

# 1 Cluster-based characterization of multi-dimensional tropospheric 2 ozone variability in coastal regions: an analysis of lidar 3 measurements and model results

4 Claudia Bernier<sup>1</sup>, Yuxuan Wang<sup>1</sup>, Guillaume Gronoff<sup>2,3</sup>, Timothy Berkoff<sup>2</sup>, K. Emma Knowland<sup>4,5</sup>, John T.  
5 Sullivan<sup>4</sup>, Ruben Delgado<sup>6,7</sup>, Vanessa Caicedo<sup>6,7</sup>, Brian Carroll<sup>2,6</sup>

6 <sup>1</sup>Department of Earth and Atmospheric Science, University of Houston, Houston, Texas, USA

7 <sup>2</sup>NASA Langley Research Center, Hampton, VA, 23666, USA

8 <sup>3</sup>Science Systems and Application Inc., Hampton, VA, 23666, USA

9 <sup>4</sup>NASA Goddard Space Flight Center, Global Modeling and Assimilation Office, Greenbelt, MD, 20771, USA

10 <sup>5</sup>Morgan State University, Goddard Earth Science Technology & Research (GESTAR) II, Baltimore, Maryland, USA

11 <sup>6</sup>Joint Center for Earth Systems Technology, Baltimore, MD, USA

12 <sup>7</sup>University of Maryland Baltimore County, Baltimore, MD, USA

13 *Correspondence:* Yuxuan Wang (ywang246@central.uh.edu)

14

15 **Abstract.** Coastal regions are susceptible to multiple complex dynamic and chemical mechanisms and emission sources that  
16 lead to frequently observed large tropospheric ozone variations. These large ozone variations occur on a meso-scale which  
17 have proven to be arduous to simulate using chemical transport models (CTMs). We present a clustering analysis of multi-  
18 dimensional measurements from ozone Light Detection And Ranging (LiDAR) in conjunction with both an offline GEOS-  
19 Chem CTM simulation and the online GEOS-Chem simulation GEOS-CF, to investigate the vertical and temporal variability  
20 of coastal ozone during three recent air quality campaigns: 2017 Ozone Water-Land Environmental Transition Study  
21 (OWLETS)-1, 2018 OWLETS-2, and 2018 Long Island Sound Tropospheric Ozone Study (LISTOS). We developed and  
22 tested a clustering method that resulted in 5 ozone profile curtain clusters. The established 5 clusters all varied significantly in  
23 ozone magnitude vertically and temporally which allowed us to characterize the coastal ozone behavior. The lidar clusters  
24 provided a simplified way to evaluate the two CTMs for their performance of diverse coastal ozone cases. An overall evaluation  
25 of the models reveals good agreement ( $R \approx 0.70$ ) in the low-level altitude range (0 to 2000 m), with a low and unsystematic  
26 bias for GEOS-Chem and high systemic positive bias for GEOS-CF. The mid-level (2000 – 4000 m) performances show a  
27 high systematic negative bias for GEOS-Chem and an overall low unsystematic bias for GEOS-CF and a generally weak  
28 agreement to the lidar observations ( $R = 0.12$  and  $0.22$ , respectively). In evaluating the cluster specific performances additional  
29 model insight is revealed as cluster-by-cluster model performance is more convoluted than the overall performances suggest.  
30 Utilizing the full vertical and diurnal ozone distribution information specific to lidar measurements, this work provides new  
31 insights on model's proficiency in complex coastal regions.

## 32 1. Introduction

33 Tropospheric ozone ( $O_3$ ) is an important secondary pollutant created by multiple reactions involving sunlight, nitrogen  
34 oxides ( $NO_x = NO + NO_2$ ), and volatile organic compounds (VOCs) which, in accumulation, can have damaging effects on  
35 human and plant health. In addition to its photochemical growth,  $O_3$  can easily be influenced by local and regional transport  
36 mechanisms. For coastal regions, surface  $O_3$  is highly variable in time and space due to its susceptibility to many factors such  
37 as local ship emissions, long range transport, and sea/bay breeze processes. Multiple studies have proven the strong influence  
38 that sea/bay breeze and wind flow patterns can have on the accumulation of coastal  $O_3$  and can often lead to poor air quality  
39 (e.g., Tucker et al., 2010; Martins et al., 2012; Stauffer et al., 2012; Li et al., 2020). Loughner et al. (2014) highlighted the  
40 importance for understanding the ability for bay breeze events to cause  $O_3$  differences not only spatially but vertically in coastal  
41 regions.

42 This variability is challenging for air quality models to capture as high-resolution measurements are necessary to fully  
43 understand and simulate this  $O_3$  behavior in coastal regions. For example, Dreessen et al. (2019) tested the U.S. Environmental  
44 Protection Agency (EPA) Community Multiscale Air Quality (CMAQ) model's ability, configured at 12 km, to simulate  $O_3$   
45 exceedances at Hart Miller Island in Maryland (HMI) revealing high bias and 'false alarms' due to several reasons such as  
46 emission transport over water and the coarse model resolution's inability to capture fine-scale meteorology and transport.  
47 Cases such as sea/bay breeze events, which directly contribute to high coastal  $O_3$  cases, are denoted by local meteorological  
48 mechanisms such as surface wind speed deceleration, wind direction convergence and recirculation (Banta et al., 2005). Air  
49 quality models with coarse horizontal and vertical resolutions are not able to capture such fine developments (Caicedo et al.,  
50 2019). Ring et al. (2018) also used CMAQ to estimate the impact of ship emissions on the air quality in eastern U.S. coastal  
51 regions indicating that an understanding of the vertical profiles of emissions was significant for improving air quality  
52 simulations. These are consistent and unanimous issues with air quality modeling in coastal regions. Since offshore sites within  
53 coastal regions are historically under sampled due to the difficulty of water-based measurements, this problem is still pertinent  
54 today.

55 Recently, three associated air quality campaigns set out to address this issue (<https://www-air.larc.nasa.gov/index.html>):  
56 2017 & 2018 NASA Ozone Water-Land Environmental Transition Study (OWLETS-1 & OWLETS-2) and Long Island Sound  
57 Tropospheric Ozone Study (LISTOS) (e.g., Sullivan et al., 2019). These three campaigns were each conducted in highly  
58 populated coastal regions along the Chesapeake Bay in Virginia and Maryland and Long Island Sound in the New  
59 England/Middle Atlantic region, respectively, that are vulnerable to  $O_3$  exceedances with the goal of filling the measurement  
60 gaps in these regions. During these campaigns, a suite of detailed airborne and ground measurements were taken during the  
61 course of highly polluted summer months (end of May through August) to capture the variability of pollutants, including  $O_3$   
62 and its precursor species, and the distinct meteorological processes specific to land-water regions that affect them.

63 The three campaigns strategically placed multi-dimensional tropospheric lidar measurements of  $O_3$  on and offshore in  
64 order to capture critical land-water gradients and to fill the deficit of measurements in these under monitored areas. These  
65 measurements were supported as part of NASA's Tropospheric Ozone Lidar Network (TOLNet). Continuous profile

66 measurements from O<sub>3</sub> lidars highlight important regional transport and temporal variations of O<sub>3</sub> in the lower and middle  
67 levels of the troposphere that are usually difficult to capture by most satellite-based remote-sensing instruments (Thompson et  
68 al., 2014). Lidar measurements are unique in their ability to capture high resolution full O<sub>3</sub> 2-D profile curtains over a period  
69 of time that indicate pollutant transport and can help in understanding O<sub>3</sub> behavior in coastal regions. In Gronoff et al. (2019),  
70 the co-located lidar at the Chesapeake Bay Tunnel Bridge (CBBT) during OWLETS-1 successfully captured a near-surface  
71 maritime ship plume emission event on August 01, 2017. An ensemble of other instruments (e.g., drones, Pandora spectrometer  
72 systems, etc.) launched near the shipping channel captured elevated NO<sub>2</sub> concentrations while the lidar instrument captured a  
73 depletion of O<sub>3</sub> simultaneously. The lidar was able to capture the unique low range altitude O<sub>3</sub> concentrations which elucidated  
74 the evolution of the trace-gas concentrations during this ship plume event.

75 Several studies have thoroughly evaluated the results from the air quality campaigns used in this study but were focused  
76 more on specific case studies (Dacic et al., 2019; Sullivan et al., 2019; Gronoff et al., 2019). Dacic et al. (2019) used lidar  
77 measurements of a high O<sub>3</sub> episode during OWLETS-1 to evaluate the ability of two NASA coupled chemistry-meteorology  
78 models (CCMMs), the GEOS Composition Forecast (“GEOS-CF”; Keller et al., 2021) and MERRA2-GMI (Strode et al.,  
79 2019), to simulate this high O<sub>3</sub> event. They found that the GEOS-CF model performed fairly in simulating O<sub>3</sub> in the lower  
80 level (between 400 to 2000 m ASL) and outperformed MERRA2-GMI based on surface observations at multiple monitoring  
81 sites and by a median difference of -6 to 8 % +/- 7 % at both lidar sites. In the case of this event, GEOS-CF was able to simulate  
82 the 2-D O<sub>3</sub> profile curtains at small scales. At the time of the Dacic et al. (2019) study, only processed observational data from  
83 OWLETS-1 was available.

84 For this study, we took advantage of 91 captured 2-D (vertical and diurnal) O<sub>3</sub> profile curtains from all three air quality  
85 campaigns (Sect. 2). To characterize the different behaviors of O<sub>3</sub> in coastal regions, we developed a novel clustering method  
86 based on the altitude and time dimensions of the lidar measurements that organized the profile curtains (Sect. 2). We used the  
87 developed clusters to evaluate the ability of both offline and online GEOS-Chem and GEOS-CF simulations to reproduce the  
88 coastal O<sub>3</sub> and wind characteristics highlighted by each cluster (Sect. 3).

89

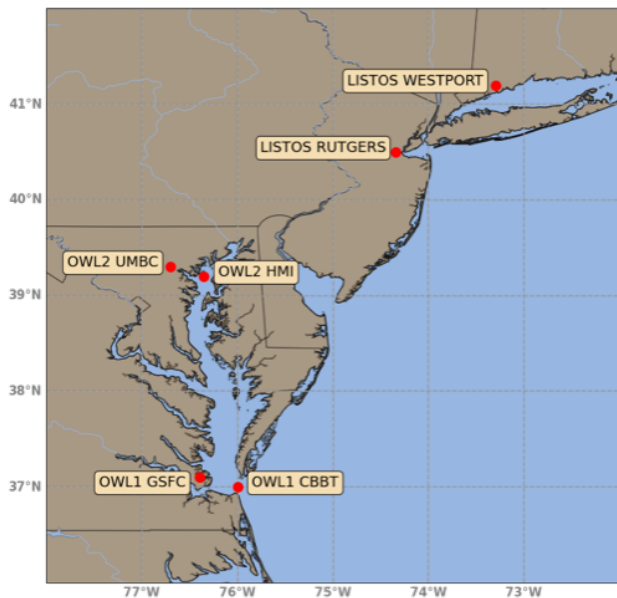
## 90 **2. Materials & Method**

### 91 **2.1. Air quality campaigns**

92 During the years 2017 and 2018, NASA in partnership with other U.S. national agencies and university research groups  
93 orchestrated three air quality campaign studies that focused on key land and water observations: OWLETS-1, OWLETS-2,  
94 and LISTOS. OWLETS-1 was conducted in 2017 from July 5 to August 3, while OWLETS-2 and LISTOS were conducted in  
95 2018 from June 6 to July 6 and July 12 to August 29, respectively. All campaigns took advantage of a multitude of ground,  
96 aircraft, and remote sensing measurements. For the sake of this study, we will focus on measurements from the two lidars from  
97 the TOLNet: NASA Langley Mobile Ozone Lidar (LMOL) (De Young et al. 2017; Farris et al. 2018; Gronoff et al, 2019,  
98 2021) and NASA Goddard Space Flight Center (GSFC) Tropospheric Ozone (TROPOZ) Differential Absorption Lidar (DIAL)  
99 (Sullivan et al. 2014, 2015a), which ran simultaneously at the marked positions in Figure 1. The TOLNet data from all three

100 campaigns are available on the NASA LaRC Airborne Science Data for Atmospheric Composition archive ([https://www-](https://www-air.larc.nasa.gov/missions.htm)  
101 [air.larc.nasa.gov/missions.htm](https://www-air.larc.nasa.gov/missions.htm); accessed – 20 January 2021).

102



103 **Figure 1.** An inset map of the Chesapeake Bay airshed in Maryland, Virginia, and Long Island Sound in New York with the  
104 six lidar monitoring locations used for OWLETS-1, OWLETS-2, and LISTOS highlighted and labeled.

105 The two lidars were placed strategically for each campaign (Figure 1), so that one lidar was closest to over-water  
106 measurements while the other was farther inland with the goal of examining how  $O_3$  transport and concentration is influenced  
107 by specific coastal mechanisms such as the land–water breezes. For OWLETS-1, the LMOL lidar was used at the CBBT  
108 [37.0366°N, 76.0767°W], depicting the real time over water  $O_3$  measurements while the GSFC TROPOZ lidar was stationed  
109 at NASA Langley Center [37.1024°N, 76.3929°W] further inland. Similarly, for OWLETS-2, the LMOL lidar was stationed  
110 for the over water measurements at Hart Miller Island [39.2449° N, 76.3583° W] and GSFC TROPOZ was stationed at the  
111 University of Maryland, Baltimore County (UMBC) [39.2557° N, 76.7111° W]. Finally, for LISTOS, LMOL was at the  
112 Westport site [41.1415° N, 73.3579° W] and TROPOZ at Rutgers [40.2823° N, 74.2525° W]. For the sake of this study the  
113 unique benefits due to the different placements (onshore versus offshore) of the co-located lidars are not specifically evaluated.  
114 Instead, the study focuses on the benefits of detailed and multi-dimensionality of both lidar instrument data in general.

115 Routine lidar measurements were taken for the duration of the campaigns providing 91 multi-dimensional  $O_3$  profile  
116 curtains. Both lidars retrieve data at a 5-min temporal resolution and use a common processing scheme to produce a final  $O_3$   
117 product which was used for this study. In this study, the individual profile curtains refer to the “full day”, vertical and diurnal

118 lidar measurements. In this study, 91 individual 2-D profile curtains were used from both lidars from the three campaigns: 26  
119 profile curtains from OWLETS-1, 28 profile curtains from OWLETS-2, and 37 profile curtains from LISTOS.

