

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18

**DECODING MIDWEST JUNE PM<sub>2.5</sub> EVENTS:  
A SELF-ORGANIZING MAP APPROACH TO  
METEOROLOGICAL ANALYSIS**

Victor Geiser

Advisor: Tsengel Nergui

LADCO Intern - Summer 2024

Prepared for Lake Michigan Air Directors Consortium (LADCO)

19 **ABSTRACT:** The Midwest is a land-locked mid-latitude geographical setting where complex  
20 atmospheric processes take place in conjunction with local emissions and transported air pollutants.  
21 Periodically, upwind wildland and prescribed fire smoke is transported into the region and results in  
22 unhealthy concentrations of fine particulate matter (PM<sub>2.5</sub>) and at the surface. Comparisons of the  
23 meteorological conditions associated with typical high pollution days, versus those of fire smoke  
24 influenced days, are useful to forecasters and air quality planners.

25 To better understand the meteorological setting and pollutant transport pathways bringing fire  
26 smoke into the region, LADCO applied a spatial classification technique called a Self-Organizing  
27 Map (SOM) to daily average PM<sub>2.5</sub> concentrations using the 3-km resolution High Resolution Rapid  
28 Refresh (HRRRv4) reanalysis dataset for 2019-2023. The objective of the analysis is to identify the  
29 primary features of the physical and dynamical atmospheric conditions associated with air pollution  
30 episodes with and without the influence of smoke.

31 We will present the results of our SOM analysis of pollution episodes caused by wildland fires  
32 originating in the southwestern US and southwestern Canada. We used the SOM to identify the  
33 synoptic scale meteorological conditions and the anticipated increases in PM<sub>2.5</sub> during fire events. In  
34 addition, we investigated a key aspect of whether the long-range transported smoke aloft reached the  
35 surface. Vertical atmospheric characteristics such as wind shear, stability, and 24 changes in the  
36 geopotential height and temperature for fire-influence SOM nodes, highlight key upper-air features  
37 for vertical mixing and indicate whether air masses ascend or descend along the transport path  
38 between the fire smoke source and receptor monitors.

39 Our study offers two practical applications for air quality forecasters. First, using SOM to identify  
40 the weather patterns associated with typical high-pollution days provides historical data for similar-  
41 day analysis for exceptional event applications. Secondly, the identified synoptic weather patterns  
42 linked to fire smoke-influenced days provide insights into the expected increases in PM<sub>2.5</sub>  
43 concentrations due to fire smoke in the Midwest.

44

## 45 1. INTRODUCTION

46 Periodically, wildland fire smoke is transported into the Midwest and results in unhealthy  
47 concentrations of fine particulate matter with a diameter less than 2.5 micrometers (PM<sub>2.5</sub>). Elevated  
48 PM<sub>2.5</sub> concentrations have been associated with a wide range of human health hazards and also  
49 affect many meteorological and chemical processes in the atmosphere. Moreover, the US Midwest's  
50 central placement within the North American continent and the Great Lakes makes it a common  
51 place for a diverse range of meteorological and chemical process to occur and converge. The  
52 identification of common meteorological conditions associated with significantly above normal  
53 PM<sub>2.5</sub> concentrations is applicable to both the fields of air quality and meteorology.

54 In this study, we examine the above idea using a Self-Organizing Map (SOM). Self-Organizing Maps  
55 were originally proposed in (Kohonen 1982) and are a type of artificial neural network that aim to  
56 find lower dimensional relationships in high dimensionality data whilst preserving the original  
57 structure (topology) of its input data. Unlike its more traditional counterparts, such as principal  
58 component analysis, it makes no underlying assumptions about relationships within the input data,

59 such as linear relationships. SOMs have been applied within a wide variety of fields such as  
60 genomics (Törönen et al. 1999), astrophysics (Carrasco Kind and Brunner 2014), and economics  
61 (Deboeck and Kohonen 2013), and are commonly used tools in the fields of data mining (Vesanto  
62 and Alhoniemi 2000) and non-linear manifold learning and representation (Forest et al. 2021). More  
63 recently SOMs have been applied within the physical sciences as well and have been proven to be  
64 successful when applied to air quality and meteorology data (Hrust et al. 2009) (Hewitson and Crane  
65 2002).

