *Atmos. Meas. Tech.*, 11, 6651–6678, [https://doi.org/10.5194/amt-11-](https://doi.org/10.5194/amt-11-6651-2018)  
616 [6651-2018](https://doi.org/10.5194/amt-11-6651-2018), 2018.
- 617 Bryan, G. L.: Fluids in the universe: adaptive mesh refinement in cosmology, *Comput. Sci. Eng.*,  
618 1, 46–53, <https://doi.org/10.1109/5992.753046>, 1999.
- 619 Bucsela, E. J., Krotkov, N. A., Celarier, E. A., Lamsal, L. N., Swartz, W. H., Bhartia, P. K.,  
620 Boersma, K. F., Veefkind, J. P., Gleason, J. F., and Pickering, K. E.: A new stratospheric  
621 and tropospheric NO<sub>2</sub> retrieval algorithm for nadir-viewing satellite instruments:  
622 applications to OMI, *Atmos. Meas. Tech.*, 6, 2607–2626, [https://doi.org/10.5194/amt-6-](https://doi.org/10.5194/amt-6-2607-2013)  
623 [2607-2013](https://doi.org/10.5194/amt-6-2607-2013), 2013.
- 624 Chilès, J.-P. and Delfiner, P.: *Geostatistics: Modeling Spatial Uncertainty*, John Wiley & Sons,  
625 718 pp., 2009.
- 626 Choi, S., Lamsal, L. N., Follette-Cook, M., Joiner, J., Krotkov, N. A., Swartz, W. H., Pickering,  
627 K. E., Loughner, C. P., Appel, W., Pfister, G., Saide, P. E., Cohen, R. C., Weinheimer, A.  
628 J., and Herman, J. R.: Assessment of NO<sub>2</sub> observations during DISCOVER-AQ and  
629 KORUS-AQ field campaigns, *Atmos. Meas. Tech.*, 13, 2523–2546,  
630 <https://doi.org/10.5194/amt-13-2523-2020>, 2020.
- 631 Goldberg, D. L., Saide, P. E., Lamsal, L. N., de Foy, B., Lu, Z., Woo, J.-H., Kim, Y., Kim, J.,  
632 Gao, M., Carmichael, G., and Streets, D. G.: A top-down assessment using OMI NO<sub>2</sub>  
633 suggests an underestimate in the NO<sub>x</sub> emissions inventory in Seoul, South Korea, during  
634 KORUS-AQ, *Atmos. Chem. Phys.*, 19, 1801–1818, [https://doi.org/10.5194/acp-19-1801-](https://doi.org/10.5194/acp-19-1801-2019)  
635 [2019](https://doi.org/10.5194/acp-19-1801-2019), 2019.
- 636 Herman, J., Cede, A., Spinei, E., Mount, G., Tzortziou, M., and Abuhassan, N.: NO<sub>2</sub> column  
637 amounts from ground-based Pandora and MFDOAS spectrometers using the direct-sun  
638 DOAS technique: Intercomparisons and application to OMI validation, *J. Geophys. Res.*  
639 *Atmos.*, 114, <https://doi.org/10.1029/2009JD011848>, 2009.
- 640 Janjić, T., Bormann, N., Bocquet, M., Carton, J. A., Cohn, S. E., Dance, S. L., Losa, S. N.,  
641 Nichols, N. K., Potthast, R., Waller, J. A., and Weston, P.: On the representation error in

642 data assimilation, *Q. J. R. Meteorol. Soc.*, 144, 1257–1278,  
643 <https://doi.org/10.1002/qj.3130>, 2018.

644 Judd, L. M., Al-Saadi, J. A., Szykman, J. J., Valin, L. C., Janz, S. J., Kowalewski, M. G., Eskes,  
645 H. J., Veefkind, J. P., Cede, A., Mueller, M., Gebetsberger, M., Swap, R., Pierce, R. B.,  
646 Nowlan, C. R., Abad, G. G., Nehrir, A., and Williams, D.: Evaluating Sentinel-5P  
647 TROPOMI tropospheric NO<sub>2</sub> column densities with airborne and Pandora spectrometers  
648 near New York City and Long Island Sound, *Atmos. Meas. Tech.*, 13, 6113–6140,  
649 <https://doi.org/10.5194/amt-13-6113-2020>, 2020.

650 Kendon, E. J., Roberts, N. M., Fowler, H. J., Roberts, M. J., Chan, S. C., and Senior, C. A.:  
651 Heavier summer downpours with climate change revealed by weather forecast resolution  
652 model, *Nature Clim. Change*, 4, 570–576, <https://doi.org/10.1038/nclimate2258>, 2014.

653 Kim, H. C., Lee, S.-M., Chai, T., Ngan, F., Pan, L., and Lee, P.: A Conservative Downscaling of  
654 Satellite-Detected Chemical Compositions: NO<sub>2</sub> Column Densities of OMI, GOME-2,  
655 and CMAQ, *Remote Sens.*, 10, 1001, <https://doi.org/10.3390/rs10071001>, 2018.

656 Laughner, J. L., Zhu, Q., and Cohen, R. C.: The Berkeley High Resolution Tropospheric NO<sub>2</sub>  
657 product, *Earth Syst. Sci. Data*, 10, 2069–2095, [https://doi.org/10.5194/essd-10-2069-](https://doi.org/10.5194/essd-10-2069-2018)  
658 2018, 2018.

659 Li, R., Cui, L., Meng, Y., Zhao, Y., and Fu, H.: Satellite-based prediction of daily SO<sub>2</sub> exposure  
660 across China using a high-quality random forest-spatiotemporal Kriging (RF-STK) model  
661 for health risk assessment, *Atmos. Environ.*, 208, 10–19,  
662 <https://doi.org/10.1016/j.atmosenv.2019.03.029>, 2019.

663 Matheron, G.: Principles of geostatistics, *Econ. Geol.*, 58, 1246–1266,  
664 <https://doi.org/10.2113/gsecongeo.58.8.1246>, 1963.

665 Mazzuca, G. M., Ren, X., Loughner, C. P., Estes, M., Crawford, J. H., Pickering, K. E.,  
666 Weinheimer, A. J., and Dickerson, R. R.: Ozone production and its sensitivity to NO<sub>x</sub>  
667 and VOCs: results from the DISCOVER-AQ field experiment, Houston 2013, *Atmos.*  
668 *Chem. Phys.*, 16, 14463–14474, <https://doi.org/10.5194/acp-16-14463-2016>, 2016.

669 Nowlan, C. R., Liu, X., Janz, S. J., Kowalewski, M. G., Chance, K., Follette-Cook, M. B., Fried,  
670 A., González Abad, G., Herman, J. R., Judd, L. M., Kwon, H.-A., Loughner, C. P.,  
671 Pickering, K. E., Richter, D., Spinei, E., Walega, J., Weibring, P., and Weinheimer, A. J.:  
672 Nitrogen dioxide and formaldehyde measurements from the GEOstationary Coastal and  
673 Air Pollution Events (GEO-CAPE) Airborne Simulator over Houston, Texas, *Atmos.*  
674 *Meas. Tech.*, 11, 5941–5964, <https://doi.org/10.5194/amt-11-5941-2018>, 2018.

675 Nowlan, C. R., Liu, X., Leitch, J. W., Chance, K., González Abad, G., Liu, C., Zoogman, P.,  
676 Cole, J., Delker, T., Good, W., Murcray, F., Ruppert, L., Soo, D., Follette-Cook, M. B.,  
677 Janz, S. J., Kowalewski, M. G., Loughner, C. P., Pickering, K. E., Herman, J. R., Beaver,  
678 M. R., Long, R. W., Szykman, J. J., Judd, L. M., Kelley, P., Luke, W. T., Ren, X., and  
679 Al-Saadi, J. A.: Nitrogen dioxide observations from the Geostationary Trace gas and  
680 Aerosol Sensor Optimization (GeoTASO) airborne instrument: Retrieval algorithm and  
681 measurements during DISCOVER-AQ Texas 2013, *Atmos. Meas. Tech.*, 9, 2647–2668,  
682 <https://doi.org/10.5194/amt-9-2647-2016>, 2016.

683 Onn, F. and Zebker, H. A.: Correction for interferometric synthetic aperture radar atmospheric  
684 phase artifacts using time series of zenith wet delay observations from a GPS network, *J.*  
685 *Geophys. Res. Solid Earth*, 111, <https://doi.org/10.1029/2005JB004012>, 2006.

686 Pan, S., Choi, Y., Jeon, W., Roy, A., Westenbarger, D. A., and Kim, H. C.: Impact of high-  
687 resolution sea surface temperature, emission spikes and wind on simulated surface ozone



688 in Houston, Texas during a high ozone episode, *Atmos. Environ.*, 152, 362–376,  
689 <https://doi.org/10.1016/j.atmosenv.2016.12.030>, 2017b.

690 Pan, S., Choi, Y., Roy, A., and Jeon, W.: Allocating emissions to 4 km and 1 km horizontal  
691 spatial resolutions and its impact on simulated NO<sub>x</sub> and O<sub>3</sub> in Houston, TX, *Atmos.*  
692 *Environ.*, 164, 398–415, <https://doi.org/10.1016/j.atmosenv.2017.06.026>, 2017a.