120 To evaluate meteorological impacts on the lidar O<sub>3</sub> clusters and distinguish certain model discrepancies we used various  
121 temperature and wind measurements. Hourly observed temperature, wind speed and wind direction, and O<sub>3</sub> from surface  
122 monitors pertaining to the study area were obtained from the Air Quality System (AQS) (data can be accessed at  
123 <https://aqs.epa.gov/aqsweb/airdata/>). Along with the O<sub>3</sub> lidar instruments, we utilized high resolution vertical and horizontal  
124 wind speed and direction data monitored by Doppler wind lidar Leosphere WINDCUBE 200s instruments deployed at HMI  
125 during OWLETS-2 during LISTOS (e.g., Couillard et al., 2021; Coggon et al., 2021; Wu et al., 2021).

126

## 127 **2.2. Clustering lidar data**

### 128 **2.2.1 Description of the ozone lidar measurements**

129 The lidar instrument is unique in that it provides high dimensional profile measurements of O<sub>3</sub>, as opposed to one  
130 dimensional surface measurements from air quality monitoring sites. The two TOLNet lidars used during the campaigns have  
131 been evaluated for their accuracy during previous air quality campaigns (DISCOVER-AQ; [https://www-  
132 air.larc.nasa.gov/missions/discover-aq](https://www-air.larc.nasa.gov/missions/discover-aq) and FRAPPÉ; <https://www2.acom.ucar.edu/frappe>) and have also been compared  
133 against each other (e.g., Sullivan et al., 2015; Wang et al., 2017). The two lidars have different transmitter and retrieval  
134 components but produce O<sub>3</sub> profiles within 10 % of each other as well as compared to ozonesondes (Sullivan et al., 2015). In  
135 comparison with other in situ instrument measurements, the TOLNet lidars were found to have an accuracy better than ±15 %  
136 for capturing high temporal tropospheric O<sub>3</sub> vertically proving their capability of capturing high temporal tropospheric O<sub>3</sub>  
137 variability (Wang et al., 2017; Leblanc et al., 2018).

138 To characterize coastal O<sub>3</sub> during the summer months, we use a multitude of lidar profile curtains obtained during the  
139 OWLETS-1, 2, and LISTOS campaigns. The two lidars used in the campaigns produced profile curtains of O<sub>3</sub> from 0 – 6000  
140 m above ground level (AGL) with some days beginning as early as 06:00 local time (EDT) and ending measurements as late  
141 as the last hour of the day. One of the challenges is that the multiple lidar datasets are not always uniform; although most of  
142 the profile curtains began at or around 08:00 EDT, the lidar measurements commence and conclude at different times. At the  
143 time of these campaigns, the lidar data retrieval was constrained by the availability of personnel as well as the availability of  
144 electricity in remote areas (at time of writing, the lidar instrument systems have been updated and are now more fully  
145 automatized for use during succeeding campaigns removing such constraints). Due to this constraint, the 91 lidar curtains  
146 range from as short as a 6-hour window to a full 24-hour window. Similarly, the profile curtains do not have an exact uniform  
147 altitude range either. In the processing of the lidar data, some measurements may be filtered out and removed due to issues,  
148 such as clouds, which can influence and degrade the retrieval leaving some blocks of empty data within the vertical altitude  
149 dimension. When the cloud conditions are perfect, the limiting factor for the altitude is the solar background: the UV from the  
150 sun is a source of noise that prevents the detection of the low level of backscattered photons. For LMOL, this means that the  
151 maximum altitude is about 10 km AGL at night (Gronoff et al., 2021) and lowered to about 4 km AGL at solar noon (worse

152 conditions possible for the summer in the continental U.S. resulting in below 4 km AGL). This results in a general scarcity of  
153 O<sub>3</sub> measurements above 4000 m AGL for most of the vertical profile curtains. Lidars still have limitations that prove to be a  
154 complication e.g., noise signal and manual operations. At the time of writing, the operative limitation has been addressed and  
155 the lidars are now more fully automated which removes some of the difficulty.

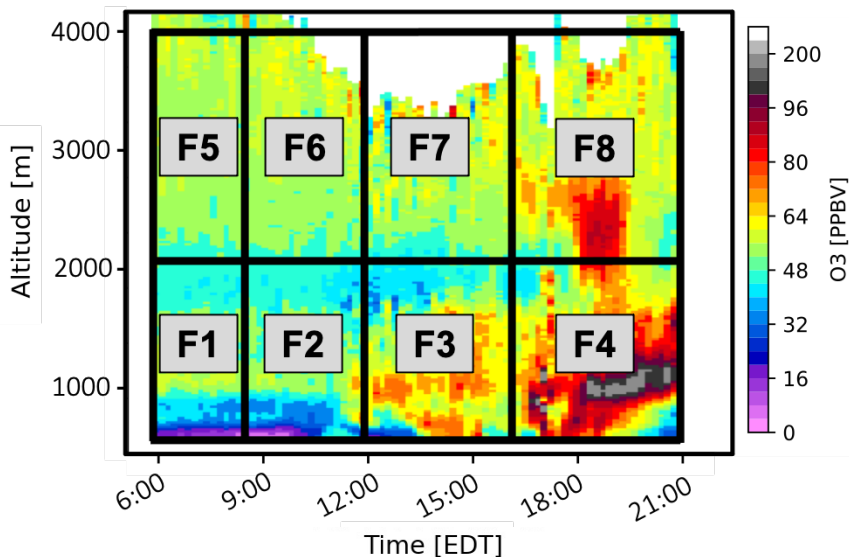
156

### 157 **2.2.2 Clustering approach and application**

158 To facilitate the comparison of the 2-D O<sub>3</sub> profile curtains and the air quality model simulations we used a cluster analysis  
159 that categorized the behavior of the tropospheric O<sub>3</sub> captured in the profile curtains. Clustering methods are commonly used  
160 in air quality and atmospheric studies to group and characterize large datasets (Darby, 2005; Alonso et al., 2006; Christiansen,  
161 2007; Davis et al., 2010; Stauffer et al., 2018). In our previous work, we have successfully used clustering methods to  
162 automatically characterize diurnal patterns of surface winds and surface O<sub>3</sub> in the Houston-Galveston-Brazoria area that proved  
163 to perform better than a rudimentary quantile method to reveal the dependence of surface O<sub>3</sub> variability on local and synoptic  
164 circulation patterns on the Gulf Coast (Bernier et al., 2019; Li et al., 2020)

165 In evaluating the structure of the lidar measurements and working within measurement limitations (described in Sect.  
166 2.2.1) from the three air quality campaigns, we developed a method to cluster multi-dimensional O<sub>3</sub> profile curtains using K-  
167 Means clustering algorithm. Input features (seed values) were rationally established to best represent the behavior of O<sub>3</sub>  
168 temporally and vertically without including an excessive amount of input features, which can weaken the results of clustering  
169 (discussed in detail in Sect. S1). With the goal of evaluating lower level tropospheric O<sub>3</sub> and based on description of the  
170 structure and constraints of the lidar measurements, the features were tailored to the altitude range 0 – 4000 m AGL and time  
171 range of 06:00 EDT – 21:00 EDT.

172 Figure 2 illustrates the 8 features that represent slabs of altitude and time used in the cluster analysis. For each O<sub>3</sub> profile  
173 curtain (total of 91), we calculated the average O<sub>3</sub> from the following time and altitude range: Features 1 – 4 altitudes range  
174 from 0 – 2000 m; Features 5 – 8 altitudes range from 2000 – 4000 m. The two altitude ranges were determined to best represent  
175 different O<sub>3</sub> transport events although they do not explicitly represent these layers. For Features 1 – 4, O<sub>3</sub> would most likely  
176 primarily be affected by local production and pollution transport while for Features 5 – 8, O<sub>3</sub> would more likely be associated  
177 with long range transport (e.g. interstate). As planetary boundary layer growth (PBL) in coastal regions do not usually reach  
178 altitudes greater than 2000 m, mixing between the boundary layer and free troposphere would presumably take place within  
179 the low-level altitude bin. Additional attention to the PBL in the selecting of low versus mid-level features for the clustering  
180 will be investigated in future work. For clarity, we will use the terms low-level and mid-level features to address the two  
181 altitude subsets e.g., Features 1 – 4 and 5 – 8, respectively. Feature 1 and 5 time range from 06:00 – 08:00 EDT; Feature 2 and  
182 6 from 08:00 – 12:00 EDT; Feature 3 and 7 from 12:00 – 16:00 EDT; and Feature 4 and 8 from 16:00 – 21:00 EDT. The four  
183 subset time ranges were indicated to best represent features that characterize the common diurnal behavior of O<sub>3</sub>.



184

185 **Figure 2.** Clustering method developed for clustering vertical O<sub>3</sub> profiles taken from lidar measurements. The color coding  
 186 shows a typical day of lidar measurements of O<sub>3</sub> profiles on August 6, 2018, from the LMOL at Westport, CT during the  
 187 LISTOS Campaign. F1 – F8 indicate the time and altitude range of the eight features used for the clustering algorithm.

188

189 The features were evaluated for cluster tendency, essentially to confirm our dataset contained meaningful clusters  
 190 (discussed in detail in Sect. S2). One statistical approach was used to test the dataset called Hopkins statistic which measures  
 191 whether there is uniform distribution (spatial randomness) within the dataset (Lawson and Jurs, 1990). The results calculated  
 192 using the Hopkins statistic concluded a value higher than 0.75 (actual = 0.77) which by this standard indicates a clustering  
 193 tendency at the 90 % confidence level. Evaluating different feature options did not lead to better statistical results than with  
 194 the final chosen features. To visualize the cluster tendency of our dataset, we applied the algorithm of the visual assessment of  
 195 cluster tendency (VAT) approach (Bezdek and Hathaway, 2002) which uses the Euclidean distance measure to compute the  
 196 dissimilarity matrix in the dataset and creates an ordered dissimilarity matrix image. Figure S1 shows the VAT approach results  
 197 which indicates high similarity (red) and low similarity (blue) and confirms a cluster structure (not random) within our dataset.

198

199 Since the choice of clustering algorithm is subjective, we chose K-means clustering for its simplicity and widespread use.  
 200 To use the K-Means clustering algorithm, the optimal number of clusters based on your dataset must be chosen beforehand.  
 201 For this study, the package Nbclust (Charrad et al., 2014) in R was used, which applies 30 indices for determining the optimal  
 202 number of clusters. Using this package, as well as testing the quality of the clustering results using the silhouette method  
 203 (Kaufman & Rousseeuw, 1990), we selected six clusters as the optimal number of clusters. Since the K-Means clustering  
 204 algorithm is based on the Euclidean distance to each centroid, the input data was normalized (to a mean of zero and standard  
 deviation of one) to ensure each feature is given the same importance in the clustering (Aksoy & Haralick, 2001; Larose,

205 2005). The resulting six clusters (described fully in Sect. 3.1) represent clusters of regularly observed lidar O<sub>3</sub> curtains for the  
206 regions of our study during the campaign periods.

207

### 208 **2.2.3 Missing data**

209 Although the input features were tailored based on the structure of the lidar measurements, the remaining data still had  
210 missing data points. In performing a quick evaluation on the 8 input features (Figure S5), we found that Features 1, 4, 5, and  
211 8 had the most missing data while Features 2, 3, 6, and 7 had few or zero cases of missing data. This means that the earlier  
212 morning measurements (06:00 – 12:00 EDT) and the later evening measurements (16:00 – 21:00 EDT) had the most cases of  
213 missing data points. This is plausible as the campaign teams were best able to retrieve clear measurement during  
214 midday/evening hours (12:00 – 16:00 EDT). As a result, 51 out of 91 O<sub>3</sub> profile curtains had at least one missing data point  
215 (feature) throughout the individual profile curtain.

216 A common practice for dealing with missing data is complete case analysis (CCA), in which observations with missing  
217 values are completely ignored, leaving only the complete data to cluster. CCA can be inefficient as it introduces selection bias  
218 since the sample data no longer retains the state of the original full dataset (Donders et al., 2006; Little & Rubin, 2014). When  
219 we applied CCA, there were only 40 O<sub>3</sub> profile curtains of complete data, removing over half of the study profiles. Instead,  
220 we used a more comprehensive solution – imputation - that yields unbiased results (Donders et al., 2006). For this study we  
221 used the single imputation (SI) technique *knnImputation* in R (Torgo, 2010), which uses the k-nearest neighbors and searches  
222 for the most similar cases and uses the weighted average of the values of those neighbors to fill the missing data. Essentially,  
223 this method selects the days that have the most similar profile curtain to any profile which has missing data points and uses  
224 those real data points to calculate a weighted mean that will fill in the missing data. We acknowledge using an imputation  
225 method on the dataset will possibly introduce a bias which is difficult to quantify, but this allows the use of the full 91 profile  
226 curtains of O<sub>3</sub> data. The silhouette method was used to test the quality of the newly imputed dataset and proved to be no worse,  
227 nor better, than the CCA (*real data*) results. Therefore, the dataset was first imputed using SI to create a complete dataset and  
228 then the clustering method described in the sect. before (2.2.2) was applied to the complete imputed dataset.

229

### 230 **2.3. Model simulations**

231 The offline GEOS-Chem chemical-transport model (CTM) was utilized to simulate the spatial and temporal variability  
232 of coastal O<sub>3</sub> in the Chesapeake Bay and Long Island Sound during the time of the campaigns. The GEOS-Chem model is a  
233 global 3-D CTM driven by assimilated meteorological data from the NASA Global Modeling and Assimilation Office  
234 (GMAO). Our simulations were driven by reanalysis data from Modern-Era Retrospective analysis for Research and  
235 Applications, Version 2 (MERRA-2; Gelaro et al., 2017). We ran a nested GEOS-Chem (v12-09) simulation at 0.5° x 0.625°  
236 horizontal resolution over the eastern portion of North America and adjacent ocean (90 – 60°W, 20 – 50°N), using lateral  
237 boundary conditions updated every three hours from a global simulation with 2° x 2.5° horizontal resolution. The nested  
238 GEOS-Chem simulation was run with 72 vertical levels from 1013 to 0.01 hPa. Since the study focuses on the altitude range



239 0 – 4000 m, the first 20 vertical levels from GEOS-Chem were used with 14 levels within the boundary layer ( $\leq 2000$  m). The  
240 nested simulation was conducted for the study periods June – September 2017 and April – August 2018. We used the standard  
241 “out-of-the-box” unmodified default settings from the tropospheric chemistry chemical mechanism (tropchem) with global  
242 anthropogenic emissions from the Community Emissions Data System (CEDS) inventory (McDuffie et al, 2020) and U.S.  
243 Environmental Protection Agency (EPA) National Emissions Inventory (NEI) 2011 for monthly mean North American  
244 regional emissions (EPA NEI, 2015).