66

## 67 1.1 GOALS AND MOTIVATIONS

68 This project is motivated by persistent  $PM_{2.5}$  episodes in the Great Lakes region and a need to  
69 further understand the nature of their origins and impacts. Additionally, given the importance of  
70 atmospheric aerosols on the earth's global radiation budget, insights gleaned from understanding  
71 where the pollutants occur and what effect they have on the global environment is relevant within a  
72 changing climate.

73 Leading to our overall research question:

74 **How can Self Organizing Maps (SOMs) be used to identify meso-scale meteorological**  
75 **conditions associated with high  $PM_{2.5}$  and fire smoke impacted conditions in the LADCO**  
76 **region?**

77 The goals laid out for this project are as follows:

- 78 1. To enrich the conceptual model regarding high concentrations of  $PM_{2.5}$  in the Midwest by  
79 incorporating meteorological settings identified through the Self-Organizing Map (SOM)  
80 method.
- 81 2. To establish a basis for determining whether the overhead smoke observed by satellites  
82 descended to the surface and impacted concentrations of  $PM_{2.5}$  at surface monitors.
- 83 3. Compare the synoptic weather conditions in the Midwest during air pollution episodes with  
84 and without the influence of wildfire smoke.

85

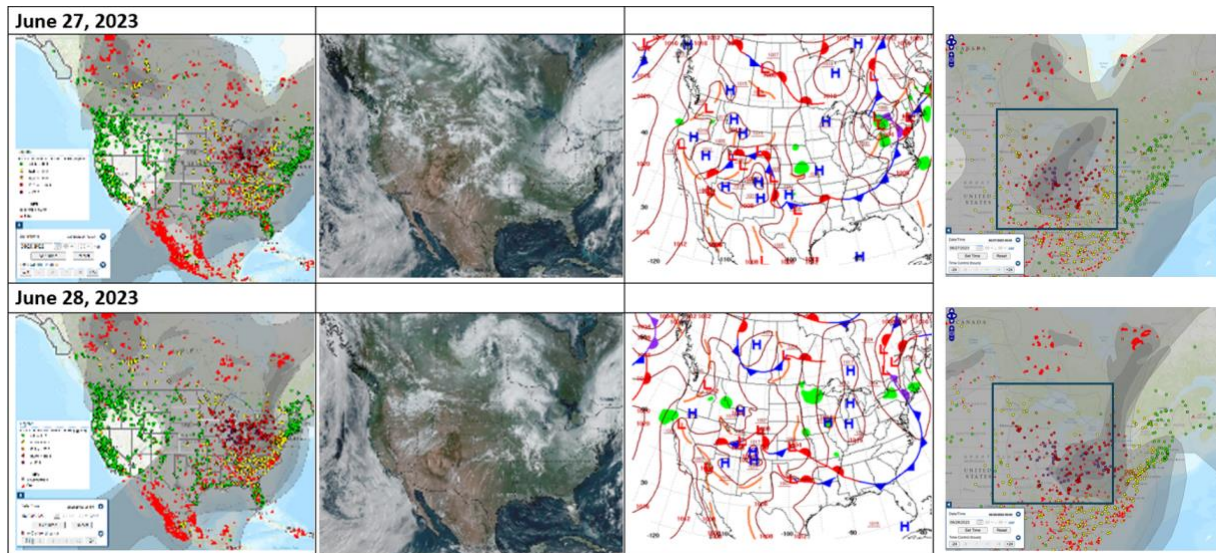
## 86 1.2 THE OBSERVATION OF RESULTS BY CONSIDERING A CASE STUDY

87 Affirmations complementing our SOM analysis can be observed clearly when following a  $PM_{2.5}$   
88 event that occurred over the LADCO region on June 25-30, 2023. The primary reason as to why  
89 this event was so anomalous was due to the impacts caused by wildfire smoke originating in Ontario  
90 and Quebec Canada. **Figure 1** illustrates the meteorological conditions during the transition  
91 between the first and second phases of the event that were dominated by an initial low-pressure  
92 system that aided in transporting polluted air into the US Midwest, followed by a high-pressure  
93 system event that led to stagnation conditions and greatly above average  $PM_{2.5}$  and impacts.