693 Pan, S., Choi, Y., Roy, A., Li, X., Jeon, W., and Souri, A. H.: Modeling the uncertainty of  
694 several VOC and its impact on simulated VOC and ozone in Houston, Texas, *Atmos.*  
695 *Environ.*, 120, 404–416, <https://doi.org/10.1016/j.atmosenv.2015.09.029>, 2015.

696 Platt, U., Wagner, T., Kuhn, J., and Leisner, T.: The “ideal” spectrograph for atmospheric  
697 observations, *Atmos. Meas. Tech.*, 14, 6867–6883, [https://doi.org/10.5194/amt-14-6867-](https://doi.org/10.5194/amt-14-6867-2021)  
698 2021, 2021.

699 Pyrcz MJ, Deutsch CV. The whole story on the hole effect. *Geostatistical Association of*  
700 *Australasia, Newsletter*. 2003 May;18:3-5.

701 Rennen, G.: Subset Selection from Large Datasets for Kriging Modeling, *Social Science*  
702 *Research Network*, Rochester, NY, <https://doi.org/10.2139/ssrn.1104595>, 2008.

703 Russell, A. R., Perring, A. E., Valin, L. C., Bucseca, E. J., Browne, E. C., Wooldridge, P. J., and  
704 Cohen, R. C.: A high spatial resolution retrieval of NO<sub>2</sub> column densities from OMI:  
705 method and evaluation, *Atmos. Chem. Phys.*, 11, 8543–8554,  
706 <https://doi.org/10.5194/acp-11-8543-2011>, 2011.

707 Shah, V., Jacob, D. J., Li, K., Silvern, R. F., Zhai, S., Liu, M., Lin, J., and Zhang, Q.: Effect of  
708 changing NO<sub>x</sub> lifetime on the seasonality and long-term trends of satellite-observed  
709 tropospheric NO<sub>2</sub> columns over China, *Atmos. Chem. Phys.*, 20, 1483–1495,  
710 <https://doi.org/10.5194/acp-20-1483-2020>, 2020.

711 Souri, A. H., Choi, Y., Jeon, W., Li, X., Pan, S., Diao, L., and Westenbarger, D. A.: Constraining  
712 NO<sub>x</sub> emissions using satellite NO<sub>2</sub> measurements during 2013 DISCOVER-AQ Texas  
713 campaign, *Atmos. Environ.*, 131, 371–381,  
714 <https://doi.org/10.1016/j.atmosenv.2016.02.020>, 2016.

715 Souri, A. H., Choi, Y., Kodros, J. K., Jung, J., Shpund, J., Pierce, J. R., Lynn, B. H., Khain, A.,  
716 and Chance, K.: Response of Hurricane Harvey’s rainfall to anthropogenic aerosols: A  
717 sensitivity study based on spectral bin microphysics with simulated aerosols, *Atmos.*  
718 *Res.*, 242, 104965, <https://doi.org/10.1016/j.atmosres.2020.104965>, 2020a.

719 Souri, A. H., Choi, Y., Pan, S., Curci, G., Nowlan, C. R., Janz, S. J., Kowalewski, M. G., Liu, J.,  
720 Herman, J. R., and Weinheimer, A. J.: First Top-Down Estimates of Anthropogenic NO<sub>x</sub>  
721 Emissions Using High-Resolution Airborne Remote Sensing Observations, *J. Geophys.*  
722 *Res. Atmos.*, 123, 3269–3284, <https://doi.org/10.1002/2017JD028009>, 2018.

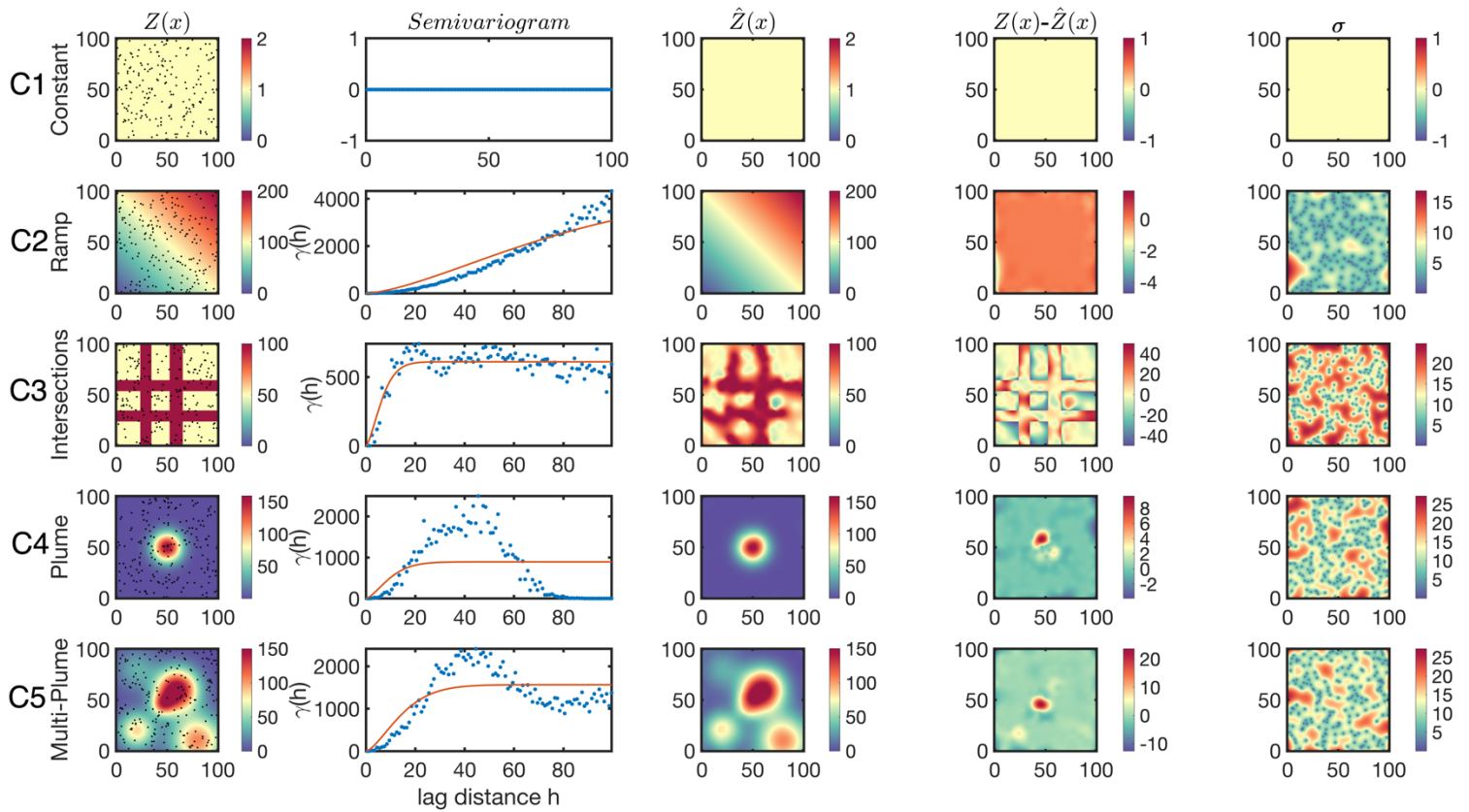
723 Souri, A. H., Nowlan, C. R., Wolfe, G. M., Lamsal, L. N., Chan Miller, C. E., Abad, G. G., Janz,  
724 S. J., Fried, A., Blake, D. R., Weinheimer, A. J., Diskin, G. S., Liu, X., and Chance, K.:  
725 Revisiting the effectiveness of HCHO/NO<sub>2</sub> ratios for inferring ozone sensitivity to its  
726 precursors using high resolution airborne remote sensing observations in a high ozone  
727 episode during the KORUS-AQ campaign, *Atmos. Environ.*, 224, 117341,  
728 <https://doi.org/10.1016/j.atmosenv.2020.117341>, 2020b.

729 Sun, K., Zhu, L., Cady-Pereira, K., Chan Miller, C., Chance, K., Clarisse, L., Coheur, P.-F.,  
730 González Abad, G., Huang, G., Liu, X., Van Damme, M., Yang, K., and Zondlo, M.: A  
731 physics-based approach to oversample multi-satellite, multispecies observations to a  
732 common grid, *Atmos. Meas. Tech.*, 11, 6679–6701, [https://doi.org/10.5194/amt-11-6679-](https://doi.org/10.5194/amt-11-6679-2018)  
733 2018, 2018.