245 We also used results from NASA’s near real-time forecasting system, GEOS-CF, an online GEOS-Chem simulation (v12-  
246 0-1) from GMAO ([https://gmao.gsfc.nasa.gov/-weather\\_prediction/GEOS-CF/](https://gmao.gsfc.nasa.gov/-weather_prediction/GEOS-CF/)) with GEOS coupled to the GEOS-Chem  
247 tropospheric-stratospheric unified chemistry extension (UCX) and run at a high spatial resolution of  $0.25^\circ$ , roughly 25 km  
248 (Keller et al., 2021, Knowland et al., 2021). The vertical resolution for GEOS-CF is interpolated onto 72 vertical levels from  
249 1000 to 10 hPa. Since the study focuses on the altitude range 0 – 4000 m, the first 21 vertical levels from GEOS-CF were used  
250 with 14 levels within the boundary layer ( $\leq 2000$  m). Prior to the launch of the 12z five-day forecast, GEOS-CF produces daily  
251 global, 3-D atmospheric composition distributions using the GEOS meteorological replay technique (Orbe et al., 2017), and  
252 this study makes use of these historical estimates, made available to the public for the period since January 2018. Therefore,  
253 the GEOS-CF results shown in this study only include the dates from OWLETS-2 and LISTOS campaigns, since they both  
254 occurred in 2018.

255 While both model simulations use similar versions of GEOS-Chem chemistry, there are noteworthy differences to keep  
256 in mind during the analysis of the clustering. The main differences between the two models are (1) GEOS-Chem is an offline  
257 CTM using archived meteorology, while GEOS-CF simulates atmospheric composition simultaneously with meteorology  
258 (online); (2) the spatial resolution of the GEOS-CF model ( $0.25^\circ$ ) is higher than GEOS-Chem ( $0.5^\circ \times 0.625^\circ$ ); and (3) the  
259 GEOS-CF model runs with Harmonized Gridded Air Pollution (HTAP; v2.2; base year 2010) anthropogenic emissions from  
260 the Emission Database for Global Atmospheric Research (EDGAR), while GEOS-Chem was run with CEDS anthropogenic  
261 emissions (base year 2014). These imperative differences can lead to disparities in the following results.

262

### 263 **3. Results & Discussion**

#### 264 **3.1 Overview of the 2-D O<sub>3</sub> curtain clusters**

265 The clustering results reveal distinctive characterized O<sub>3</sub> behavior during the three campaigns in which O<sub>3</sub> concentrations  
266 vary across the clusters. As previously mentioned in Sect. 2.2.3, the clustering analysis initially identified six cluster groups  
267 from the O<sub>3</sub> profile curtains. Only one date was assigned to Cluster 6 (16 June 2018): the lidar profile curtain on this day  
268 (Figure S6) shows a large fraction of data missing, and the available data have relatively high O<sub>3</sub> throughout the lowest 3 km,  
269 which is different from other clusters. Therefore, we consider Cluster 6 to be an outlier and will not include it in the subsequent  
270 analysis.

271 Various O<sub>3</sub> and surface meteorological parameter cluster statistics for the remaining five clusters are summarized in Table  
272 1. With only 5 of the 2-D profile curtains assigned, Cluster 5 depicts the least common O<sub>3</sub> behavior during the campaigns. On

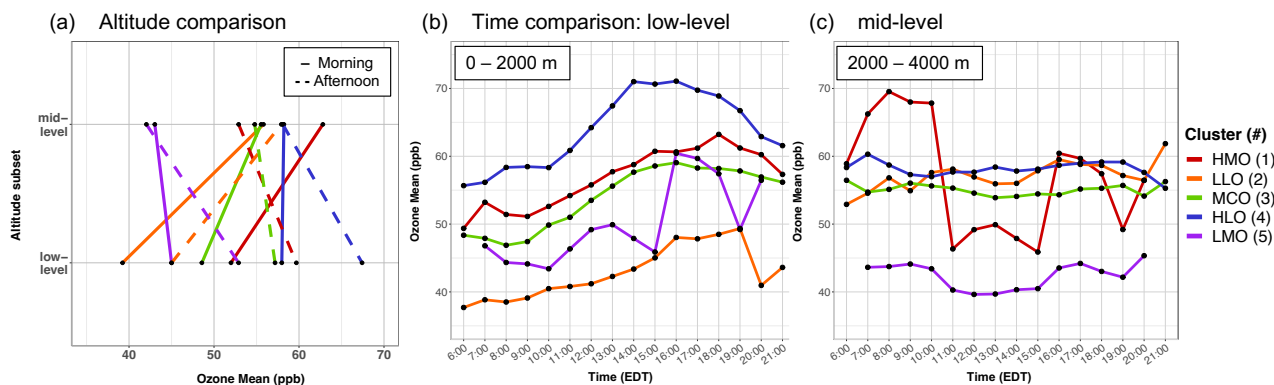
273 the other hand, Cluster 3 is the most common O<sub>3</sub> behavior during the campaigns with 28 profile curtains assigned to this  
 274 cluster. Following Cluster 3, Cluster 1 is the next most common cluster with 25 profile curtains. Cluster 2 and Cluster 4 fall in  
 275 the middle with 14 and 18 profile curtains assigned to the cluster numbers, respectively.

Cluster #	a) No. of vertical profiles	b) O <sub>3</sub> Max (ppb)	c) O <sub>3</sub> Min (ppb)	d) T avg. (min; max) (°F)	e) WS avg. (min; max) (m s <sup>-1</sup> )
1	25	86.5	42.2	74.1 (67.8; 86.4)	1.5 (0.5; 2.8)
2	14	72.8	28.9	71.6 (64.0; 83.9)	1.6 (0.6; 2.9)
3	28	86.6	34.2	77.2 (67.0; 87.6)	1.3 (0.5; 2.4)
4	18	97.8	44.1	78.4 (68.0; 90.4)	1.2 (0.4; 2.3)
5	5	67.7	29.1	74.5 (66.8; 74.5)	1.2 (0.3; 3.4)

285 **Table 1.** Lidar vertical O<sub>3</sub> profile cluster statistics: a) total number of vertical profiles; b) O<sub>3</sub> maximum; c) O<sub>3</sub> minimum O<sub>3</sub>;  
 286 AQS monitoring station cluster mean d) surface temperature and e) wind speed; minimum and maximums in parenthesis. The  
 287 statistics and averages were derived from the total number of profile curtains assigned to each cluster.

288  
 289 The five clusters were distinguished by the varying O<sub>3</sub> concentrations between the low-level and mid-level as well as  
 290 diurnal variations (Figure 3). Figure 3a quantifies the between-cluster differences. We separate the data by the two altitude  
 291 subsets (low and mid-level) and by two time subsets (morning = 6:00 – 12:00 and afternoon = 12:00 – 21:00) for lucidity as  
 292 the majority of the cluster differences are contrasted between these subsets. In the low-level, all five clusters exhibit the  
 293 common O<sub>3</sub> diurnal pattern where surface O<sub>3</sub> is titrated overnight and reaches a minimum but then is quickly exacerbated with  
 294 the increase of sunlight throughout the day and typically peaks after midday (Figure 3b). The extent of this common diurnal  
 295 pattern varies by cluster.

296



297

298 **Figure 3.** Lidar O<sub>3</sub> cluster average comparisons (five clusters depicted in colors). a) Altitude comparison of mean O<sub>3</sub> averaged  
299 over time: morning hours from 6:00 – 12:00 (solid line) and afternoon hours from 12:00 – 21:00 (dashed lines). Time  
300 comparison of mean hourly O<sub>3</sub> split between the b) low-level and c) mid-level.

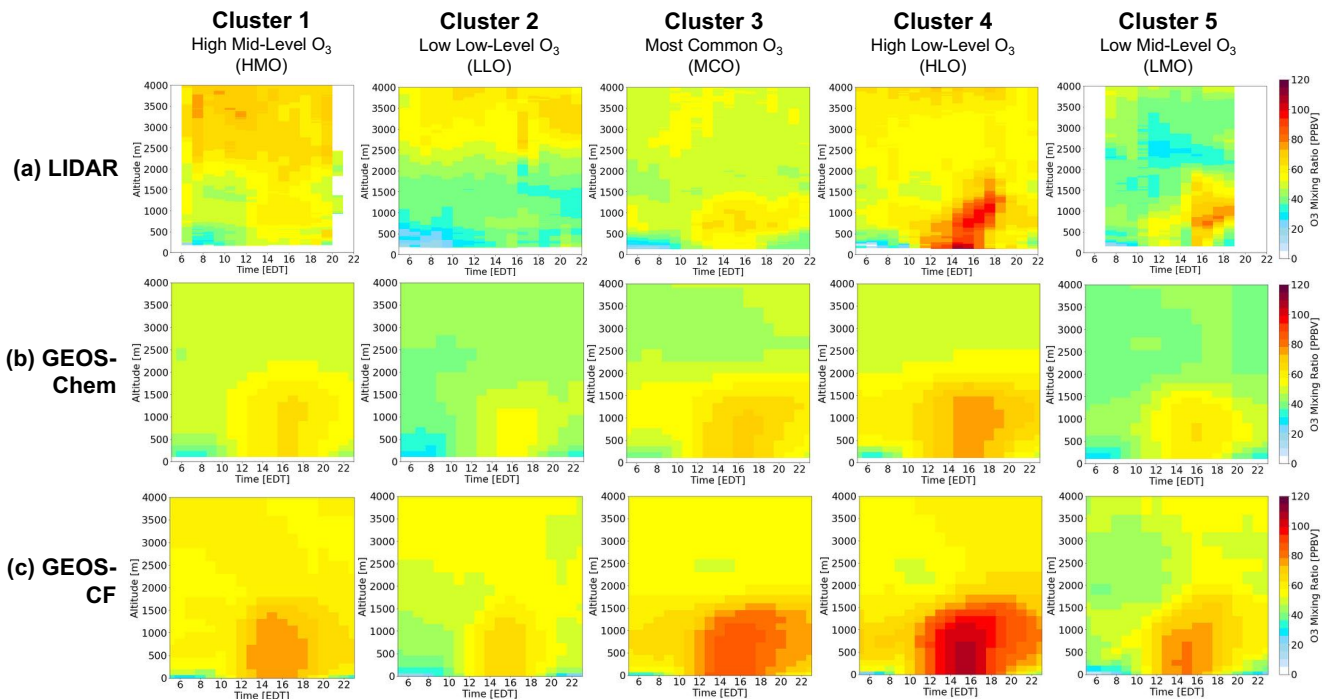
301

302 Cluster 1 in the low-level has the second highest morning and afternoon O<sub>3</sub> average (52 and 59 ppb) and in the mid-level  
303 the highest morning O<sub>3</sub> average (64 ppb) (Figure 3a). Cluster 1 also exhibits the most unique pattern of mid-level O<sub>3</sub> (Figure  
304 3c), with the highest concentrations found in the early morning and an uncharacteristic plunge to lower O<sub>3</sub> concentrations from  
305 11:00 – 15:00 EDT. This is contrary to the other clusters which do not show much O<sub>3</sub> variation temporally in the mid-level.  
306 The majority of the individual profile curtains assigned to Cluster 1 show concentrated early morning residual layers in the  
307 mid-level that diffuse after the morning, which is distinctive to the other clusters. In the low-level, Cluster 2 has the lowest  
308 morning and afternoon O<sub>3</sub> average among the clusters (39 and 45 ppb) with moderate mid-level O<sub>3</sub> concentrations. Cluster 3  
309 has the most uniform vertical O<sub>3</sub> extent between the low and mid-level (Figure 3a), in contrast to the other clusters that differ  
310 greatly in O<sub>3</sub> concentrations between the two altitude subsets. Cluster 4 has the highest morning and afternoon O<sub>3</sub> averages (59  
311 and 68 ppb) in the low-level, reaching > 70 ppb temporally (Figure 3b). Finally, Cluster 5 has, considerably, the lowest morning  
312 and afternoon O<sub>3</sub> averages (42 and 43 ppb) in the mid-level, almost 10 ppb lower than the other clusters. Cluster 5 does not  
313 have a smooth-evolving O<sub>3</sub> diurnal pattern in the lower level (Figure 3b), which can be attributed to the averaging of only five  
314 different profile curtains that were assigned to this cluster (Table 1).

315 Figure 4a illustrates the mean lidar O<sub>3</sub> 2-D profile curtains for each of the clusters. For Cluster 1, 3, 4, and 5, higher O<sub>3</sub>  
316 concentrations in the low-level are captured during afternoon/evening time (12:00 – 21:00 EDT), with the highest low-level  
317 O<sub>3</sub> in Cluster 4 (> 70 ppb). This behavior follows the common diurnal pattern of O<sub>3</sub>, that was distinguishable in Figure 3b. This  
318 common O<sub>3</sub> growth reaches vertically to approximately 1500 m for each of the clusters but is generally contained below 2000  
319 m. Differing from the low-level O<sub>3</sub> behavior, mid-level O<sub>3</sub> is generally less variable in magnitude throughout the entire profile  
320 curtain (except for Cluster 1; see Figure 3a). The highest O<sub>3</sub> concentrations for the mid-level are exhibited in Cluster 1, 2, 3,  
321 and 4, with the highest mid-level O<sub>3</sub> in Cluster 1 during the early morning hours ( $\geq 70$  ppb).

322 Following the descriptions above, each cluster is given a nomenclature according to their unique characteristics. Cluster  
323 1 is termed as the highest mid-level O<sub>3</sub> (HMO) cluster; Cluster 2 as the lowest low-level O<sub>3</sub> (LLO) cluster; Cluster 3 is the  
324 most common O<sub>3</sub> (MCO) cluster; Cluster 4 is the highest low-level O<sub>3</sub> (HLO); Cluster 5 is the least common and lowest mid-  
325 level O<sub>3</sub> (LMO) cluster. The O<sub>3</sub> variability represented and justified above is what led to the successful clustering of the lidar  
326 O<sub>3</sub> 2-D profile curtains.

327



328

329 **Figure 4.** Cluster-mean O<sub>3</sub> vertical profile results by cluster assignment (1- 5) and arranged: a) LIDAR; b) GEOS-Chem  
 330 simulation; and c) GEOS-CF simulation.