94 Although not the primary topic of this study, throughout the remainder of this report we will  
95 provide extra visual elements considering this event. This is not only to display how the results of  
96 the SOM are observable when applied in a real-world context, but also as a quick reference to

97 certain known conditions within our input data that will point to positive signatures regarding our  
 98 SOM's functionality and performance.

99



100

101 **Figure 1** The meteorological conditions surrounding June 27<sup>th</sup> and 28<sup>th</sup> 2023 during which “very  
 102 unhealthy” air quality was observed.

103

## 104 2. METHODOLOGY

105 The primary method contained within this study is the self-organizing maps algorithm itself. While  
 106 there are a multitude of implementations for self-organizing maps written in many different  
 107 programming languages, the implementation used in this study is the “MiniSOM” implementation.  
 108 MiniSOM is an open-source and purely pythonic implementation of self-organizing maps that is  
 109 available on [GitHub](https://github.com). It gets its name from “minimalistic SOM” as its only dependency is the  
 110 NumPy library, and it is generally used for small to medium sized datasets.

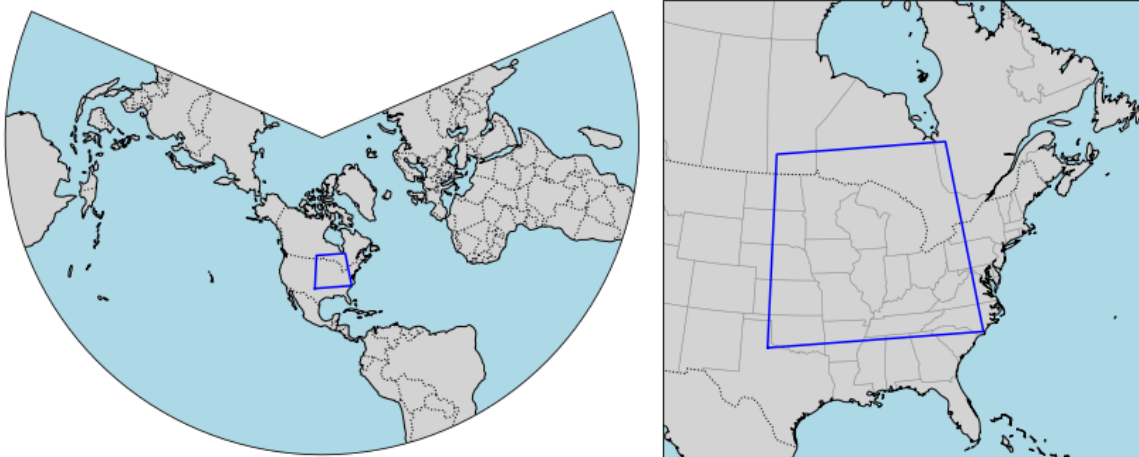
111 The self-organizing maps algorithm seeks to produce a low-dimensional (usually two-dimensional)  
 112 representation of the input space while preserving the topological properties of the original data.