- 734 Swall, J. L. and Foley, K. M.: The impact of spatial correlation and incommensurability on  
735 model evaluation, *Atmos. Environ.*, 43, 1204–1217,  
736 <https://doi.org/10.1016/j.atmosenv.2008.10.057>, 2009.
- 737 Tadić, J. M., Michalak, A. M., Iraci, L., Ilić, V., Biraud, S. C., Feldman, D. R., Bui, T., Johnson,  
738 M. S., Loewenstein, M., Jeong, S., Fischer, M. L., Yates, E. L., and Ryoo, J.-M.: Elliptic  
739 Cylinder Airborne Sampling and Geostatistical Mass Balance Approach for Quantifying  
740 Local Greenhouse Gas Emissions, *Environ. Sci. Technol.*, 51, 10012–10021,  
741 <https://doi.org/10.1021/acs.est.7b03100>, 2017.
- 742 Tang, W., Edwards, D. P., Emmons, L. K., Worden, H. M., Judd, L. M., Lamsal, L. N., Al-Saadi,  
743 J. A., Janz, S. J., Crawford, J. H., Deeter, M. N., Pfister, G., Buchholz, R. R., Gaubert, B.,  
744 and Nowlan, C. R.: Assessing sub-grid variability within satellite pixels over urban  
745 regions using airborne mapping spectrometer measurements, *Atmos. Meas. Tech.*, 14,  
746 4639–4655, <https://doi.org/10.5194/amt-14-4639-2021>, 2021.
- 747 Valin, L. C., Russell, A. R., Hudman, R. C., and Cohen, R. C.: Effects of model resolution on the  
748 interpretation of satellite NO<sub>2</sub> observations, *Atmos. Chem. Phys.*, 11, 11647–11655,  
749 <https://doi.org/10.5194/acp-11-11647-2011>, 2011.
- 750 Vinken, G. C. M., Boersma, K. F., Jacob, D. J., and Meijer, E. W.: Accounting for non-linear  
751 chemistry of ship plumes in the GEOS-Chem global chemistry transport model, *Atmos.*  
752 *Chem. Phys.*, 11, 11707–11722, <https://doi.org/10.5194/acp-11-11707-2011>, 2011.
- 753 Wang, P., Piters, A., van Geffen, J., Tuinder, O., Stammes, P., and Kinne, S.: Shipborne MAX-  
754 DOAS measurements for validation of TROPOMI NO<sub>2</sub> products, *Atmos. Meas. Tech.*,  
755 13, 1413–1426, <https://doi.org/10.5194/amt-13-1413-2020>, 2020.
- 756 Wang, Y., Sabatino, S. D., Martilli, A., Li, Y., Wong, M. S., Gutiérrez, E., and Chan, P. W.:  
757 Impact of land surface heterogeneity on urban heat island circulation and sea-land breeze  
758 circulation in Hong Kong, *J. Geophys. Res. Atmos.*, 122, 4332–4352,  
759 <https://doi.org/10.1002/2017JD026702>, 2017.
- 760 Wolfe, G. M., Nicely, J. M., Clair, J. M. S., Hanisco, T. F., Liao, J., Oman, L. D., Brune, W. B.,  
761 Miller, D., Thames, A., Abad, G. G., Ryerson, T. B., Thompson, C. R., Peischl, J.,  
762 McKain, K., Sweeney, C., Wennberg, P. O., Kim, M., Crouse, J. D., Hall, S. R.,  
763 Ullmann, K., Diskin, G., Bui, P., Chang, C., and Dean-Day, J.: Mapping hydroxyl  
764 variability throughout the global remote troposphere via synthesis of airborne and  
765 satellite formaldehyde observations, *PNAS*, 116, 11171–11180,  
766 <https://doi.org/10.1073/pnas.1821661116>, 2019.
- 767 Wu, C.-D., Zeng, Y.-T., and Lung, S.-C. C.: A hybrid kriging/land-use regression model to  
768 assess PM<sub>2.5</sub> spatial-temporal variability, *Sci. Total Environ*, 645, 1456–1464,  
769 <https://doi.org/10.1016/j.scitotenv.2018.07.073>, 2018.
- 770 Yu, K., Jacob, D. J., Fisher, J. A., Kim, P. S., Marais, E. A., Miller, C. C., Travis, K. R., Zhu, L.,  
771 Yantosca, R. M., Sulprizio, M. P., Cohen, R. C., Dibb, J. E., Fried, A., Mikoviny, T.,  
772 Ryerson, T. B., Wennberg, P. O., and Wisthaler, A.: Sensitivity to grid resolution in the  
773 ability of a chemical transport model to simulate observed oxidant chemistry under high-  
774 isoprene conditions, *Atmos. Chem. Phys.*, 16, 4369–4378, <https://doi.org/10.5194/acp-16-4369-2016>, 2016.
- 776 Zhan, Y., Luo, Y., Deng, X., Zhang, K., Zhang, M., Grieneisen, M. L., and Di, B.: Satellite-  
777 Based Estimates of Daily NO<sub>2</sub> Exposure in China Using Hybrid Random Forest and  
778 Spatiotemporal Kriging Model, *Environ. Sci. Technol.*, 52, 4180–4189,  
779 <https://doi.org/10.1021/acs.est.7b05669>, 2018.

780 Zhao, X., Griffin, D., Fioletov, V., McLinden, C., Cede, A., Tiefengraber, M., Müller, M.,  
781 Bognar, K., Strong, K., Boersma, F., Eskes, H., Davies, J., Ogyu, A., and Lee, S. C.:  
782 Assessment of the quality of TROPOMI high-spatial-resolution NO<sub>2</sub> data products in the  
783 Greater Toronto Area, *Atmos. Meas. Tech.*, 13, 2131–2159, [https://doi.org/10.5194/amt-](https://doi.org/10.5194/amt-13-2131-2020)  
784 13-2131-2020, 2020.  
785

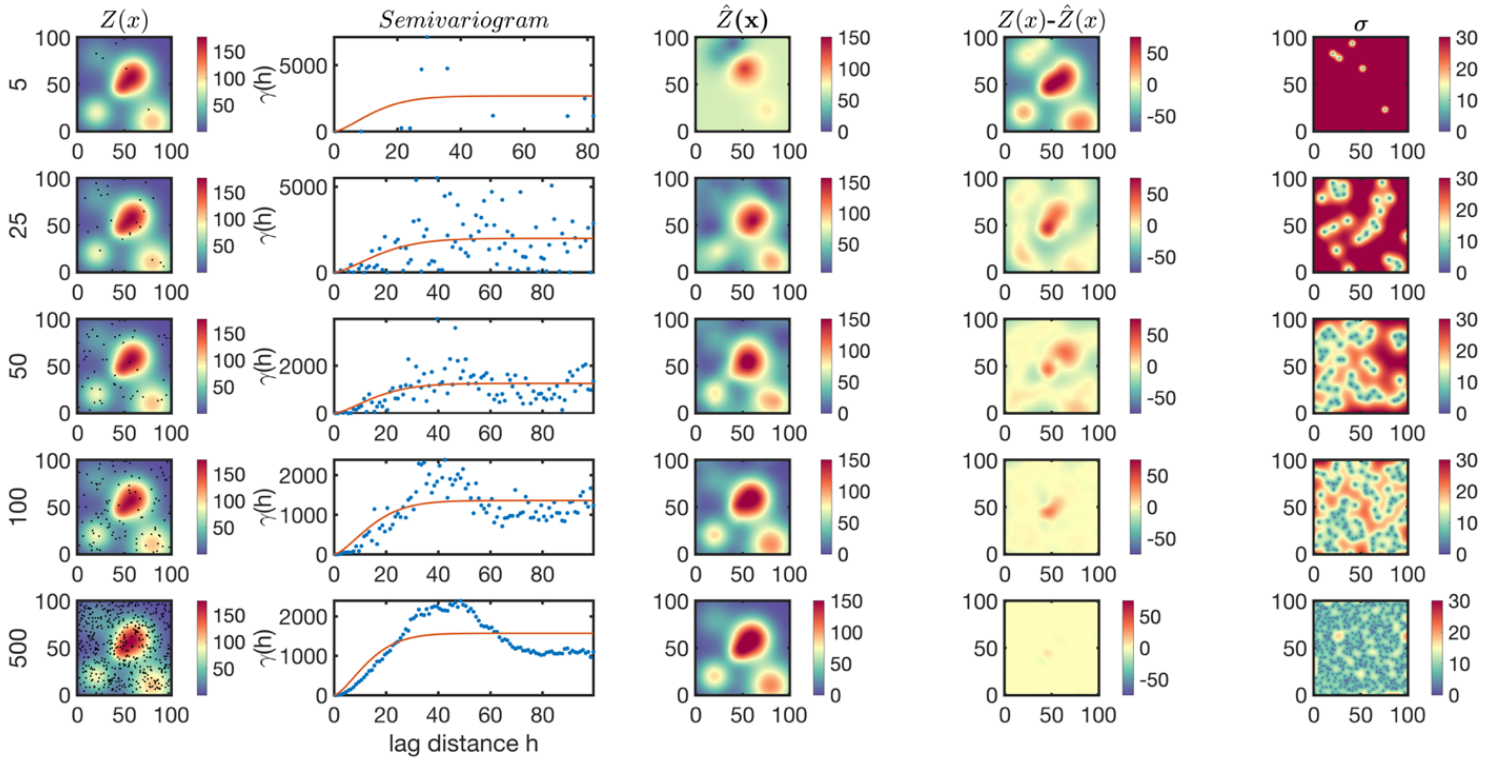
786 Figures:



788 **Figure 1.** (first column) Five theoretical fields randomly sampled with 200 points (dots), namely,  
 789 a constant field (C1), a ramp starting from zero in the lower left to higher values in the upper right  
 790 (C2), an intersection with concentrated values in four corridors (C3), a Gaussian plume placed in  
 791 the center (C4), and multiple Gaussian plumes spread over the entire domain (C5). (second column)  
 792 the corresponding isotropic semivariograms computed based on Eq.2; the red line shows the stable  
 793 Gaussian fitted to the semivariogram based on Levenberg-Marquardt method. (third column) The  
 794 kriging estimate at the same resolution of the truth (i.e.,  $1 \times 1$ ) based on Eq.6. (fourth column) The  
 795 difference between the estimate and the truth. (fifth column) the kriging standard error based on  
 796 Eq.11.