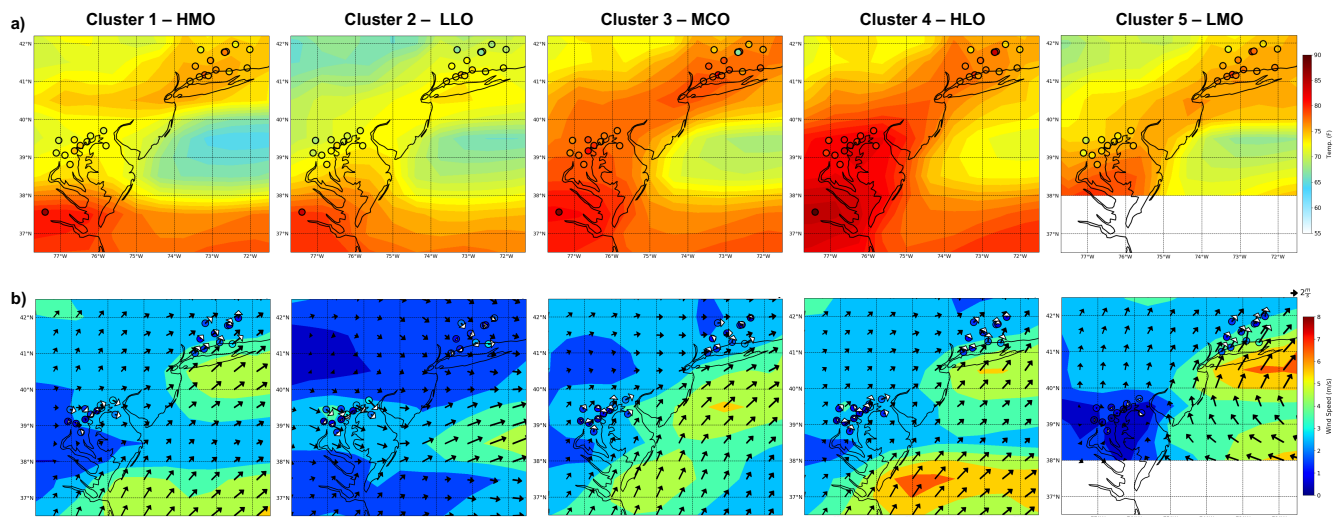
331

332 The clustering analysis results provided a characterization of O<sub>3</sub> behavior that transpired during these three campaigns.  
 333 Figure 3b and 3c indicate each cluster represents a different O<sub>3</sub> evolution pattern, likely related to different photochemical or  
 334 transport regimes. This kind of evaluation is useful in that it combines O<sub>3</sub> information from both temporal and vertical  
 335 dimensions. For example, the HLO cluster reveals the specific case in which higher O<sub>3</sub> is captured early in the temporal profile  
 336 in the low-level and translates to the higher O<sub>3</sub> captured in the low-level as well. The profile curtains show higher background  
 337 O<sub>3</sub>, indicating these cases did not have “clean air” to begin with which can allow a greater accumulation in the low-level in the  
 338 afternoon. This is an example of how this type of clustering analysis, if applied, could demonstrate background O<sub>3</sub> in the  
 339 similar case studies. In another example, several profile curtains assigned to the HMO cluster indicate concentrated residual  
 340 layers in the mid-level and possible entrainment to the surface as the day progressed. To prove this feature, vertical velocity  
 341 and vertical velocity variance data would be needed but the knowledge that a clustering approach is able to pinpoint these  
 342 features that could only be discernible through lidar measurements proves to be useful. The clustering results was valuable in  
 343 recognizing a significant large pollution related cluster (HLO), a total of 18 out of the 91 curtain profiles which correspond  
 344 with the highest daily surface maxima measured at these sites (= 97.77 ppb) (Table 1). This cluster, on average, exhibited a  
 345 daily surface maxima up to 10 ppb greater than any of the other clusters. Discerning these higher O<sub>3</sub> cases is imperative for  
 346 mitigating severe air pollution.

347

348 **3.2. Cluster surface analysis**

349 To support the lidar clustering results, daily averaged meteorological surface observations from AQS stations pertaining  
 350 to the campaign period and GEOS-Chem surface model output were evaluated in regard to the five clusters. Figure 5 shows  
 351 the cluster mean surface temperature from AQS stations and GEOS-Chem model as well as the simulated wind speed and  
 352 direction. The average surface temperature from each station is represented as the circular markers while the simulated  
 353 temperatures are represented as the spatial contour and the simulated wind speed ( $\text{m s}^{-1}$ ) and direction as arrows. Cluster  
 354 average, minimum, and maximum AQS surface temperature and wind speed can be found in Table 1d, e.  
 355



356

357 **Figure 5.** Cluster averaged meteorological surface AQS station observations and GEOS-Chem model results. a) Surface  
 358 temperature observations represented as the circular markers and simulated surface temperatures represented as the spatial  
 359 contour (top-panel). b) Surface wind speed and direction observations represented as the circular markers and white arrows  
 360 and simulated wind speed and direction represented as spatial contour and black arrows (bottom-panel).

361

362 In general, the surface meteorological conditions agree with our knowledge of transport and  $\text{O}_3$  production that would  
 363 lead to each of the five clustered lidar  $\text{O}_3$  profile curtains. It is evident that the clusters with the highest surface  $\text{O}_3$  (HMO,  
 364 MCO, and HLO) all share a predominant offshore, westerly wind. Furthermore, MCO and HLO presented higher overall  
 365 observed and simulated surface temperatures compared to the other clusters (Figure 5a). Observed and simulated wind speeds  
 366 reveal slightly lower average wind speeds and primarily continental wind flow for both clusters as well (Figure 5b). These  
 367 meteorological conditions are conducive to a higher production of surface  $\text{O}_3$  concentrations which validates the higher  $\text{O}_3$   
 368 found in the low-level results (Figure 3b, 4a).

369 Conversely, the lowest surface temperatures are found in LLO. Lower surface temperatures are also indicative of low  
370 vertical mixing due to less generation of convection. Relatively calm wind speeds and lower temperatures indicate other  
371 possible meteorological factors such as high cloud cover that could have contributed to the lower O<sub>3</sub> concentrations in LLO.  
372 Although surface O<sub>3</sub> concentrations in LMO reach higher levels later in the day, first at 13:00 EDT and then again at 16:00  
373 EDT, the rest of the temporal profile stays below moderate levels. Average temperatures for LMO are moderately high but, in  
374 contrast, the average wind speed is higher (specifically over the Long Island Sound) and unique to the other clusters, wind  
375 direction is predominantly onshore (Easterly – Southerly). This prevalent onshore flow indicates a transport of cleaner marine  
376 air which corroborates the lower surface O<sub>3</sub> levels. LMO did not have any profile curtains assigned from OWLETS-1 which  
377 is why data for the lower Chesapeake Bay area is not shown in Figure 5.

378 There was only one occurrence during the dates in which the lidar instruments were operating in which there was a  
379 recorded maximum daily 8-hour average (MDA8) O<sub>3</sub> exceedance (> 70 ppbv). This exceedance date is 25 May 2018 in which  
380 3 AQS sites in the LISTOS region measured MDA8 O<sub>3</sub> of 73, 72, and 72 ppbv. This curtain profile was assigned to the HMO  
381 cluster (Cluster 1), the cluster with high O<sub>3</sub> in the mid-level and moderate O<sub>3</sub> in the low-level and near the surface.

382

### 383 **3.3. Evaluating the GEOS-Chem and GEOS-CF model**

384 In this sect. the model results from GEOS-Chem and GEOS-CF will be compared to the lidar data using the five lidar O<sub>3</sub>  
385 profile clusters discussed in Sect. 3.1. Both model results were sampled in an equal manner, in which we extracted the same  
386 cluster date assignments from the lidar clusters and created mean vertical profiles based on the model results. This allowed us  
387 to evaluate the model performance based on the five characterized O<sub>3</sub> lidar clusters. As mentioned previously, the GEOS-CF  
388 simulation data is not available for 2017. Thus, the results shown subsequently will only include GEOS-CF results from 2018  
389 (only dates from the OWLETS-2 and LISTOS campaigns). The GEOS-Chem simulation results include both years thus all  
390 three campaign duration periods.

391

#### 392 **3.3.1 Overall model performance**

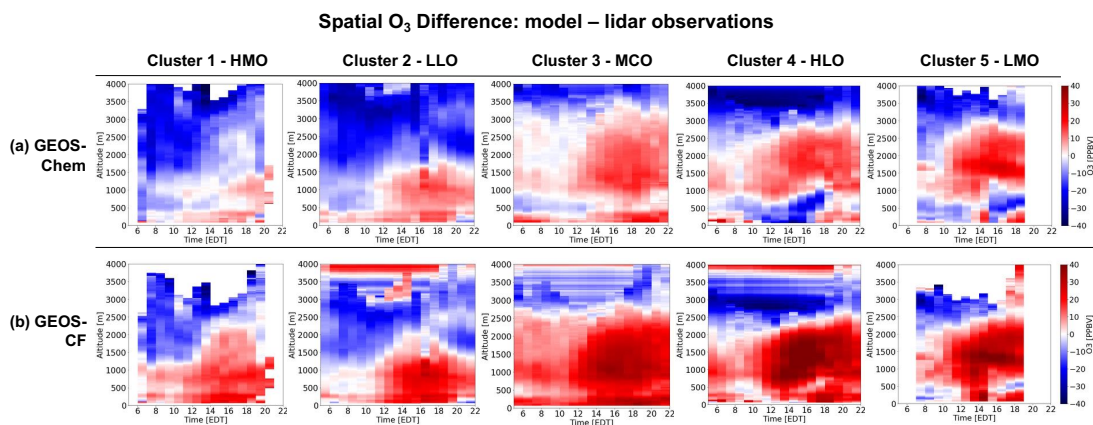
393 Figure 4b and 4c depict the simulated cluster-mean O<sub>3</sub> profile curtains from GEOS-Chem and GEOS-CF, mirroring the  
394 mean lidar profile curtains in Figure 4a. For all clusters in the low-level, both models simulate a consistent accumulation of  
395 O<sub>3</sub> near the surface after 12:00 EDT, mirroring the O<sub>3</sub> common diurnal pattern depicted in mean lidar profile curtains in Figure  
396 4a. However, the extent the models simulate is often higher in magnitude than the observations, specifically GEOS-CF  
397 consistently predicting the accumulation at a higher magnitude than GEOS-Chem. In the mid-level, both models simulate  
398 much less O<sub>3</sub> variability than what is captured in the lidar observations. Figure 4b and 4c clearly show how the models struggle  
399 to reproduce any mid-level O<sub>3</sub> pattern or variability that is relayed in the lidar observations. This is in contrast to the low-level  
400 where the models are able to reproduce the common diurnal pattern of O<sub>3</sub>. With the lidar data providing a full temporal and  
401 vertical profile curtain of O<sub>3</sub> behavior and development, we are able to indicate areas where the models struggle such as in this  
402 case in the mid-level.

403 We first evaluate overall correlation and biases between the model and lidar data. The overall correlation between both  
 404 models and the lidar data, disregarding the specific clusters, based on the two altitude subsets as the performances differ  
 405 between low-level and mid-level for both GEOS-Chem (Figure S7a) and GEOS-CF (Figure S7b). The mean normalized biases  
 406 for the five clusters displayed in Table S1 (in Supplementary Material) were calculated from the total vertical and diurnal  
 407 averages separated by low-level and mid-level. For both models, overall low-level O<sub>3</sub> correlation rounds to 0.70, signifying a  
 408 strong relationship between the model simulations and the lidar observations (Figure S7 - top panel). This indicates that both  
 409 models can simulate the development and pattern of O<sub>3</sub> well in the low-level. Overall, GEOS-Chem performs well in  
 410 simulating low-level O<sub>3</sub> with a lower non-systematic normalized bias ranging from -0.10 to +0.13 for the five clusters. Thus,  
 411 based on the lower bias, GEOS-Chem fairs well simulating the magnitude of low-level O<sub>3</sub> as well. For all clusters, GEOS-CF  
 412 overestimates the average magnitude of low-level O<sub>3</sub> with a systematic high positive normalized bias ranging from +0.30 to  
 413 +0.67. This consistently high bias reveals that GEOS-CF generally is unable to simulate low-level O<sub>3</sub> magnitude well.

414 For the mid-level, the overall correlation reveals that GEOS-CF and GEOS-Chem both have a weak relationship with the  
 415 lidar (R = 0.22 and R = 0.12, respectively) (Figure S7 - bottom panel). This indicates that neither model is able to simulate  
 416 mid-level O<sub>3</sub> pattern well. GEOS-Chem consistently underestimates the magnitude of mid-level O<sub>3</sub> with a systematic high  
 417 negative normalized bias ranging from -0.44 to -0.18, for all clusters, while GEOS-CF has a lower and non-systematic  
 418 normalized bias ranging from -0.22 to 0.28. Overall, both models are not able to simulate the variability of O<sub>3</sub> nor the magnitude  
 419 well in the mid-level. The overall analysis in this sect. provides a fundamental but condensed assessment of model  
 420 performance. In the next sect., the cluster specific differences reveal additional model performance insight that would be  
 421 conceivably overlooked when evaluating overall performance.