113

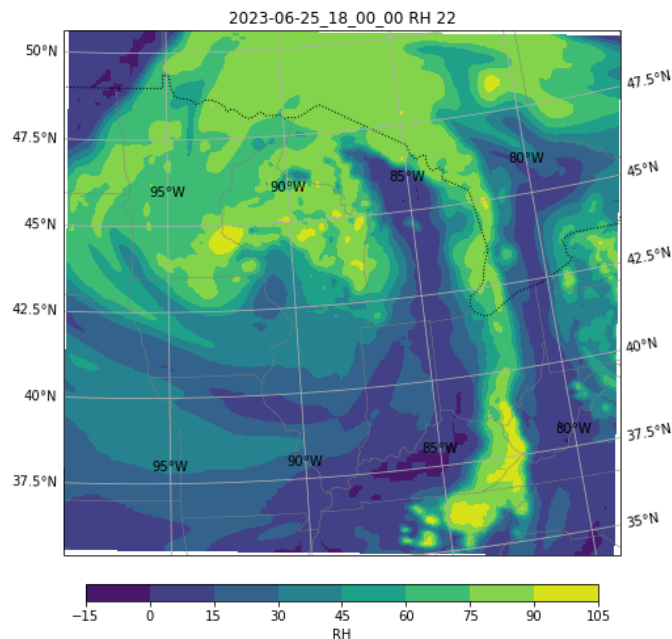
### 114 2.2 INPUT DATA

115 Meteorological data: This study contains data from a variety of sources, however the primary data  
 116 source that is used when training the SOM is daily meteorological reanalysis data, which is a blend of  
 117 the 3-km resolution HRRRv4 (High Resolution Rapid Refresh) surface reanalysis and, 12-km  
 118 resolution NAM (North American Model) reanalysis data. The dataset contains data for all June days  
 119 between 2019 and 2023. The meteorological dataset has a spatial resolution (grid spacing) of 4km  
 120 and is using the conditions at 18:00 UTC (12:00pm CST). The files were originally output in  
 121 NetCDF format and are read into python through use of the Xarray package in python. Each

122 meteorological file also contained the necessary projection information that allowed the data to be  
 123 plotted on a 420 latitude by 444 longitude extent on a Lambert Conformal Conic projection. **Figure**  
 124 **2** displays the extent of our data ranges from a latitude and longitude of (34.163, -100.316) in the SW  
 125 corner to (50.644, -78.027) in the NE corner. **Figure 3** provides an example visualization of one of  
 126 our data variables, relative humidity at the 500hPa level.



127  
 128 **Figure 2** The (LADCO) region of interest for this study on a Lambert Conformal Conic projection.  
 129



130  
 131 **Figure 3** 500hPa Relative Humidity (%) for 06-25-2023 at 18:00 UTC.

132  
 133 Although the meteorological data used for this study contains over 115 variables, only the variables  
 134 that are used as inputs into the SOM will be discussed in this report. Mentioned here briefly are the



135 names of these variables, their available vertical levels (model levels), their units, and associated  
 136 abbreviations within the dataset. Motivations for why these variables were selected for SOM analysis  
 137 can be found in section 4.

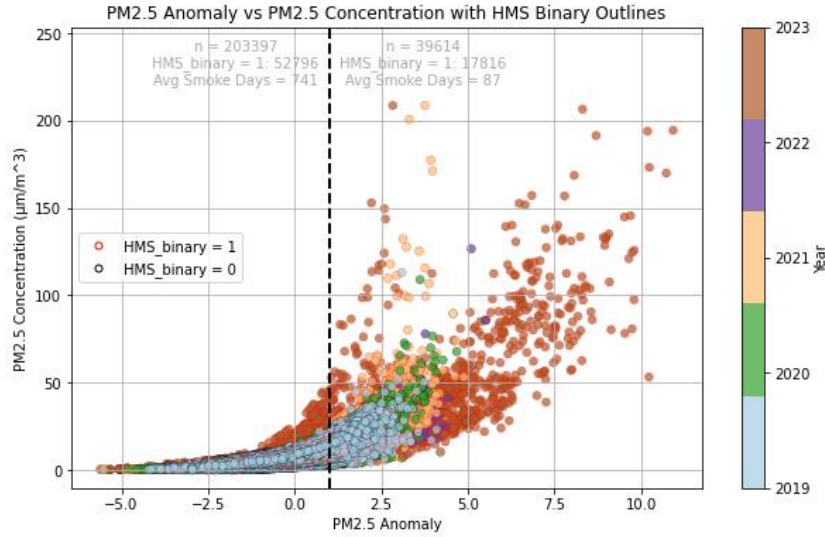
- 138 1. “PMSL” – Pressure at mean sea level (surface only). Units: Pascals (Pa)
- 139 2. “RH” – Relative humidity (model levels 1-40). Units: Percentage (%)
- 140 3. “TT” – Temperature (model levels 1-40). Units: Degrees Kelvin (K)
- 141 4. “UU” – Horizontal “U” wind component (model levels 1-40). Units: meters per second
- 142 (m/s)
- 143 5. “VV” – Vertical “V” wind component (model levels 1-40). Units: meters per second (m/s)
- 144 6. “GHT” – Geopotential height (model levels 1-40). Units: meters (m)