797

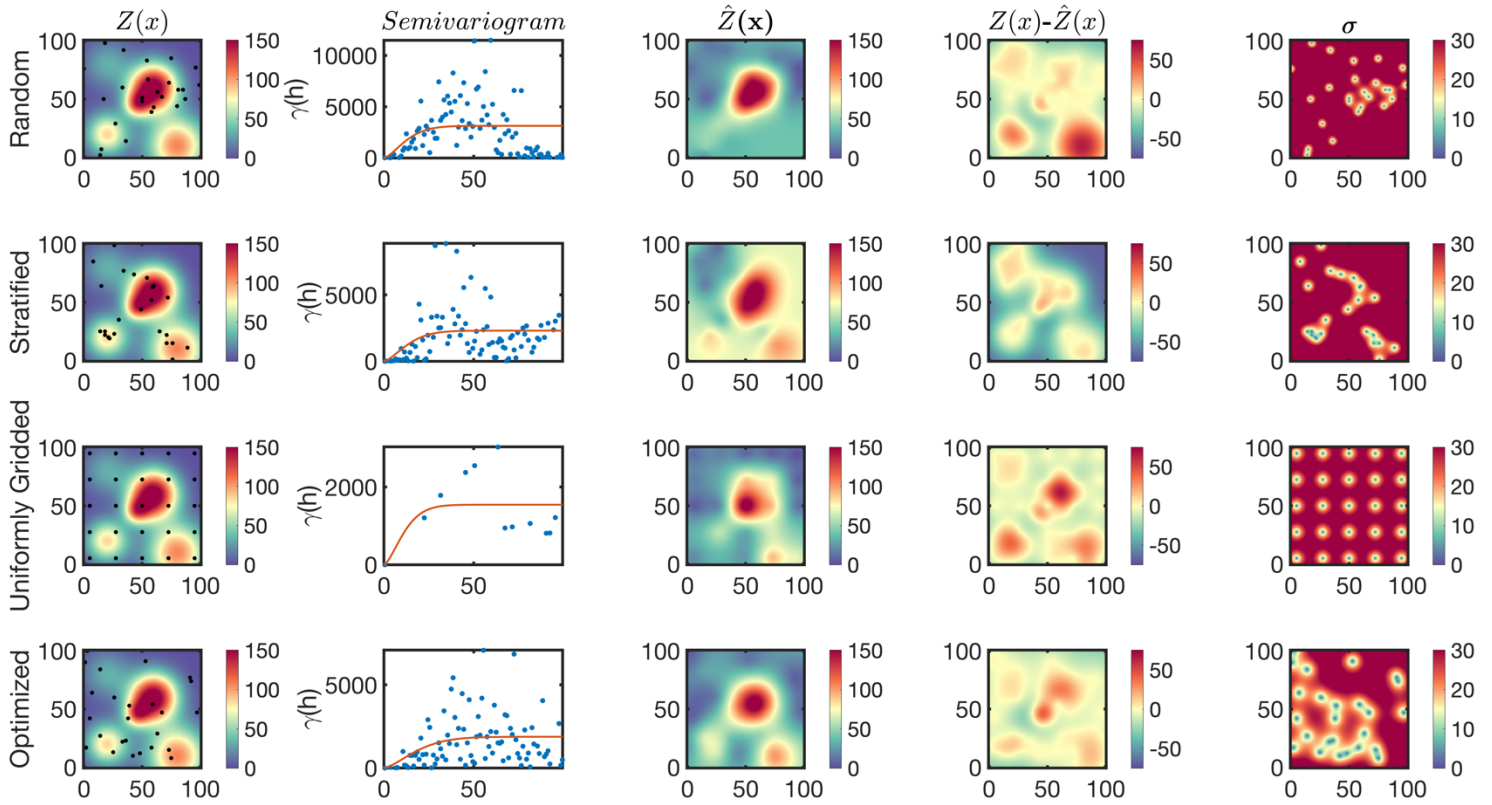
798



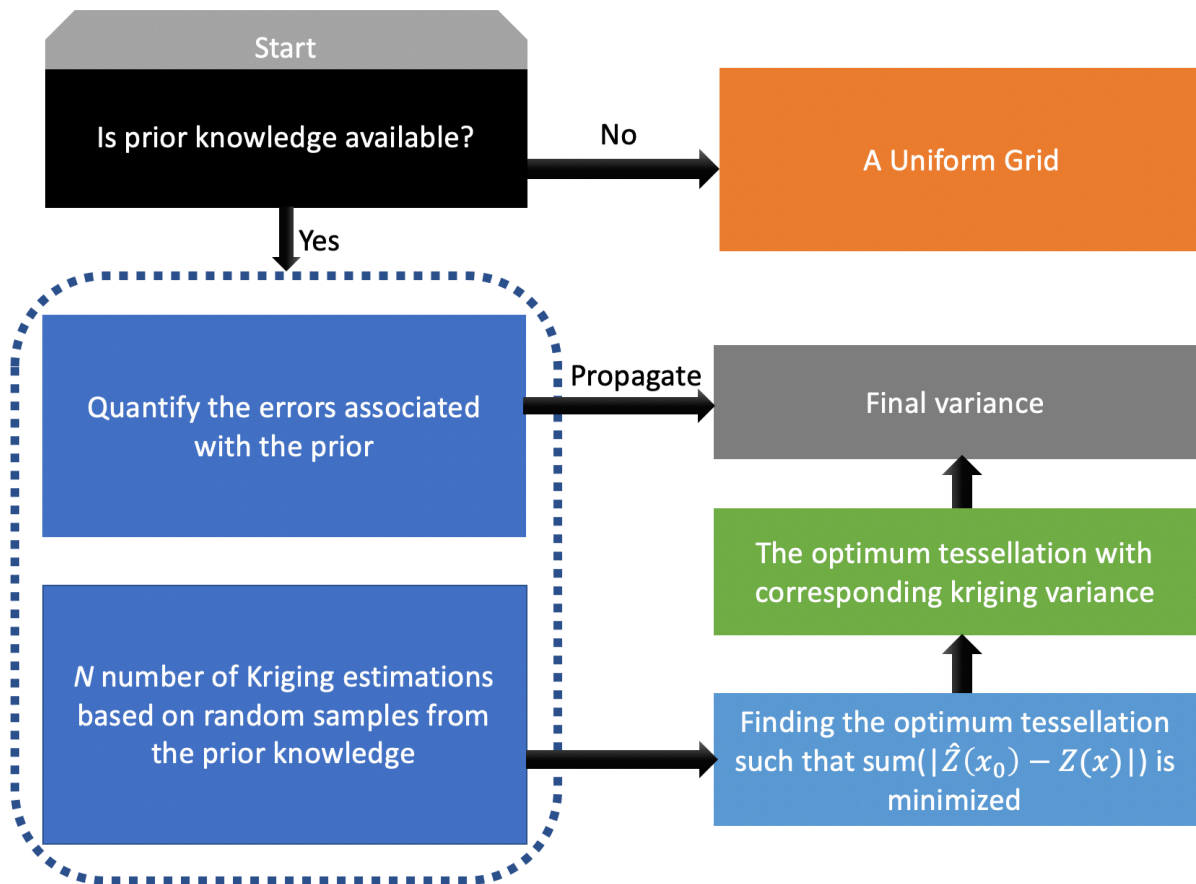
800 **Figure 2.** (first column) The multi-plume case (C5) randomly sampled with different number of  
801 samples (5, 25, 50, 100, and 500), (second column) the corresponding isotropic semivariogram,  
802 (third column) the kriging estimate, (fourth column) the difference between the estimate and the  
803 truth, and (fifth column) the kriging standard error.