422

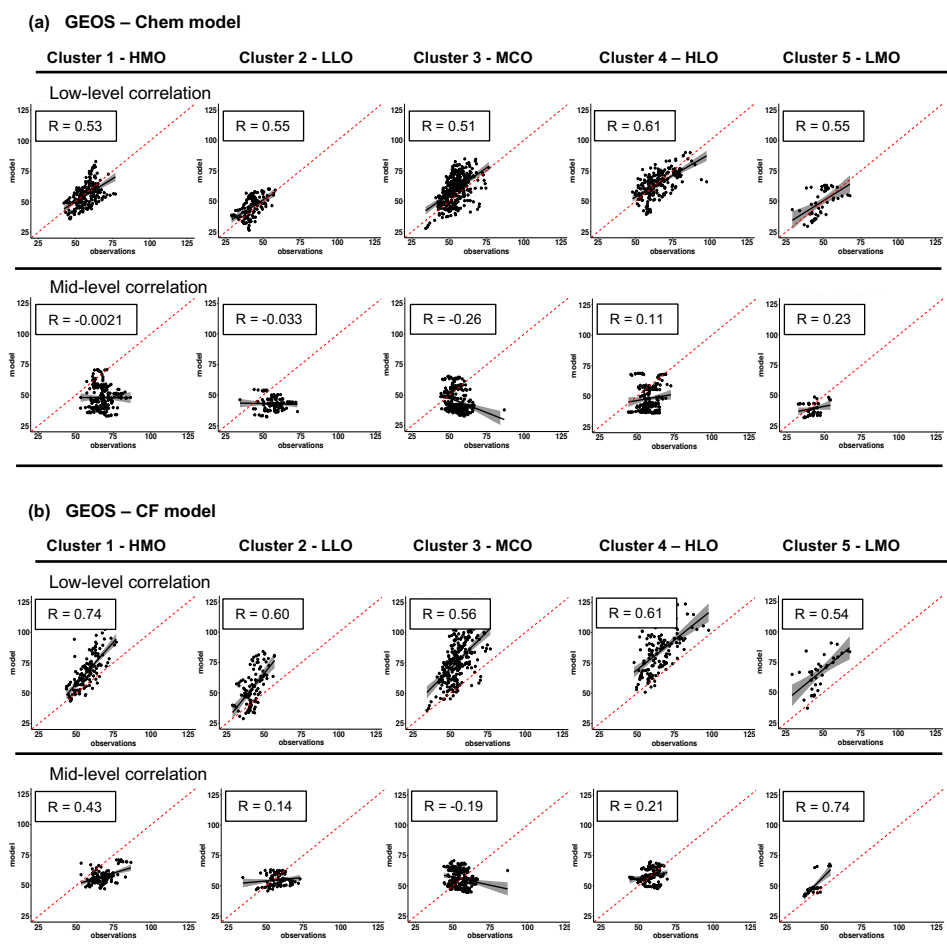
### 423 3.3.2 Model evaluation based on lidar clusters



424

425 **Figure 6.** Mean profile curtain spatial O<sub>3</sub> difference (model – lidar observations) for each cluster (1 – 5). GEOS-Chem  
 426 differences (a) and GEOS-CF differences (b).

427



429

430 **Figure 7.** O<sub>3</sub> correlation between lidar observations and a) GEOS-Chem model simulation results and b) GEOS-CF model  
 431 results by each cluster split by low-level (top panel) and mid-level (bottom panel).

432

433 Significant cluster by cluster differences are unmasked in evaluating the models based on the established O<sub>3</sub> behavior  
 434 cases. To quantify the results illustrated in Figure 4, we show spatial O<sub>3</sub> differences (model – lidar observations) for each  
 435 cluster (Figure 6) as well as individual cluster correlation (Figure 7) (subsequent cluster calculated normalized biases and  
 436 correlation can be found in Table S1). Evaluating the individual cluster biases and correlation reveal more in-depth model  
 437 discrepancies as well as areas where the models perform well.

438

439 In the low-level, GEOS-CF has a similar performance ability for the HMO, HLO, and LMO clusters with high positive  
 440 biases at + 0.30, + 0.41, and + 0.45 respectively. These higher biases imply GEOS-CF has difficulty capturing moderate O<sub>3</sub>  
 441 concentrations below 2000 m (HMO and LMO) as well as the in the high O<sub>3</sub> cases (HLO). GEOS-CF also has a high positive  
 bias (+ 0.50) in the LLO cluster indicating that GEOS-CF struggles to capture the lower O<sub>3</sub> concentrations in the low-level.



442 This is warranted as models are intended to approximate and are not usually able to capture extremes (high or low) but GEOS-  
443 CF also seems to struggle capturing moderate cases as well. In the low-level, GEOS-Chem has the best performance (minimal  
444  $-0.04$  bias and strong correlation,  $R = 0.61$ ) in HLO, which is the cluster with the highest low-level  $O_3$  accumulation (refer to  
445 Figure 4a). The second-best performance for GEOS-Chem in the low-level follows closely behind (minimal  $+0.07$  bias and  
446 fair correlation,  $R = 0.55$ ) in LLO, the cluster with the lowest  $O_3$  accumulation. These results suggest GEOS-Chem actually  
447 performs well in cases of high  $O_3$  as well as cases of low  $O_3$  with a slight tendency to overpredict lower  $O_3$  concentrations and  
448 underpredict higher  $O_3$  concentrations. This challenges the overall assumption that models struggle to capture extreme cases  
449 since GEOS-Chem actually performs best in simulating both extreme cases of high  $O_3$  in HLO and, again, low  $O_3$  in LLO.  
450 GEOS-Chem has a similar performance for the LMO and HMO clusters with negative biases of  $-0.10$  and  $-0.09$ , respectively.  
451 GEOS-Chem is also able to capture the moderate  $O_3$  in both of these clusters well with slight underestimations.

452 Both models perform the worst (in comparison with the other clusters) in the low-level in the MCO cluster with a  $+0.13$   
453 bias for GEOS-Chem and  $+0.67$  bias for GEOS-CF. As described in Sect. 3.1, MCO is the most common cluster with moderate  
454 - high average  $O_3$  concentrations in the low-level (refer to Figure 3b). Although GEOS-Chem has the worst performance in  
455 the MCO cluster, it is not necessarily a poor performance. The performance follows the conclusion previously made that  
456 GEOS-Chem can fairly simulate moderate  $O_3$  in the low-level although, in this case, with slight overestimations. Contrarily,  
457 the GEOS-CF performance in the MCO cluster reveals a more substantially high positive bias. This stands out as models are  
458 usually able to capture moderate levels (e.g., non-extreme cases). Evaluating the full temporal and vertical profile indicates  
459 that the higher GEOS-CF bias in the MCO cluster is additionally influenced by the greater overestimation of morning  $O_3$ , not  
460 solely the afternoon  $O_3$ . This is different to the performance in the LLO and LMO clusters where GEOS-CF also had a high  
461 positive bias in the low-level but does better simulating the early morning  $O_3$  magnitude. A similar conclusion can be drawn  
462 when evaluating the low-level GEOS-Chem performance. HMO, LLO, MCO, and LMO all share 'higher' biases (rounding to  
463  $\pm 0.10$ ), but the highest bias is found in the MCO cluster. Analogous to GEOS-CF, this can similarly be attributed to GEOS-  
464 Chem overestimating morning  $O_3$  the worst in the MCO cluster in contrast to the better early morning estimation in the other  
465 clusters.

466 In the mid-level, GEOS-Chem underestimates  $O_3$  magnitude to the greatest extent in the HMO and the LLO cluster (both  
467 bias =  $-0.44$ ), which are both clusters with higher mid-level  $O_3$  concentrations (refer to Figure 3c). GEOS-Chem performs  
468 similarly in the HLO and MCO clusters, with a negative mean bias of  $-0.30$  and  $-0.27$ , respectively. This indicates that  
469 GEOS-Chem most struggles to simulate higher concentrations of  $O_3$  in the mid-level. The GEOS-Chem model actually never  
470 reaches  $O_3$  cluster averages greater than 50 ppb, directly divulging the greater systemic negative bias in the mid-level. GEOS-  
471 Chem simulates LMO mid-level  $O_3$  magnitude the best ( $-0.18$  bias), which is the cluster with the lowest  $O_3$  average ( $< 45$   
472 ppb). Although for the LMO cluster GEOS-Chem has a lower bias, the correlation is still poor ( $R = 0.23$ ) which indicates that  
473 the model is relatively capable of simulating mid-level  $O_3$  only when the case devises lower concentrations but still fails to  
474 replicate any  $O_3$  variability and pattern.

475 On the other hand, GEOS-CF does best simulating LLO, MCO, and HLO, which are all clusters with moderate O<sub>3</sub> in the  
476 mid-level ( $\geq 50$  and  $\leq 70$  ppb). GEOS-CF has the highest bias in the LMO cluster (+ 0.28), the cluster with the lowest mid-  
477 level O<sub>3</sub> magnitude. GEOS-CF also has the strongest correlation in the same LMO cluster (R = 0.74). This is a unique case  
478 where although GEOS-CF is not able to capture the magnitude in the mid-level, it is able to capture the pattern of low O<sub>3</sub> well.  
479 Comparing the full multi-dimensional lidar and model mean profile curtains it is evident that in the LMO cluster, the GEOS-  
480 CF model simulates a similar mid-level O<sub>3</sub> pattern in the early morning/afternoon that is captured in the mean lidar curtain  
481 profile. The second worst performance for GEOS-CF is the underestimation of mid-level O<sub>3</sub> in the HMO cluster, contrarily  
482 the cluster with the highest mid-level O<sub>3</sub> ( $\geq 70$  ppb). This supports the previous conclusion that although GEOS-CF has a  
483 relatively lower biases in the mid-level, the model still struggles to simulate the extreme O<sub>3</sub> cases. Although GEOS-CF  
484 underestimates O<sub>3</sub> magnitude in the HMO cluster, it actually has a higher correlation than most of the other clusters (R = 0.43)  
485 (Figure 7, Table S1). In comparing the full multi-dimensional lidar and model mean profile curtain (Figure 3), GEOS-CF does  
486 a fair job connecting the mid-level higher O<sub>3</sub> pattern in the early morning that develops down to the low-level later in the  
487 afternoon. From this we can draw a conclusion that GEOS-CF is better able to capture mid-level O<sub>3</sub> patterns earlier in the  
488 temporal profile leading to higher correlations with the lidar.

489

### 490 3.3.3 Cluster approach and model conclusions

491 Several studies rely on case study investigations or grouping data by altitude to evaluate model performance. As  
492 demonstrated in Sect. 3.3.1, we can evaluate the overall summarized the model profile curtains O<sub>3</sub> against the lidar profile  
493 curtains and come to the simple conclusion that both models fairly simulate low-level O<sub>3</sub> but struggle to simulate mid-level  
494 O<sub>3</sub>. However, a systematic and comprehensive understanding of the different photochemical regimes in coastal regions does  
495 not only require case studies and overall summaries. The clustering approach allows for a comprehensive yet still detailed  
496 evaluation of the different photochemical regimes in coastal regions utilizing the lidar derived full profile curtains.  
497 Additionally, using the clusters, we can efficiently evaluate the ability of the models to simulate many different cases of O<sub>3</sub>.  
498 This approach revealed specific O<sub>3</sub> cases in which the models perform well and others where the models fail that would have  
499 been overlooked by solely considering the overall results. Using the clustering, we are able evaluate how the cluster specific  
500 differences (Figure 6, Figure 7, and Table S1) reveal additional model performance insight and specific gaps that would be  
501 conceivably overlooked when evaluating overall performance.

502 It is warranted that models struggle simulating extreme events/cases such as seen in the low-level in the HLO cluster and  
503 in the LLO cluster. However, GEOS-Chem performs best in both clusters with minimal biases and strong to fair correlations.  
504 Our result suggest that GEOS-Chem does a much better job simulating extreme O<sub>3</sub> cases in the low-level than expected. This  
505 specific model feature is not eminent when evaluating overall performance. Additionally, overall GEOS-Chem performs  
506 poorly in the mid-level. The detailed analysis granted by the cluster approach reveals GEOS-Chem has the lowest bias in the  
507 LMO cluster signifying the model is better able to capture low O<sub>3</sub> conditions in the mid-level. The overall high systemic  
508 positive bias for GEOS-CF in the low-level is further dissected when evaluating the individual clusters. GEOS-CF

509 systematically overestimates low-level O<sub>3</sub>, but the individual clusters indicate that the model has a better correlation with O<sub>3</sub>  
510 in HMO cases. An even more profound case is exposed in which GEOS-CF has a strong correlation with mid-level O<sub>3</sub> in the  
511 LMO cases despite having a low correlation overall. This concludes that in cases where the GEOS-CF model struggles to  
512 reproduce O<sub>3</sub> concentrations, the model can still capture the O<sub>3</sub> variability seen by the lidar measurements.

513 The clustering approach also reveals more discrepancies in the models such as in the MCO cluster. The advantage of  
514 evaluating full temporal and vertical profile curtains indicates that overestimation of early morning O<sub>3</sub> throughout the low-  
515 level leads to the poorer performances in MCO for both models. The overestimation of morning O<sub>3</sub> in GEOS-CF adds to the  
516 systemic overestimation in afternoon O<sub>3</sub> contributing the greater bias and poorer correlation. The same case can be found in  
517 the GEOS-Chem MCO cluster performance but to a lesser extent as GEOS-Chem has a much lower positive bias. Previous  
518 studies have found that excessive vertical mixing leads to overestimation of O<sub>3</sub> near the surface as well as underestimation of  
519 O<sub>3</sub> night-time depletion resulting in overestimation of O<sub>3</sub> the next day (Dacic et al., 2020; Keller et al., 2021; Travis &  
520 Jacob, 2019). The titration that occurs at night after the initial afternoon build up requires successful simulation to prevent the  
521 model beginning the following day with higher O<sub>3</sub> than is observed which can lead to the overprediction of O<sub>3</sub> later that day.  
522 Therefore, in the given case where there is an O<sub>3</sub> event that lasts more than one day (at the same lidar location), the model will  
523 likely underestimate O<sub>3</sub> night-time depletion, overpredict morning O<sub>3</sub>, and subsequently overpredict the afternoon build-up.  
524 Given multiple cases of multi-day high O<sub>3</sub> events from the lidar measurements (17 total from HMO, MCO, and HLO), this is  
525 likely one of the reasons for GEOS-CF overestimating early and therefore afternoon O<sub>3</sub> in these high O<sub>3</sub> cases in the low-level.  
526 In Figure 6, GEOS-CF exhibits the greatest afternoon O<sub>3</sub> overprediction in MCO and HLO. In HLO alone, there were 4 (out  
527 of 18) of the profiles that were consecutive while in MCO there were 8 (out of 28). This gives explanation for upwards of 22  
528 – 29 % of the overestimation of O<sub>3</sub> in the profile curtains of these clusters. These multi-day O<sub>3</sub> events are particularly important  
529 as they can indubitably lead the models to higher overprediction of afternoon O<sub>3</sub>. As the full lidar profile curtains reveal, the  
530 models tend to overestimate early morning O<sub>3</sub> in the MCO cases which links to the overestimation in afternoon O<sub>3</sub> as well.