145 Air quality data: In addition to the HRRR meteorological data, two other datasets were used for  
 146 analysis purposes. The first is a tabular dataset containing observed  $PM_{2.5}$  and data from the US EPA  
 147 Air Quality System (AQS). We calculated additional SOM node metrics based off four columns  
 148 contained within this dataset:

- 149 1. “value” – An observed  $PM_{2.5}$  concentration in  $\mu g/m^3$ .
- 150 2. “std\_log\_value” (or  $PM_{2.5}$  “anomaly”) – A standardized value of the measured  $PM_{2.5}$   
 151 concentration. Standardization (i.e., normalization) was done using the monthly mean and  
 152 standard deviation of the log-transformed measured values at a monitor over the 2019-2022  
 153 period. This standardized value (i.e., anomaly) provides a measure for how much  $PM_{2.5}$   
 154 concentration deviates from its typical mean.
- 155 3. “HMS\_binary” – A binary flag variable (either 0 or 1) that determined if overhead smoke  
 156 was identified at the location of a monitor through a satellite-driven product called the  
 157 Hazard Mapping System (HMS).
- 158 4. “res\_1sigma\_std\_log\_value” (or “res1”) – The residual value of  $PM_{2.5}$  concentrations above  
 159 and below 1 standard deviation, high indicates how much the measure value was beyond the  
 160 typically observed values a monitor.

161 **Figure 4** shows a visualization of the above data variables for the  $PM_{2.5}$  dataset: “value” (on the y  
 162 axis) and “std\_log\_value” (on the x axis) with HMS\_binary outlines. Values for  
 163 “res\_1sigma\_std\_log\_value” would then be the data on the right side of the vertical dashed black  
 164 line and with red outlines.

165

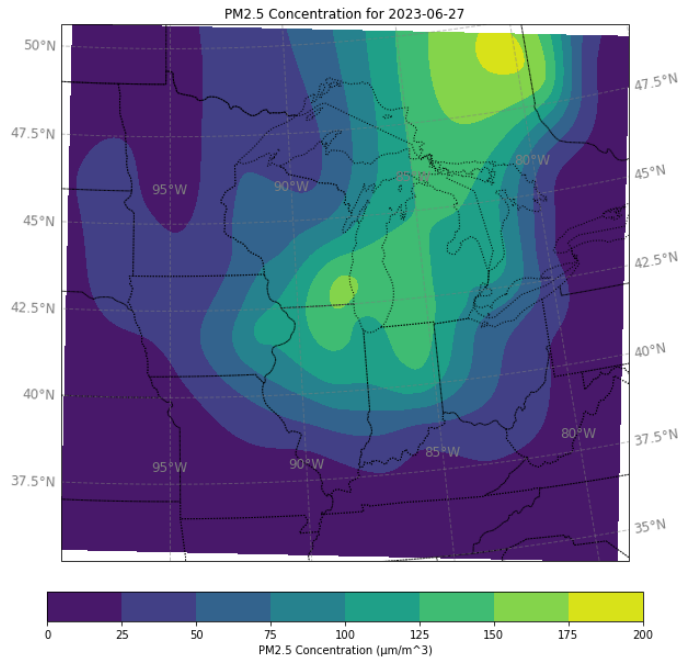


166

167 **Figure 4** Scatter plot of PM<sub>2.5</sub> concentration vs. its standardized anomaly with outlines for overhead  
 168 smoke.