804

805



807 **Figure 3.** The multi-plume case (C5) randomly sampled by four different sampling strategies  
 808 using a constant number of samples (25). The sampling strategies include purely random (first  
 809 row), stratified random (second row), uniform grids (third row), and an optimized tessellation  
 810 proposed based on kriging (fourth row). Columns represent the truth, the isotropic semivariogram,  
 811 the kriging estimate, the difference between the estimate and the truth, and the kriging standard  
 812 error.  
 813

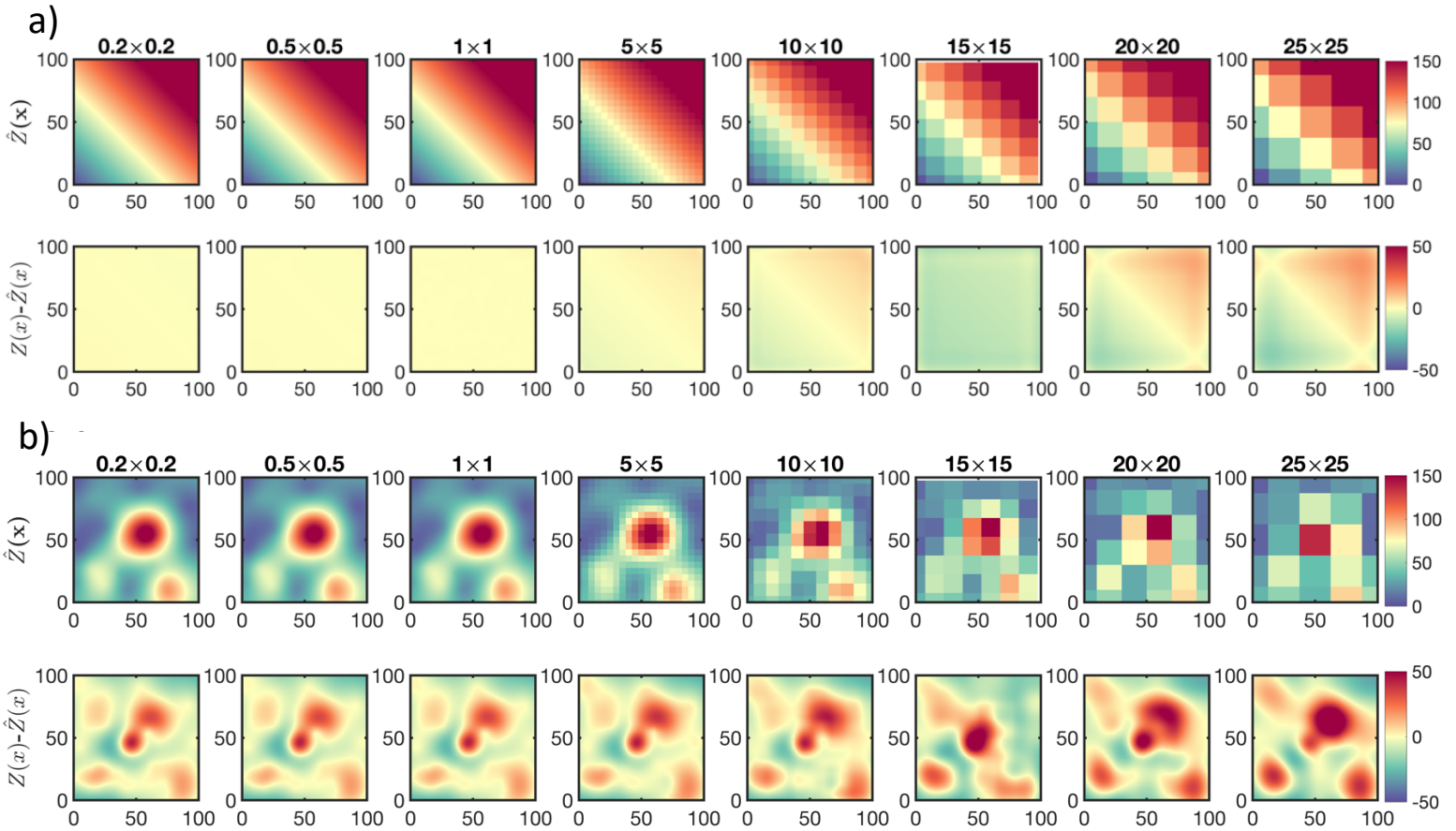


814

815 **Figure 4.** A schematic illustrating a framework for optimum sampling (tessellation) strategy. The  
 816 prior knowledge refers to any data being able of describing our quantity of interest including site-  
 817 visits, theoretical models, satellite observations, emissions, and etc.

818

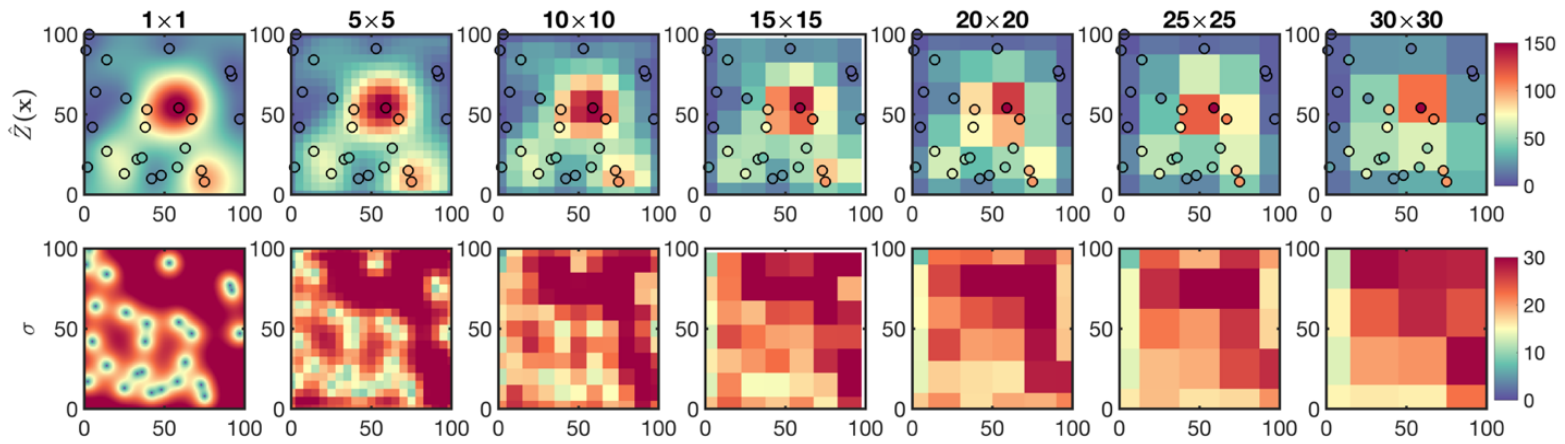
819



822 **Figure 5.** Finding an optimum grid size for kriging. (a) The kriging estimates of the ramp (C2) at  
 823 different grid resolutions ranging from  $25 \times 25$  pixel to  $0.2 \times 0.2$ . (b) The kriging estimates of the  
 824 multi-plume (C5) with optimized samples shown in Figure 3 for different grid resolutions. C2 is  
 825 more homogeneous than C5, as a result, it is less sensitive to the resolution of the kriging  
 826 estimate. The optimum grid resolution for C2 is  $10 \times 10$ , whereas it is  $1 \times 1$  for C5. These numbers  
 827 are based on observing negligible difference ( $< 1\%$ ) between the kriging estimate at the optimum  
 828 resolution and the one computed at a finer resolution step. We call the optimum output for C5 as  
 829 C5opt.