531 Both models have a better ability to simulate early morning O<sub>3</sub> magnitude and pattern for other clusters than the MCO.  
532 For example, GEOS-CF does best simulating morning low-level O<sub>3</sub> in cases of lower O<sub>3</sub> extent (LLO and LMO). Excluding  
533 MCO, GEOS-Chem does not have such an issue overestimating low-level O<sub>3</sub> in the afternoon. In the other clusters, GEOS-  
534 Chem actually underpredicts early morning low-level O<sub>3</sub> in the full vertical profile. An underestimation of early morning O<sub>3</sub>  
535 does not warrant the same build-up up of afternoon O<sub>3</sub>. This gives some explanation to why GEOS-Chem underpredicts the  
536 other clusters with higher O<sub>3</sub> concentrations in the low-level (HMO and HLO). In the mid-level GEOS-Chem has a systemic  
537 high negative bias for all clusters, consistently underestimating O<sub>3</sub> but the clusters reveal a better performance in LMO, the  
538 cluster with lowest mid-level O<sub>3</sub> extent. It is evident that the model cannot simulate cases with higher O<sub>3</sub> concentrations in the  
539 mid-level but simulates low O<sub>3</sub> cases better. On the other hand, GEOS-CF results indicate a lower non-systemic bias in the  
540 mid-level. Since the version of GEOS-Chem used in this study was run with the tropchem chemistry mechanism which  
541 excludes stratospheric chemistry (now obsolete with current GEOS-Chem developments) and GEOS-CF uses the UCX  
542 chemistry mechanism that includes stratospheric chemistry, this may allude to better performance of GEOS-CF in simulating

543 higher O<sub>3</sub> concentrations in the mid-level. Both models indicate weak correlations with the lidar observations in the mid-level  
544 and it is apparent that both models struggle to capture the pattern of O<sub>3</sub> behavior in the mid-level. This could be due to multiple  
545 model inefficiencies such as the coarse model resolutions. Although GEOS-CF has a finer resolution than GEOS-Chem, it still  
546 may not be sufficient in horizontal and vertical grid resolution to replicate the O<sub>3</sub> variations captured in the 2-D lidar  
547 observations.

548 There are additional model discrepancies that can lead to underestimations of O<sub>3</sub> in GEOS-Chem in the mid-level that  
549 was found in all 5 clusters. One gap in the GEOS-Chem model could be the representation of tropospheric halogen chemistry  
550 which has a large effect of coastal O<sub>3</sub> production. Newer updates to the GEOS-Chem model (v12.9) have included updated  
551 tropospheric halogen chemistry mechanisms (iodine, bromine, and chlorine) (Wang et al., 2021). This study found that the  
552 updated halogen chemistry actually worsens the overall underestimation of O<sub>3</sub> throughout the troposphere, specifically in the  
553 northern hemisphere, indicating further investigation of halogen chemistry is needed for better model representation. Another  
554 study finds a similar conclusion in the proper representation of cloud uptake and tropospheric chemistry (Holmes et al., 2019).  
555 This study found that implementing an updated, more accurate, and stable cloud entrainment-limited uptake in the GEOS-  
556 Chem model reduces the sensitivity of oxidants and aerosol chemistry in the troposphere but still had little effect on O<sub>3</sub> model  
557 comparison to observations (such as sonde and aircraft). This is due to the environmental variability being much higher than  
558 the effect of NO<sub>x</sub> and O<sub>3</sub> cloud chemistry but still warrants further testing. The role lightning plays in tropospheric oxidation  
559 is another feature that is commonly misrepresented in global models and can affect O<sub>3</sub> simulation (Mao et al., 2021). These  
560 are all examples of features that if not simulated correctly can lead to misestimations of O<sub>3</sub>. The clustering approach allows us  
561 to organize the detailed lidar measurements to scope out specific cases where these misrepresentations occur. These previous  
562 studies also highlight the importance of lidar measurements and their ability to depict tropospheric emission development and  
563 behavior throughout the vertical profile and diurnal cycle which can be used to constrain model emissions and improve  
564 simulations.

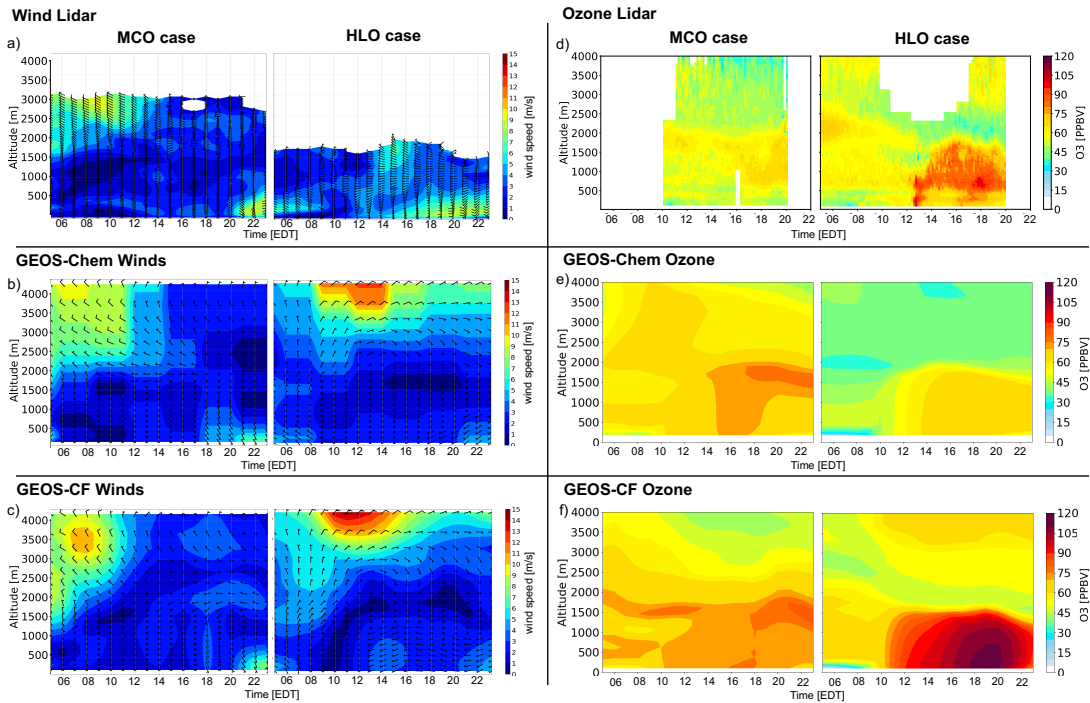
565 Although this analysis proves to be a useful technique to characterize the largely variably O<sub>3</sub> behavior in coastal regions  
566 and evaluate the subsequent model performance, there are also limitations. In this study we are comparing single point lidar  
567 versus model output, therefore we cannot simply state that the model is incorrect. We make conclusions and draw biases based  
568 on the ability to subset a grid point and compare that to a single point lidar curtain to the best ability but that still leaves an  
569 uncertainty. The high vertical and spatiotemporal resolution reveal intricate details about the behavior of O<sub>3</sub> during these  
570 campaigns. O<sub>3</sub> lidars have a unique advantage, compared to traditional surface measurements, in measuring vertical  
571 distribution of O<sub>3</sub> with respect to time. This advantage is of great value when investigating model ability in simulating the  
572 spatial and temporal distribution of O<sub>3</sub> and can provide crucial information in understanding surface O<sub>3</sub> events.

573

### 574 **3.4 Cluster derived case studies to evaluate modeled wind and ozone**

575 Meteorological factors such as wind speed and direction can directly impact whether a coastal region will experience  
576 clean air or O<sub>3</sub> exceedances. When local meteorological processes such as sea/bay breeze occur at such a fine scale, equally

577 fine resolution measurements are essential in capturing this. The Doppler wind lidar offers a focus on fine details that are only  
 578 revealed in the multi-dimensional data which allows for such a comprehensive evaluation of the established O<sub>3</sub> cluster profile  
 579 curtains. In this sect., we evaluate the 2-D relationship between wind and O<sub>3</sub> to assess model performance using lidar and  
 580 model derived profile curtains (Figure 8). We derived two specific case studies, each from a different cluster: MCO = 17 June  
 581 2018 and HLO = 30 June 2018. Utilizing the derived clusters, the case studies were chosen to focus on high low-level O<sub>3</sub>  
 582 behavior cases with a goal of evaluating possible sea/bay breeze events. The two case studies are both from the HMI location  
 583 during the OWLETS-2 campaign. There are consistent Doppler lidar measurements throughout the low-level (<2000 m) which  
 584 allows for a direct comparison with the simulated profiles; therefore, the focus of the following analysis will be on the low-  
 585 level altitudes. The deficit of mid-level observed wind data disallows for a conclusive and concrete evaluation of simulated  
 586 mid-level O<sub>3</sub>.



587  
 588 **Figure 8.** Profile curtains of wind speed/direction (a-c) and ozone (d-f) from the lidar (top panel), GEOS-Chem (middle panel),  
 589 and GEOS-CF (bottom panel). Results from OWLETS-2 at HMI.

590

### 591 3.4.1 Sea breeze event interpretation

592 GEOS-Chem and GEOS-CF both struggle to capture low-level wind speed and direction in both MCO and HLO cases  
 593 (Figure 8a-c). In the MCO case, the Doppler wind lidar captures a wind direction shift from westerly to easterly winds  
 594 beginning at 06:00 EDT accompanied by calm winds (approximately 0 m s<sup>-1</sup>) indicating a likely common sea/bay breeze event.  
 595 The timing of the start of this event is simulated well but the models fail to predict an actual well-defined wind shift, instead

596 merely simulating  $0 \text{ m s}^{-1}$  winds after 05:00 EDT. It is apparent that the models struggle to capture the finer processes such as  
597 a sea/bay breeze which could have likely led the underprediction of wind speed. It is important to note that GEOS-Chem runs  
598 with offline meteorology, averaged every 3 hours. Since sea/bay breezes often happen at a finer temporal resolution, the GEOS-  
599 Chem model is at a disadvantage in modelling such fine processes. A wind direction shift is also depicted in the HLO case,  
600 with westerly winds early in the morning and a shift to south-easterly winds later in the temporal profile (at about 10:00 EDT).  
601 This could also likely be an early onset sea breeze event which could have contributed to the high observed  $\text{O}_3$  concentrations  
602 in the afternoon. Again, the exact timing of the start of the wind shift is captured by the models but then no defined directional  
603 shift and little to no winds are simulated after. Both the MCO case and HLO case observe increased wind speeds near the  
604 surface, first before 08:00 EDT then again in the evening. Both models underestimate the extent of the increased wind speeds.  
605

### 606 **3.4.2 Relation to ozone cases and clustering**

607 In this sect., the wind lidar curtains will be assessed in relation to the  $\text{O}_3$  lidar profile curtains and the model performance.  
608 The results in sect. 3 revealed that both models had the highest bias and lowest correlation simulating low-level  $\text{O}_3$  in MCO.  
609 Evaluating the wind and  $\text{O}_3$  lidar profile curtains against the model simulations helps paint a better picture as to why. Similar  
610 to the MCO cluster mean curtain profile, early morning low-level  $\text{O}_3$  in each case is overestimated by both models (Figure 8e,  
611 f). There is higher  $\text{O}_3$  captured in the lidar curtain profile, but it is constrained between 1500 – 2000 m. Both models bring this  
612 higher  $\text{O}_3$  pattern down to the surface (below 500 m) overestimating  $\text{O}_3$  throughout the low-level. Since both models predict  
613 little to no winds during this time, this could contribute to overestimations of  $\text{O}_3$  near the surface.

614 In the HLO case, GEOS-CF overestimates low-level  $\text{O}_3$  while GEOS-Chem underestimates low-level  $\text{O}_3$ . From sect. 3.3  
615 the results revealed that although GEOS-CF has a high positive normalized bias for low-level  $\text{O}_3$  in HLO, the model had a  
616 reasonable relationship ( $R = 0.61$ ) with the  $\text{O}_3$  lidar measurements. This is corroborated with the individual HLO case (Figure  
617 8f) as GEOS-CF is better able to simulate the development of  $\text{O}_3$  in the low-level, especially in the early morning. The GEOS-  
618 CF modeled winds mirror this performance with a better reproduction of the wind shift in HLO (Figure 8c). While GEOS-  
619 Chem has a lower normalized bias for low-level  $\text{O}_3$  in the HLO cluster, GEOS-Chem consistently underestimates wind speed  
620 and fails to reproduce any wind shifts. This reveals that in the possible sea breeze event, the two models do not perform equally.  
621 Since GEOS-Chem is an offline CTM using archived meteorology and GEOS-CF simulates atmospheric composition  
622 simultaneously with meteorology (online), the replication of a sea breeze case would not necessarily be comparable.

623 In most cases, sea/bay breeze events can contribute to high concentrated daytime  $\text{O}_3$  events in which  $\text{O}_3$  is recirculated  
624 throughout the region. Such cases would likely lead to a similar curtain profile as seen in the HLO case (Figure 8a), where  
625 high  $\text{O}_3$  in the morning is likely associated with the higher  $\text{O}_3$  at the surface in the afternoon. But it is apparent that the cases  
626 for MCO and HLO are dissimilar. We would expect per the clustering approach that sea breeze cases would most likely be  
627 assigned to the same cluster, but this is not the case here. Investigating the full lidar and model profile curtains for the two  
628 cases gives us more information as to why these two curtains are not in the same cluster. It is evident that the HLO case has  
629 much higher afternoon  $\text{O}_3$  near the surface (below 1000 m) than the MCO case, with peaks  $> 75$  ppb at both 12:00 and again

630 at 16:00 EDT. In contrast, the MCO case has higher afternoon O<sub>3</sub> concentrations captured above 2000 m than the HLO case.  
631 The HLO case has high O<sub>3</sub> in the afternoon, but it is constrained to the lower 2000 m and just above this high O<sub>3</sub> plume, there  
632 is an O<sub>3</sub> deficit of almost 50 ppb. Although the MCO case also reveals lower O<sub>3</sub> above 2000 m, the vertical gradient in this  
633 case is not as stark. This is also replicated in both models which simulate lower O<sub>3</sub> directly above the high surface O<sub>3</sub> in the  
634 HLO cluster but simulate much higher O<sub>3</sub> above 2000 m in the MCO cluster. From their distinct vertical and temporal behavior,  
635 it is easy to conclude why these two cases were not assigned to the same cluster.

636 The cases elected for MCO and HLO give reason to address the difficulty simulating complex coastal mechanisms.  
637 Despite the fact that MCO and HLO both indicated prospective sea/bay breeze cases, the results of the simulated winds and  
638 O<sub>3</sub> were distinctive. Simulating complex sea/bay and land relations is imperative for correctly mitigating high O<sub>3</sub> cases. To  
639 accurately simulate such complex exchanges, high resolution vertical and horizontal simulations are needed. Because of the  
640 models' relatively coarse resolutions (nominally 50 and 25 km horizontal resolution; 72 vertical levels), the fine-scale vertical  
641 wind gradients and horizontal wind shifts are difficult to resolve and, in these cases, not fully able to replicate. This study also  
642 acknowledges the need for an evaluation of other modeled factors, such as divulged in sect. 3.3.3, considering the possible  
643 confounding effects on modeled O<sub>3</sub> outcome.