169 The last dataset that we used during the final stages of this project is a “krigged” (spatially  
 170 interpolated) PM<sub>2.5</sub> dataset. This dataset is an interpolated product based off the PM<sub>2.5</sub> ground sensor  
 171 network. **Figure 5** is an example visualization of the krigged PM<sub>2.5</sub> dataset.

172



173

174 **Figure 5** Krigged PM<sub>2.5</sub> field for 06-27-2023 displaying PM<sub>2.5</sub> transport into the LADCO region.

175 Although mentioned here for completeness, the krigged PM<sub>2.5</sub> dataset is not used until **section 5**.

176

177 To ensure compatibility with MiniSOM two preprocessing steps needed to be applied to the data:

- 178 1. Vectorization – Due to the spatial nature of the data they needed to be vectorized in order  
179 to be input into MiniSOM.
- 180 2. Standard Scaling – Standard scaling (through Sci-kit learn) is a technique that scales the data  
181 between (-1 and 1) where each variable has a mean of “0” and a standard deviation of “1”

182 Vectorization adds to the dimensionality of our data substantially and, as we will see in our analysis,  
183 this will have lasting effects as far as our quantitative metrics and furthermore in our interpretive  
184 analysis. However, because no clear alternatives to vectorization currently exist, and vectorizing our  
185 data still allows for spatial patterns in our data to be represented, this is the standard approach. This  
186 raises the question: why not use an initial dimensionality reduction technique when applying  
187 preprocessing steps? The answer to this questions lies in the interpretive need to recreate and  
188 visualize our data. Further study could potentially look at applying techniques such as t-distributed  
189 Stochastic Neighbor Embedding or Uniform Manifold Approximation and Projection as a further  
190 preprocessing step, given the nature of the non-linear relationships we are attempting to explore,  
191 however this would make the final visualizations produced from the SOM less meaningful.

192 As a final data preprocessing step, a standardized scaler was applied to the meteorological data as  
193 preparation for input into the SOM. This study used the StandardScaler method built into Sci-Kit  
194 Learn, which scales the data to a range between -1 and 1 and a mean of 0. To account for the  
195 varying scales of the meteorological data, a standard scaler was applied individually to all variables.

196

## 197 2.2 DESCRIPTION OF THE SELF ORGANIZING MAPS ALGORITHM

198 Although descriptions of the Self Organizing Maps algorithm exist across many sources within the  
199 literature, a brief summary adapted from (Kohonen 1982) and (Hulle and Marc 2012) will be  
200 presented here.

201 Step 1: Initialization

202 The first step within the SOM algorithm happens when a grid of neurons (also called nodes) is  
203 initialized. Each neuron has a weight vector of the same dimensionality as the input data.

204 Step 2: Training algorithm

205 The self-organizing maps training algorithm has two main components:

- 206 1. The best matching unit (BMU)

207 The BMU is the neuron whose weight vector is closest to the input vector in terms of Euclidean  
208 distance. This can be mathematically expressed as:

$$209 \quad c = \operatorname{argmin}_i \| \mathbf{x}(t) - \mathbf{w}_i(t) \|$$

210 Where:

- 211 •  $c$  is the index of the BMU.





242 
$$\rho(t) = \frac{\rho_0}{1 + \frac{t}{\sigma}}$$

243 Where:

- 244 •  $\rho_0$  is the initial neighborhood radius
- 245 •  $t$  is the current iteration number
- 246 •  $\sigma$  is the time constraint that controls the rate of decay

247

## 248 2.3 SOM HYPERPARAMETERS AND CONFIGURATION

249 Within the primary Self-Organizing Maps algorithm established above, there also different  
250 configurations that can be achieved by tweaking a SOM's hyperparameters. The hyperparameters for  
251 the LADCO SOM are as follows:

- 252 • `som_size = (3,5)`
- 253 • `sigma = 1`
- 254 • `learning_rate = .3`
- 255 • `ngb_function = 'gaussian'`
- 256 • `decay_function = 'linear_decay_to_zero'`
- 257 • `sigma_decay_function = 'asymptotic_decay'`
- 258 • `init = 'random'`
- 259 • `train = 'random'`
- 260 • `iterations = 200`
- 261 • `topology = 'hexagonal'`
- 262 • `activation_distance = 'euclidean'`
- 263 • `random_state = '64'`

264 Notable deviations from the default parameters include a hexagonal topology which allows our  
265 nodes to have more neighbors as opposed to the default rectangular topology. We have a slightly  
266 lower than normal learning rate that resulted in better performance via iterative experimentation and  
267 our learning rate decay function is set to decrease linearly as opposed to the standard asymptotic  
268 decay which resulted in better clustering via the clustering metrics as described in **section 2.4**. We  
269 set the number of iterations at 200 because more iterations did not result in significantly improved  
270 performance. It should be noted that, since our learning curve visualization in **Figure 8** uses  
271 quantization error, and since our data are highly dimensional, this curve appears in a slightly atypical  
272 fashion as opposed SOMs that may occur elsewhere within the literature. More about LADCO  
273 SOM's learning curve will be discussed in **section 3**.

274 Most importantly within the topic of SOM hyperparameters is the SOM size, which controls the  
275 number of output neurons within the SOM. The determination of this hyperparameter is often  
276 crucial to the functionality of the SOM and is often a trade-off involving capturing more general  
277 trends, sensitivity to outliers, and having enough nodes to capture the more nuanced and

278 informative trends within the data (Hulle and Marc 2012). In the case of the LADCO SOM our  
279 SOM size appears to be limited primarily by the number of samples currently ingestible withing the  
280 workflow. We consider all June days between 2019 and 2023 leaving us with 149 samples (1 day of  
281 the reanalysis dataset is missing to generate due to an incomplete HRRR run for that time period). If  
282 SOM size increases in an attempt to capture harder to detect relationships within the data, we begin  
283 to observe nodes that have an activation response (or the number of samples from the input data  
284 that get classified as having that pattern, or activated that particular node during the training process)  
285 of 0. Due to this, determination the optimal SOM size for the LADCO SOM is an area for potential  
286 enhancement. However, if different climatological periods are considered or perhaps expanded  
287 upon in the future this may come naturally given the current implementation.

288

## 289 2.4 SUMMARY OF SOM AND NODE METRICS

290 In addition to the primary output of this study, which comes in the form of a visualization of the  
291 weights of LADCO SOM itself, there will also occur above or below each node, secondary node  
292 statistics calculated from averaged variables for all nodes that are included in the activation response  
293 for a particular node. By running each input vector through the SOM after the training period has  
294 completed, we are able to generate a map of which input vectors are considered to “match” that  
295 output node.

296 The secondary parameters visualized alongside the weights for each variable node are:

- 297 1. “Node (x,y)” – The node’s position within the SOM hexagonal grid
- 298 2. “n = ...” – The number of samples mapped to that particular node
- 299 3. “Smoke Days” – The number of identified days that met the condition mentioned in the  
300 “Res1” variable explanation
- 301 4. “Avg PM” – Average measured PM<sub>2.5</sub> concentrations over all monitors within the domain  
302 for days classified for a particular node
- 303 5. “Avg PM anom” – Similar to the Avg PM variable, but for standardized anomalies  
304 (std\_log\_value variable).
- 305 6. “Avg Res1PM” – A node average of the “res1” variable

## 306 2.5 VERTIAL PROFILE GENERATION PROCEDURE

307 Since our meteorological data are model derived and all vertical model levels are present within the  
308 meteorological data used for this study, we are able to produce an additional secondary (or tertiary)  
309 node analysis in the form of a visualization of a node’s averaged vertical atmospheric profile. The  
310 profiles are point soundings in one location (although an extended explanation about developing the  
311 functionality further will occur in **section 5**) for each model level temperature, relative humidity, and  
312 u and v wind vector components. Visualization of the profile is handled by the MetPy library in  
313 Python.