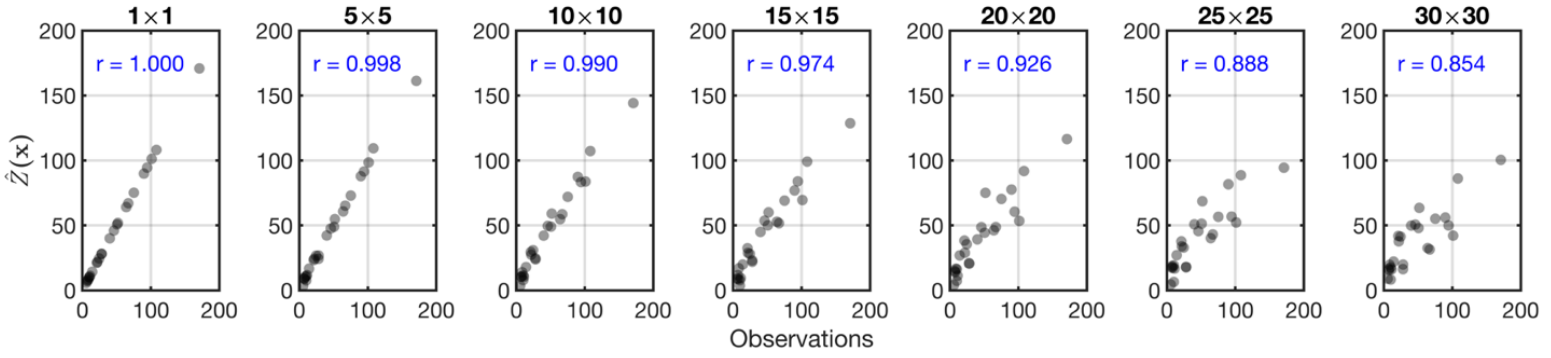
830  
 831  
 832





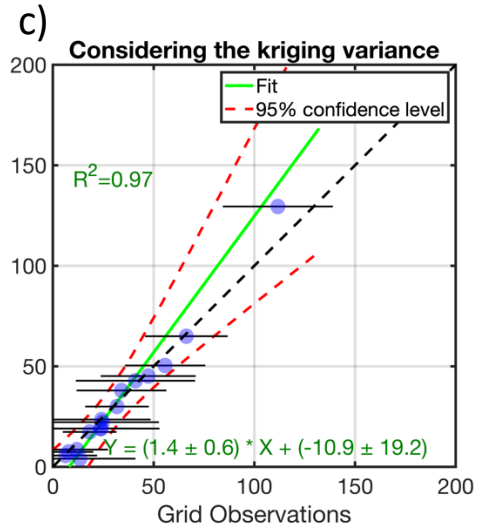
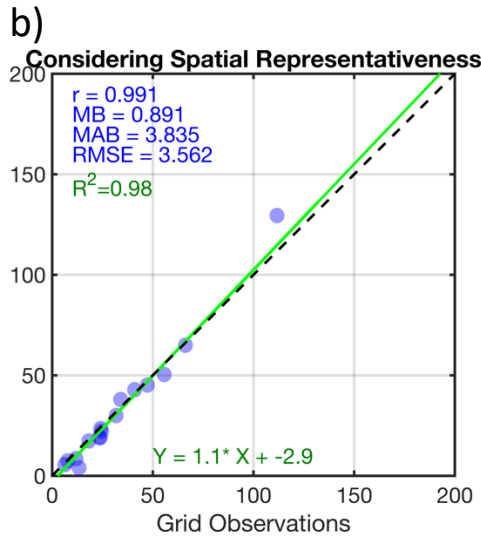
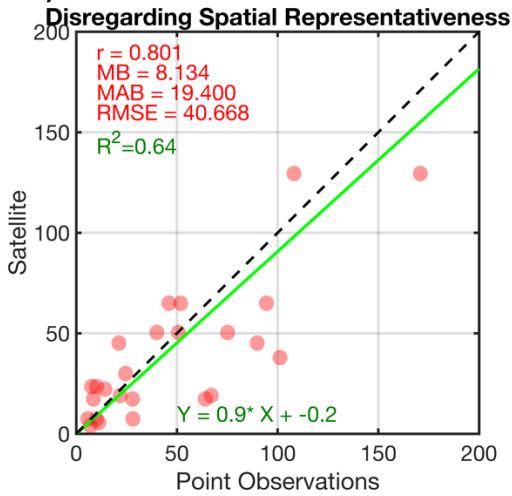
834 **Figure 6.** (first row) C5Opt outputs convolved with an ideal box kernel with different sizes ( $1 \times 1$   
 835 up to  $30 \times 30$ ) overlaid by the C5Opt optimum samples. (second row) the associated kriging errors  
 836 convolved with the same kernel. The coarser the resolution is, the larger the discrepancy between  
 837 the samples and the estimates is.

838  
 839  
 840  
 841



843 **Figure 7.** Illustrating the problem of spatial scale: comparisons of the kriging estimates at seven  
 844 different spatial scales with the samples used for the C5opt estimation. The perceived  
 845 discrepancies are purely due to the spatial representativeness.  
 846

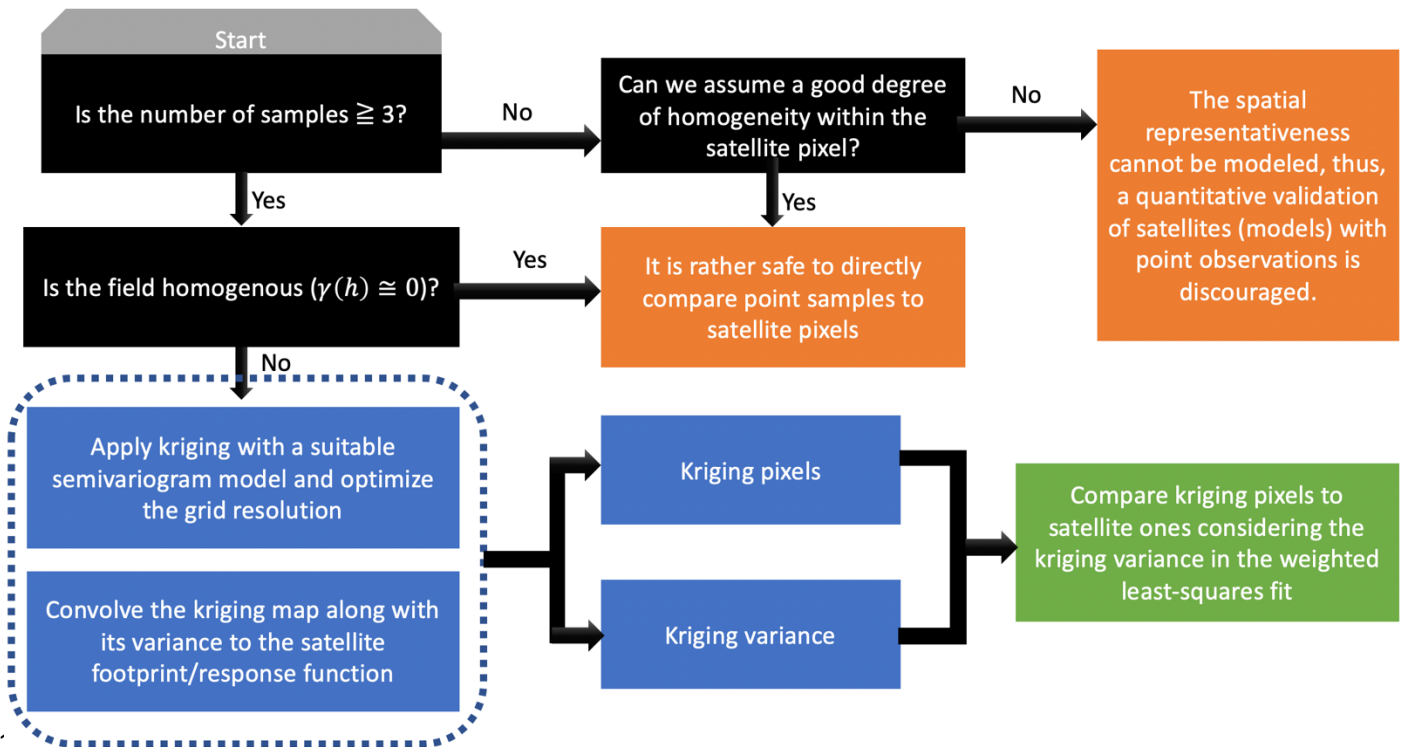
847  
848  
849  
a)