644

#### 645 **4. Conclusion**

646 We developed and tested a clustering method on a suite of 91 multi-dimensional lidar O<sub>3</sub> profile curtains retrieved  
647 from three recent land/sea campaigns (OWLETS-1, OWLETS-2, and LISTOS), during the summer months of 2017 and 2018.  
648 The K-Means clustering algorithm, driven by 8 well defined features, was applied to categorize the fine resolution O<sub>3</sub> data,  
649 revealing five distinct O<sub>3</sub> behavior cases that are distinct in pattern and magnitude vertically and temporally. We present five  
650 different clusters of O<sub>3</sub> behavior identified as: highest mid-level O<sub>3</sub> (HMO) cluster; lowest low-level O<sub>3</sub> (LLO) cluster; most  
651 common O<sub>3</sub> (MCO) cluster; highest low-level O<sub>3</sub> (HLO); lowest mid-level O<sub>3</sub> (LMO) cluster. The results indicate that fine  
652 resolution data can be used to differentiate the behavior of O<sub>3</sub> in a region and classify different cases of O<sub>3</sub> exploiting the  
653 multiple dimensions. The clustering approach allowed us to characterize the range of highly variable vertical and temporal  
654 coastal O<sub>3</sub> behavior for the duration of these campaigns which can be a good indicator of how O<sub>3</sub> behaves in general in these  
655 coastal regions during the summer months. Furthermore, this approach could be used by states to better identify different O<sub>3</sub>  
656 photochemical regimes and frequency beyond just surface sampling.

657 We evaluated the performance of two CTMs, GEOS-Chem and GEOS-CF, in these complex environments. Overall,  
658 the models have the greatest difficulty simulating the vertical extent and variability of O<sub>3</sub> concentrations in the mid-level, with  
659 weak overall relationships to the lidar observations ( $R = 0.12$  and  $0.22$ ). GEOS-Chem had a systematic high negative bias and  
660 GEOS-CF an overall lower unsystematic bias range. In the low-level, GEOS-Chem had overall low unsystematic bias range  
661 and fair relationship with the lidar observations ( $R = 0.66$ ), while GEOS-CF had a systematic high positive bias but overall  
662 fair relationship ( $R = 0.69$ ).

663 Utilizing the curated clusters reveals new model insight that is neglected in the overall performance analysis. The  
664 cluster approach divulges specific model limitations but also cases in which the models perform well. GEOS-Chem simulates  
665 low-level O<sub>3</sub> cases best in the HLO and LLO clusters and the worst in the MCO cluster. HLO and LLO are the clusters with  
666 the most extreme (low and high) O<sub>3</sub> cases while MCO is the most common cluster with moderate O<sub>3</sub>. This concludes that  
667 GEOS-Chem does best simulating extreme low-level O<sub>3</sub> but struggles to capture the frequently occurring moderate O<sub>3</sub>  
668 behavior. GEOS-CF also has the greatest overestimations for low-level O<sub>3</sub> in the MCO cluster. Evaluating the full profile  
669 curtain reveals that this overestimation can be most attributed to the greater overestimation of early morning O<sub>3</sub>. This feature  
670 is unique to the MCO cluster and warrants further investigation as O<sub>3</sub> left in the residual layer can contribute to higher O<sub>3</sub> in  
671 the afternoon and proves to be a challenge for CTMs. The value of lidar measurements is reflected in its ability to reveal these  
672 features.

673 Both models share poor performances in the mid-level but there are specific cases that stand out in the clustering  
674 results, specifically the LMO cluster, in which GEOS-CF shares a good agreement with the lidar measurements. It can be  
675 concluded that although the model struggles to simulate O<sub>3</sub> magnitude, it can relatively emulate the mid-level O<sub>3</sub> pattern in  
676 LMO. This is also apparent in the MCO cluster, in which the pattern of higher mid-level O<sub>3</sub> that suggests a relationship with  
677 the low-level O<sub>3</sub> is simulated fairly in the GEOS-CF model. This pattern is also a rare feature that is captured in the lidar that  
678 demonstrates the significance of the measurements. The greater underestimations of mid-level O<sub>3</sub> for GEOS-Chem can be  
679 alluded to multiple model discrepancies. Since the GEOS-Chem version and mechanism used in this study (tropchem) only  
680 considers tropospheric chemistry we can expect the performance in the mid-level to have deficiencies. Although GEOS-CF is  
681 run with the combined tropospheric and stratospheric chemistry mechanism, has a better grid resolution, and is an online  
682 model, there are still limitations to both models especially when simulating mid-level O<sub>3</sub>. Known model errors and coarse  
683 horizontal and vertical grid resolution contribute to the difficulty in simulating fine-scale coastal O<sub>3</sub> variability. There are many  
684 contributing model factors that can be affecting the performance of GEOS-Chem and GEOS-CF that were mentioned in this  
685 study not solely coarse model resolution.

686 A unique value of the clustering approach on multi-dimensional lidar data is that it offers a convenient way to ascertain  
687 different O<sub>3</sub> case studies. An example of this is our evaluation of two cases studies from the MCO and HLO clusters. Modeled  
688 winds were evaluated using Doppler wind lidar data observed during the OWLETS-2 campaign. The wind lidar data was  
689 mostly limited to lower altitudes (< 2000 m), which allowed for wind speed and direction validation at the low-level. The  
690 morning wind deceleration and directional shifts (onshore to offshore) illustrated in lidar profile curtains indicate a possible  
691 sea/bay breeze event in both case studies. This is likely another contributor that led to enhanced surface O<sub>3</sub> in these cases. Due  
692 to the coarser model resolution, GEOS-Chem and GEOS-CF were not able to capture the sea breeze phenomena in these cases  
693 which could have facilitated in the high O<sub>3</sub> biases for these clusters. With GEOS-CF having a finer horizontal resolution than  
694 GEOS-Chem, the results reveal minimal advantages simulating the pattern of wind speeds better but none in simulating the  
695 wind directional shifts. This affirms that the spatial resolution of GEOS-CF (~25 km) is still not fine enough for mesoscale  
696 processes such as the sea/bay breeze. Although a regional model analysis is out of the scope of this study, we propose to use



697 multi-dimensional lidar measurements to evaluate finer regional modeling in our future work. We acknowledge that other  
698 factors, aside from model resolution, contribute to discrepancies in modeled coastal O<sub>3</sub> and further warrant a deeper evaluation.  
699 The clustering approach on lidar measurements offers an unmatched ability to pinpoint these features.

700 This work is the first time that all three associated campaign lidar data have been analyzed in conjunction. In utilizing  
701 the highly detailed suite of multi-dimensional lidar data, we are able to comprehensively explore the behavior and variability  
702 of coastal O<sub>3</sub> for the duration of the campaigns. Applying the clustering analysis directly to the lidar O<sub>3</sub> data emerges as a  
703 useful and robust approach for identifying O<sub>3</sub> patterns during the highly polluted summer months in coastal environments.  
704 Since the time of the OWLETS and LISTOS campaigns, the lidar instrument systems have been updated and are now more  
705 fully automatized for use eliminating such constraints faced in this study. Further observations using lidar instruments should  
706 be especially valuable in investigating coastal O<sub>3</sub> behavior as it can divulge the finer-scale O<sub>3</sub> characteristics that remain  
707 difficult to successfully simulate in CTMs. The time-height and fine resolution measurements only available from multi-  
708 dimensional lidar instruments were vital in allowing us to form these conclusions.

709 This kind of evaluation allows for detailed model assessment of specific O<sub>3</sub> cases that are unmasked through the  
710 clustering analysis. Looking at the overall correlations, it would seem the models have a good relationship with the low-level  
711 lidar observations but looking into the cluster-by-cluster differences, the gaps within the models are elucidated. Using the  
712 cluster assignments, we are able evaluate how the cluster specific differences reveal additional model performance insight that  
713 could be conceivably overlooked when evaluating overall performance. This work is a middle ground between looking at  
714 specific cases (or dates) and summarizing overall model performance. Additionally, the clustering approach provides an  
715 abridged way to detecting distinctive case studies. We provide a new approach that allows a synopsis of summer coastal O<sub>3</sub>  
716 behavior and subsequently model performance without completely muting distinct O<sub>3</sub> features. Evaluating model performance  
717 for diverse O<sub>3</sub> behavior in coastal regions is crucial for improving the simulation and furthermore, mitigation of air quality  
718 events.

719 *Code availability.* Model code is available upon request to the first author.

720 *Data availability.* The GEOS-Chem model simulation data from this study is publicly accessible online at  
721 <https://doi.org/10.7910/DVN/V99LHT>. The GEOS-CF model data is publicly available online at their website  
722 [https://gmao.gsfc.nasa.gov/-weather\\_prediction/GEOS-CF/](https://gmao.gsfc.nasa.gov/-weather_prediction/GEOS-CF/). The lidar data is publicly available online at [https://www-](https://www-air.larc.nasa.gov/missions.htm)  
723 [air.larc.nasa.gov/missions.htm](https://www-air.larc.nasa.gov/missions.htm).

724 *Supplement.*

725 *Author contributions.* CB and YW conceived the research idea. CB wrote the initial draft of the paper and performed the  
726 analyses and model development. All authors contributed to the interpretation of the results and the preparation of the paper.

727 *Competing interests.* The authors declare that they have no conflict of interest.

728 *Acknowledgements.* This study is supported by NASA MUREP Graduate Fellowship (80NSSC19K1680). The Ozone Water-  
729 Land Environmental Transition Study (OWLETS-1, 2) and Long Island Sound Tropospheric Ozone Study (LISTOS) field  
730 measurements described here were funded by the NASA's Tropospheric Composition Program and Science Innovation Fund  
731 (SIF), Maryland Department of Environment, the National Oceanic and Atmospheric Administration (NOAA), the  
732 Environmental Protection Agency (EPA), the Northeast States for Coordinated Air Use Management (NESCAUM), and the  
733 New Jersey and Connecticut Departments of Energy and Environmental Protection. The authors acknowledge the principal  
734 investigators and data operators John Sullivan, Joel Dreessen, Ruben Delgado, William Carrion, and Joseph Sparrow as well  
735 as the guidance of the Tropospheric Ozone Lidar Network (TOLNet). LMOL and TROPOZ data are publicly available at  
736 (<https://www-air.larc.nasa.gov/missions/TOLNet/>). The OWLETS and LISTOS data are available at ([https://www-](https://www-air.larc.nasa.gov/)  
737 [air.larc.nasa.gov/](https://www-air.larc.nasa.gov/)). The Doppler wind data taken from the UMBC wind lidar and are publicly available at ([https://www-](https://www-air.larc.nasa.gov/cgi-bin/ArcView/owlets.2018)  
738 [air.larc.nasa.gov/cgi-bin/ArcView/owlets.2018](https://www-air.larc.nasa.gov/cgi-bin/ArcView/owlets.2018)). The GEOS-CF model simulation data were provided directly from the NASA  
739 Center Global Modeling and Assimilation Office (GMAO) at the Goddard Space Flight Center  
740 ([https://gmao.gsfc.nasa.gov/weather\\_prediction/GEOS-CF/](https://gmao.gsfc.nasa.gov/weather_prediction/GEOS-CF/)).

## 741 **References**

- 742 Alonso, A. M., Berrendero, J. R., Hernández, A., and Justel, A.: Time Series Clustering Based on Forecast Densities,  
743 Computational Statistics & Data Analysis, 51(2), 762–776., <https://doi.org/10.1016/j.csda.2006.04.035>, 2006.
- 744 Banta, R. M., Senff, C. J., Nielsen-Gammon, J., Darby, L. S., Ryerson, T. B., Alvarez, R. J., Sandberg, S. P., Williams, E. J.,  
745 and Trainer, M: A bad air day in Houston. Bulletin of the American Meteorological Society, 86(5), 657–  
746 670. <https://doi.org/10.1175/BAMS-86-5-657>, 2005.
- 747 Bernier, C., Wang, Y., Estes, M., Lei, R., Jia, B., Wang, S., and Sun, J.: Clustering Surface Ozone Diurnal Cycles to Understand  
748 the Impact of Circulation Patterns in Houston, TX, Journal of Geophysical Research: Atmospheres, 124(23), 13457–  
749 13474., <https://doi.org/10.1029/2019jd031725>, 2019.
- 750 Caicedo, V., Rappenglueck, B., Cuchiara, G., Flynn, J., Ferrare, R., Scarino, A. J., Berkoff, T., Senff, C., Langford, A., and  
751 Lefer, B.: Bay Breeze and Sea Breeze Circulation Impacts on the Planetary Boundary Layer and Air Quality from an  
752 Observed and Modeled Discover-AQ Texas Case Study, Journal of Geophysical Research: Atmospheres, 124(13),  
753 7359–7378, <https://doi.org/10.1029/2019jd030523>, 2019.
- 754 Charrad, M., Ghazzali, N., Boiteau, V., and Niknafs, A.: NBCLUST: An R package for Determining the Relevant Number of  
755 Clusters in a Data Set, Journal of Statistical Software, 61,(6), <https://doi.org/10.18637/jss.v061.i06>, 2014.
- 756 Christiansen, B.: Atmospheric Circulation Regimes: Can Cluster Analysis Provide the Number?, Journal of Climate, 20(10),  
757 2229–2250., <https://doi.org/10.1175/jcli4107.1>, 2007.

758 Coggon, M. M., Gkatzelis, G. I., McDonald, B. C., Gilman, J. B., Schwantes, R. H., Abuhassan, N., Aikin, K. C., Arend, M.  
759 F., Berkoff, T. A., and Brown, S. S.: Volatile chemical product emissions enhance ozone and modulate urban  
760 chemistry, *Proceedings of the National Academy of Sciences*, 118,32, National Academy of Sciences,  
761 <https://doi.org/10.1073/pnas.2026653118>, 2021.

762 Couillard, M. H., Schwab, M. J., Schwab, J. J., Lu, C. H., Joseph, E., Stutsrim, B., Shrestha, B., Zhang, J., Knepp, T. N., and  
763 Gronoff, G. P., Vertical Profiles of Ozone Concentrations in the Lower Troposphere Downwind of New York City  
764 during LISTOS 2018-2019, *Journal of Geophysical Research: Atmospheres*, 126(23), e2021JD035108,  
765 <https://doi.org/10.1029/2021JD035108>, 2021.

766 Darby, L. S.: Cluster Analysis of Surface Winds in Houston, Texas, and the Impact of Wind Patterns on Ozone, *Journal of*  
767 *Applied Meteorology*, 44(12), 1788–1806., <https://doi.org/10.1175/jam2320.1>, 2005.

768 Davis, R. E., Normile, C. P., Sitka, L., Hondula, D. M., Knight, D. B., Gawtry, S. P., and Stenger, P. J.: A Comparison of  
769 Trajectory and Air Mass Approaches to Examine Ozone Variability, *Atmospheric Environment*, 44(1), 64–74.,  
770 <https://doi.org/10.1016/j.atmosenv.2009.09.038>, 2010.

771 De Young, R., Carrion, W., Ganoe, R., Pliutau, D., Gronoff, G., Berkoff, T., and Kuang, S.: Langley Mobile Ozone LIDAR:  
772 Ozone and Aerosol Atmospheric Profiling for Air Quality Research, *Applied Optics*, 56(3), 721,  
773 <https://doi.org/10.1364/ao.56.000721>, 2017.

774 Dreessen, J., Orozco, D., Boyle, J., Szymborski, J., Lee, P., Flores, A., and Sakai, R. K.: Observed Ozone over the Chesapeake  
775 Bay Land-Water Interface: The Hart-Miller Island Pilot Project, *Journal of the Air & Waste Management Association*,  
776 69, (11), 1312–1330, <https://doi.org/10.1080/10962247.2019.1668497>, 2019.

777 EPA NEI (National Emissions Inventory v1): Air Pollutant Emission Trends Data, available at:  
778 <http://www.epa.gov/ttn/chief/trends/index.html> last access: 23 June 2015.

779 Farris, B. M., Gronoff, G. P., Carrion, W., Knepp, T., Pippin, M., and Berkoff, T. A.: Demonstration of an off-Axis Parabolic  
780 Receiver for near-Range Retrieval of Lidar Ozone Profiles, *Atmospheric Measurement Techniques*, 12(1), 363–370,  
781 <https://doi.org/10.5194/amt-12-363-2019>, 2019.

782 Gelaro, R., Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., Randles, C. A., Darmenov, A.,  
783 Bosilovich, M. G. Reichle, R., Wargan, K., Coy, L., Cullather, R., Draper, C., Akella, S., Buchard, V., Conaty, A., da  
784 Silva, A. M., Gu, W., Kim, G-K., Koster, R., Lucchesi, R., Merkova, D., Nielsen, J. E., Partyka, G., Pawson, S., Putman,  
785 W., Rienecker, M., Schubert, S. D., Sienkiewicz, M., and Zhao, B.: The Modern-Era Retrospective Analysis for Research  
786 and Applications, Version 2 (Merra-2), *Journal of Climate*, 30(14), 5419–5454., <https://doi.org/10.1175/jcli-d-16-0758.1>,  
787 2017.

788 Gronoff, G., Robinson, J., Berkoff, T., Swap, R., Farris, B., Schroeder, J., Halliday, H.S., Knepp, T., Spinei, E., Carrion,  
789 W., Adcock, E.E., Johns, Z., Allen, D., Pippin, M.: A Method for Quantifying near Range Point Source Induced O<sub>3</sub>  
790 Titration Events Using Co-Located Lidar and Pandora Measurements, *Atmospheric Environment*, 204, 43–52,  
791 <https://doi.org/10.1016/j.atmosenv.2019.01.052>, 2019

792 Gronoff, G., Berkoff, T., Knowland, K.E., Lei, L., Shook, M., Fabbri, B., Carrion, W., Langford, A. O.: Case study of  
793 stratospheric Intrusion above Hampton, Virginia: lidar-observation and modeling analysis, *Atmospheric Environment*,  
794 259,118498, 1352-2310, <https://doi.org/10.1016/j.atmosenv.2021.118498>, 2021.

795 Han, J. and Kamber, M.: *Data mining: concepts and techniques*, San Francisco: Morgan Kaufmann Publishers, 2001.

796 Holmes, C. D., Bertram, T. H., Confer, K. L., Graham, K. A., Ronan, A. C., Wirks, C. K., and Shah, V.: The Role of Clouds  
797 in the Tropospheric NO<sub>x</sub> Cycle: A New Modeling Approach for Cloud Chemistry and Its Global Implications, *Geophys.*  
798 *Res. Lett.*, 46, 4980–4990, <https://doi.org/10.1029/2019GL081990>, 2019.

799 Hu, L., Keller, C. A., Long, M. S., Sherwen, T., Auer, B., Da Silva, A., Nielsen, J. E., Pawson, S., Thompson, M. A., Trayanov,  
800 A. L., Travis, K. R., Grange, S. K., Evans, M. J., and Jacob, D. J.: Global simulation of tropospheric chemistry at 12.5 km  
801 resolution: performance and evaluation of the GEOS-Chem chemical module (v10-1) within the NASA GEOS Earth  
802 system model (GEOS-5 ESM), *Geosci. Model Dev.*, 11, 4603–4620, <https://doi.org/10.5194/gmd-11-4603-2018>, 2018.

803 Kaufman, L. and Rousseeuw, P.: *Finding groups in data: An introduction to cluster analysis*, New York, Wiley, 1990.

804 Keller, C.A., Knowland, K.E., Duncan, B.N., Liu, J., Anderson, D.C., Das, S., Lucchesi, R.A., Lundgren, E.W., Nicely, J.M.,  
805 Nielsen, E., Ott, L.E., Saunders, E., Strode, S.A., Wales, P.A., Jacob, D.J., and Pawson, S.: Description of the NASA  
806 Geos Composition Forecast Modeling System GEOS-CF v1.0, *Journal of Advances in Modeling Earth Systems*, 13(4),  
807 <https://doi.org/10.1029/2020ms002413>, 2021.

808 Knowland, K.E., Keller, C.A., Lucchesi, R.: File specification for GEOS-CF products, GMAO office note No. 17 (version  
809 1.0). available from: [https://gmao.gsfc.nasa.gov/pubs/office\\_notes.php](https://gmao.gsfc.nasa.gov/pubs/office_notes.php), 32, 2019.

810 Knowland, K. E., Keller, C. A., Wales, P. A., Wargan, K., Coy, L., Johnson, M. S., Liu, J., Lucchesi, R. A., Eastham, S. D.,  
811 Fleming, E. L., Liang, Q., Leblanc, T., Livesey, N. J., Walker, K. A., Ott, L. E., and Pawson, S.: NASA GEOS  
812 Composition Forecast Modeling System GEOS-CF v1.0: Stratospheric Composition, 14(6), e2021MS002852,  
813 <https://doi.org/10.1002/essoar.10508148.1>, 2021.

814 Lawson, R. G., and Jurs, P. C.: New Index for Clustering Tendency and Its Application to Chemical Problems, *Journal of*  
815 *Chemical Information and Computer Sciences*, 30(1), 36–41., <https://doi.org/10.1021/ci00065a010>, 1990.

816 Leblanc, T., Brewer, M. A., Wang, P. S., Granados-Muñoz, M. J., Strawbridge, K. B., Travis, M., Firanski, B., Sullivan, J. T.,  
817 McGee, T. J., Sumnicht, G. K., Twigg, L. W., Berkoff, T. A., Carrion, W., Gronoff, G., Aknan, A., Chen, G., Alvarez,  
818 R. J., Langford, A. O., Senff, C. J., Kirgis, G., Johnson, M. S., Kuang, S., and Newchurch, M. J.: Validation of the TOLNet  
819 lidars: the Southern California Ozone Observation Project (SCOOP), *Atmos. Meas. Tech.*, 11, 6137–6162,  
820 <https://doi.org/10.5194/amt-11-6137-2018>, 2018.

821 Lei, L., Berkoff, T. A., Gronoff, G., Su, J., Nehrir, A. R., Wu, Y., Moshary, F., and Kuang, S.: Retrieval of UVB aerosol  
822 extinction profiles from the ground-based Langley Mobile Ozone Lidar (LMOL) system, *Atmos. Meas. Tech.*, 15, 2465–  
823 2478, <https://doi.org/10.5194/amt-15-2465-2022>, 2022.

824 Li, W., Wang, Y., Bernier, C., and Estes, M.: Identification of Sea Breeze Recirculation and Its Effects on Ozone in Houston,  
825 TX, during Discover-Aq 2013, *Journal of Geophysical Research: Atmospheres*, 125(22),  
826 <https://doi.org/10.1029/2020jd033165>, 2020.

827 Little, R. J., A., and Rubin, D., B.: *Statistical Analysis with Missing Data*. Second ed., Wiley, 1987.

828 Mao, J., Zhao, T., Keller, C. A., Wang, X., McFarland, P. J., Jenkins, J. M., and Brune, W. H.: Global Impact of Lightning-  
829 Produced Oxidants, *Geophys. Res. Lett.*, 48, <https://doi.org/10.1029/2021GL095740>, 2021.

830 McDuffie, E. E., Smith, S. J., O'Rourke, P., Tibrewal, K., Venkataraman, C., Marais, E. A., Zheng, B., Crippa, M., Brauer,  
831 M., and Martin, R. V.: A global anthropogenic emission inventory of atmospheric pollutants from sector- and fuel-specific  
832 sources (1970–2017): an application of the Community Emissions Data System (CEDS), *Earth Syst. Sci. Data*, 12, 3413–  
833 3442, <https://doi.org/10.5194/essd-12-3413-2020>, 2020.

834 Orbe, C., Oman, L. D., Strahan, S. E., Waugh, D. W., Pawson, S., Takacs, L. L., and Molod, A. M.: Large-scale atmospheric  
835 transport in GEOS replay simulations, *J. Adv. Model. Earth Syst.*, 9, 2545–2560, <https://doi.org/10.1002/2017MS001053>,  
836 2017.

837 Ring, A. M., Canty, T. P., Anderson, D. C., Vinciguerra, T. P., He, H., Goldberg, D. L., Ehrman, S. H., Dickerson, R. R., and  
838 Salawitch, R. J.: Evaluating commercial marine emissions and their role in air quality policy using observations and the  
839 CMAQ model, *Atmospheric Environment*, 173, 96-107, <https://doi.org/10.1016/j.atmosenv.2017.10.037>, 2018.

840 Stauffer R.M., Thompson A.M., and Witte J.C.: Characterizing Global Ozoneprofile Variability from Surface to the  
841 UT/LS with a Clustering Technique and MERRA-2 Reanalysis, *J Geophys Res Atmos.* 123(11):6213-6229.  
842 <https://doi.org/10.1029/2018JD028465>, 2018.

843 Strode, S. A., Ziemke, J.R., Oman, L. D., Lamsal, L. N., Olsen, M. A., and Liu, J.: Global changes in the diurnal cycle of  
844 surface ozone, *Atmospheric Environment*, 199, 323-333, <https://doi.org/10.1016/j.atmosenv.2018.11.028>, 2019.

845 Sullivan, J. T., McGee, T. J., Leblanc, T., Sumnicht, G. K., and Twigg, L. W.: Optimization of the GSFC TROPOZ DIAL  
846 retrieval using synthetic lidar returns and ozonesondes – Part 1: Algorithm validation, *Atmos. Meas. Tech.*, 8, 4133–  
847 4143, <https://doi.org/10.5194/amt-8-4133-2015>, 2015.

848 Sullivan, J.T., McGee, T.J., DeYoung, R., Twigg, L.W., Sumnicht, G.K., Pliutau, D., Knepp, T., and Carrion, W.: Results  
849 from the NASA GSFC and LaRC Ozone Lidar intercomparison: new mobile tools for atmospheric research, *J. Atmos.*  
850 *Ocean. Technol.*, 32 (10), 1779-1795, <https://doi.org/10.1175/JTECH-D-14-00193.1>, 2015.

851 Thompson, A. M., Stauffer, R. M., Miller, S. K., Martins, D. K., Joseph, E., Weinheimer, A. J., and Diskin, G. S.: Ozone  
852 profiles in the Baltimore-Washington region (2006-2011): satellite comparisons and DISCOVER-AQ observations, *J*  
853 *Atmos Chem.*, 72(3-4), 393-422. <https://doi.org/10.1007/s10874-014-9283-z>, 2015.

854 Torgo, L.: *Data Mining with R*, Chapman & Hall/CRC Data Mining and Knowledge Discovery Series,  
855 <https://doi.org/10.1201/b10328>, 2010.

856 Wang, L., Newchurch, M. J., Alvarez, R. J., Berkoff, T. A., Brown, S. S., Carrion, W., De Young, R. J., Johnson, B. J., Ganoe,  
857 R., Gronoff, G., Kirgis, G., Kuang, S., Langford, A. O., Leblanc, T., McDuffie, E. E., McGee, T. J., Pliutau, D., Senff,

858 C. J., Sullivan, J. T., Sumnicht, G., Twigg, L. W., and Weinheimer, A. J.: Quantifying TOLNet Ozone Lidar Accuracy  
859 during the 2014 DISCOVER-AQ and FRAPPÉ Campaigns. *Atmos Meas Tech*, 10(10), 3865-3876,  
860 <http://doi.org/10.5194/amt-10-3865-2017>. 2017.

861 Wang, X., Jacob, D. J., Downs, W., Zhai, S., Zhu, L., Shah, V., Holmes, C. D., Sherwen, T., Alexander, B., Evans, M. J.,  
862 Eastham, S. D., Neuman, J. A., Veres, P. R., Koenig, T. K., Volkamer, R., Huey, L. G., Bannan, T. J., Percival, C. J.,  
863 Lee, B. H., and Thornton, J. A.: Global tropospheric halogen (Cl, Br, I) chemistry and its impact on oxidants, *Atmospheric*  
864 *Chem. Phys.*, 21, 13973–13996, <https://doi.org/10.5194/acp-21-13973-2021>, 2021.

865 Wu, Y., Nehrir, A. R., Ren, X., Dickerson, R. R., Huang, J., Stratton, P. R., Gronoff, G., Kooi, S. A., Collins, J. E., and  
866 Berkoff, T. A.: Synergistic aircraft and ground observations of transported wildfire smoke and its impact on air quality  
867 in New York City during the summer 2018 LISTOS campaign, *Science of The Total Environment*, 773,145030,  
868 <https://doi.org/10.1016/j.scitotenv.2021.145030>, 2021.