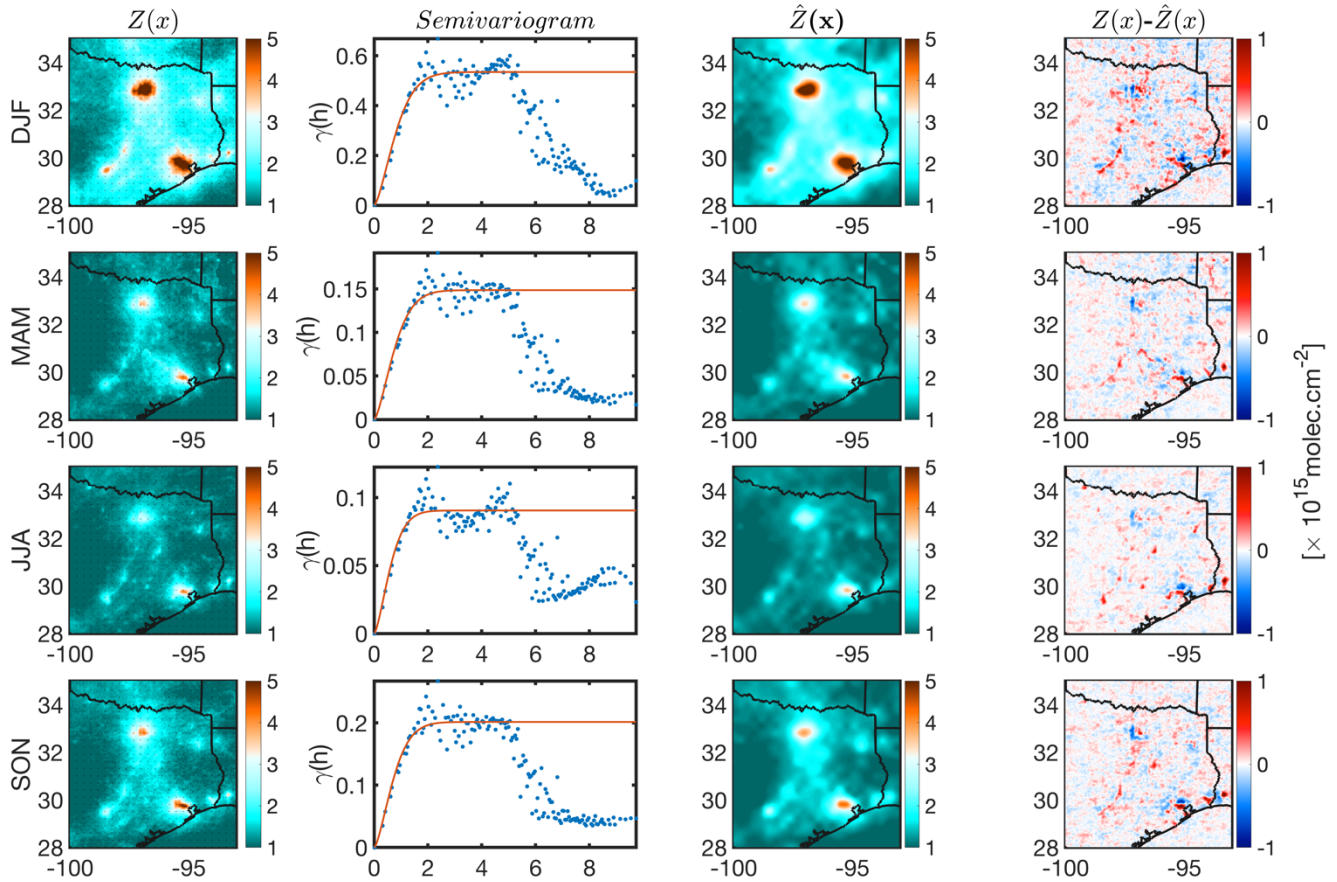
851

852 **Figure 8.** (a) the direct comparison of pseudo observations of a satellite observing the C5 case at  
853 30×30 resolution versus the 25 samples used for C5opt. (b) same for y-axis, but the point samples  
854 are transformed to grid boxes using kriging convolved with the satellite spatial response function  
855 (ideal box with 30×30 kernel size). The differences in statistics between these two experiments  
856 speak to the problem of scale. (b) ignores the kriging errors but (c) incorporates them using a  
857 Monte Carlo method. Note that the best linear fit has changed indicating that the consideration of  
858 the kriging variance is critical. MB = mean bias (point minus satellite), MAB = mean absolute  
859 bias, RMSE = root mean square error,  $R^2$  = coefficient of determination.

860

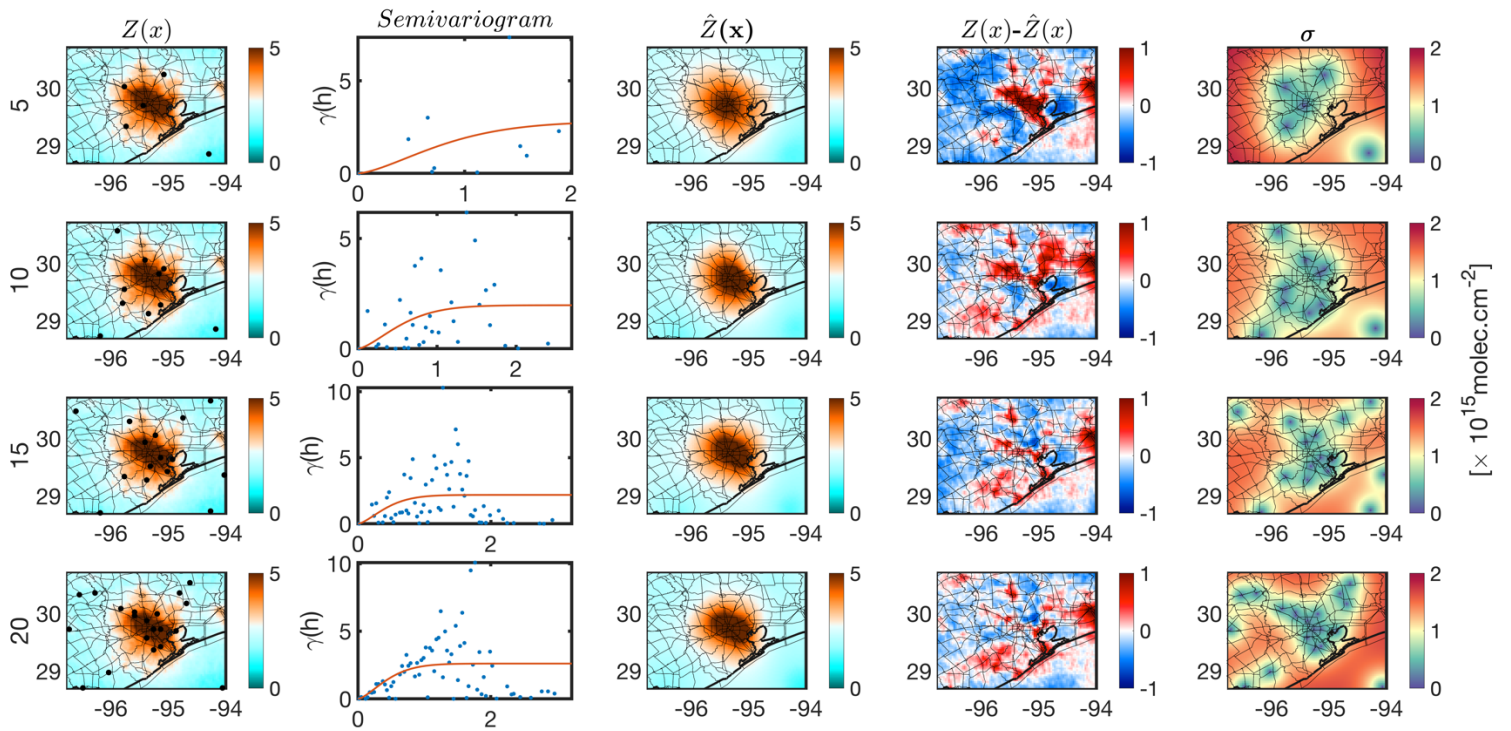


86:  
 862 **Figure 9.** The proposed roadmap for transforming pointwise measurements to gridded data in  
 863 satellite (model) validation.  
 864

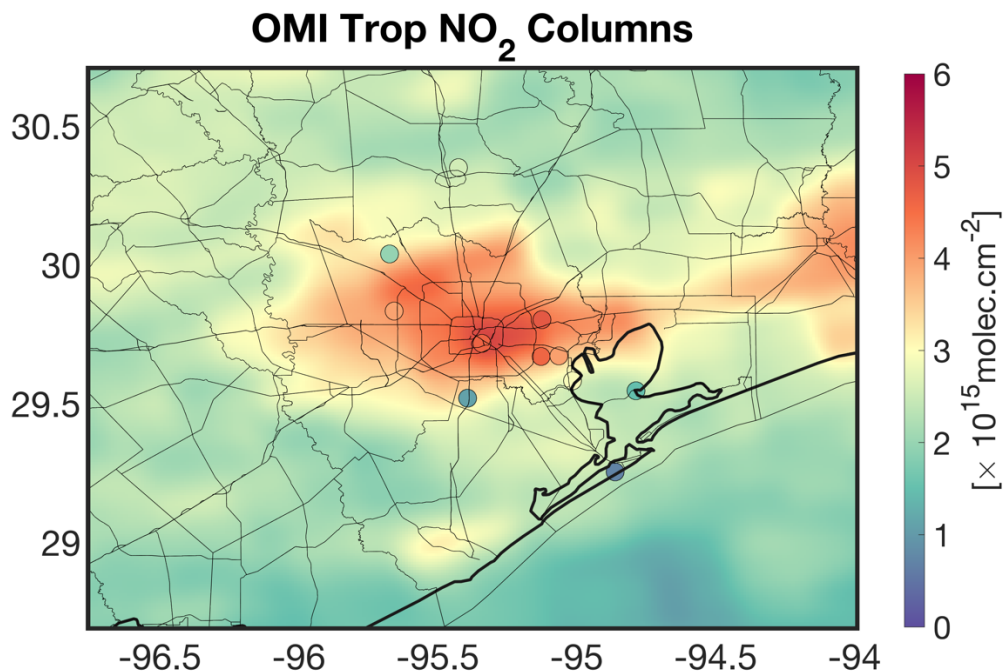


8  
 866 **Figure 10.** (first column) The spatial distribution of TROPOMI tropospheric NO<sub>2</sub> columns  
 867 oversampled in four different seasons at 3×3 km<sup>2</sup> spatial resolution. (second column) The  
 868 corresponding semivariogram from samples selected from uniform 30×30 km<sup>2</sup> blocks (shown  
 869 with black dots in the first column) along the fitted stable Gaussian model (red line). (third  
 870 column) the kriging estimates, and (fourth column) their differences with respect to the  
 871 observations.  
 872  
 873

874



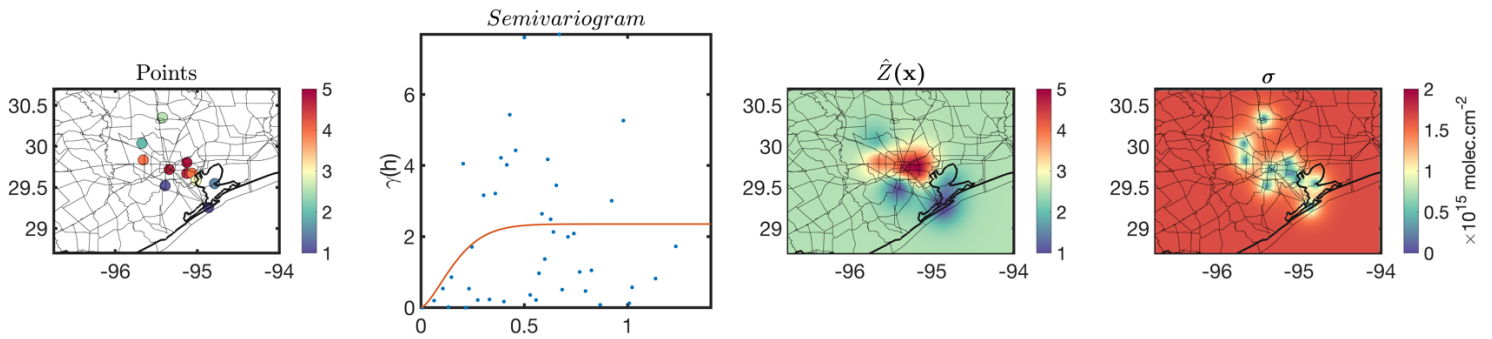
876 **Figure 11.** Finding an optimum sample tessellation for wintertime over Houston given different  
877 number of spectrometers (5, 10, 15, and 20).  
878



879  
 880  
 881  
 882  
 883  
 884  
 885  
 886  
 887

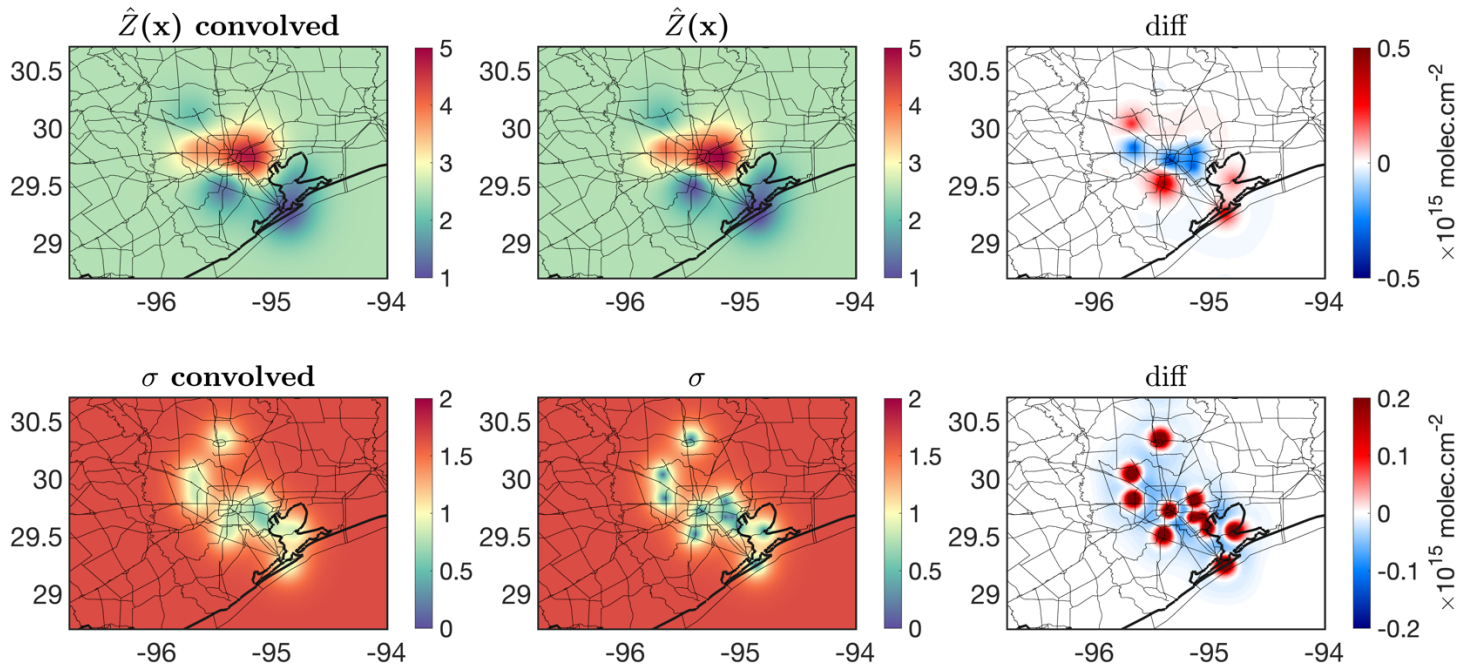
**Figure 12.** The spatial distribution of OMI tropospheric NO<sub>2</sub> columns oversampled at the resolution at  $0.2 \times 0.2^\circ$  over Houston in September 2013. The plot is overlaid by surface Pandora spectrometer instrument averaged over the same month. The surface measurements originally measured the total columns, therefore we subtract the stratospheric columns provided by the OMI data ( $2.8 \pm 0.16 \times 10^{15}$  molecules  $\text{cm}^{-2}$ ) from the total columns to focus on the tropospheric part.

888  
889

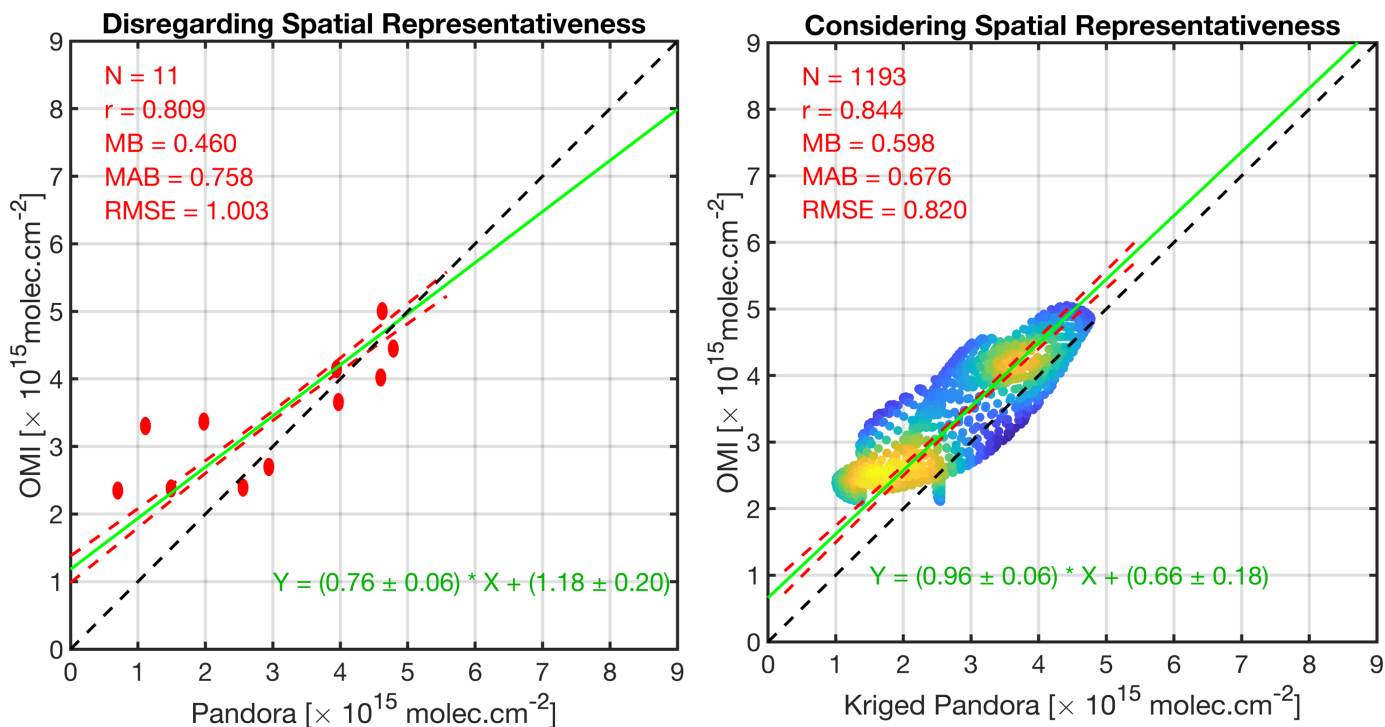


891 **Figure 13.** The Pandora tropospheric NO<sub>2</sub> measurements (made from subtracting the total columns  
892 from the OMI stratospheric NO<sub>2</sub> columns) during September 2013, the corresponding  
893 semivariogram, the kriging estimates, and the kriging standard errors. Note that the semivariogram  
894 suggests a large degree of spatial heterogeneity occurring at different spatial scales.  
895





8  
 897 **Figure 14.** Convolving both kriging estimates and errors with the OMI spatial response function  
 898 formulated in Sun et al. [2018]. The differences against the pre-convolved fields are also depicted.  
 899



901  
 902 **Figure 15.** (left): the direct comparison of OMI tropospheric NO<sub>2</sub> columns with 11 pointwise  
 903 Pandora measurements in September 2013 over Houston. (right) same for y-axis, but the PSI  
 904 measurements are translated to grid boxes using kriging convolved with the OMI spatial response  
 905 function. PSI tropospheric NO<sub>2</sub> columns are estimated based on subtracting the OMI stratospheric  
 906 NO<sub>2</sub> columns ( $2.8 \pm 0.16 \times 10^{15}$  molecules  $\text{cm}^{-2}$ ) from the total columns. We only consider kriging  
 907 estimates whose errors are below  $1.2 \times 10^{15}$  molecules  $\text{cm}^{-2}$ . The kriging variance is also considered  
 908 using the Monte Carlo method applied on  $\chi^2$ . The slope has improved after considering the  
 909 modeled spatial representativeness. MB = mean bias (OMI vs Pandora), MAB = mean absolute  
 910 bias, RMSE = root mean square error.  
 911