314 With these visualizations of the vertical profile, we also display calculated environmental statistics  
315 based on the profiles generated for each node.

316 These statistics include:

- 317 1. “Node (x,y)” – Same as in **section 2.4**  
 318 2. “n = ...” – Same as in **section 2.4**  
 319 3. “SH01” – The 0-1km environmental shear vector  
 320 4. “SH06” – The 0-6km environmental shear vector  
 321 5. “850hPa  $\omega$ ” – Vertical velocity at 850hPa  
 322 6. “CIN” – Convective Inhibition  
 323 7. “CAPE” – Convective Available Potential Energy  
 324 8. “#TI” – Number of temperature inversions (2°C / 100hPa)  
 325 9. “850-950hPa avg temp” – Average temperature difference between the 850 hPa and 950hPa  
 326 levels  
 327 10. “Avg PM2.5 Anomaly” – Same as in **section 2.4**

328

### 329 3. RESULTS

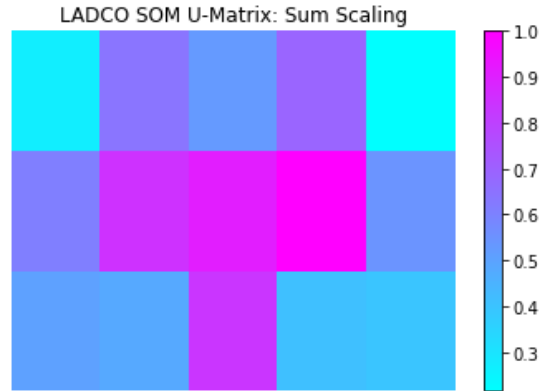
330 This chapter includes the following three sections:

- 331 • **Section 3.1** includes SOM diagnostic plots that support the primary analysis.  
 332 • **Section 3.2** covers the primary results of the study, the visualization of the weights of  
 333 LADCO SOM, and the conclusions that can be reached as a result.  
 334 • **Section 3.3** will present the vertical profile results described in **section 2.5**

335

#### 336 3.1 SOM DIAGNOSTIC PLOTS

337 Presented in **Figure 6** and **Figure 7** are the LADCO SOM distance matrix (or u-matrix) with sum  
 338 scaling and mean scaling respectively. The distance matrix is used to measure the distances between  
 339 the nodes in a SOM grid. This distance can be scaled in various ways, two of which are a mean  
 340 (average) scaled distance matrix and a sum scaled distance matrix. The sum scaled distance matrix,  
 341 visualized in **Figure 6**, represents the total sum of distances between (the vector values of) a  
 342 particular node and all other nodes in the SOM. The sum scaled distance matrix indicates the overall  
 343 quality and separation of clusters in the SOM and it can also be used to inform where potential  
 344 boundaries between clusters appear within the SOM. All weight visualizations in section 3.2 will use  
 345 the mean scaled distance matrix as a background color, or “frame color” for reference.



346

347

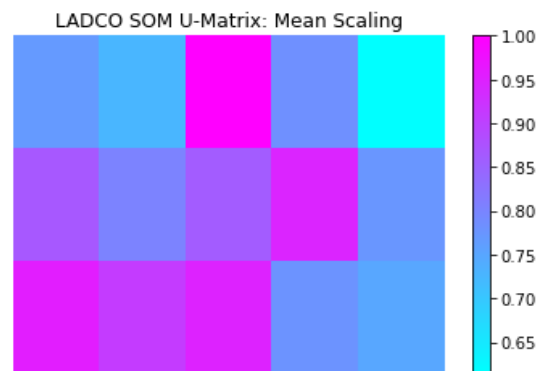
**Figure 6** LADCO SOM distance matrix (sum scaled)

348

349 The mean scaled distance matrix, visualized in **Figure 7**, is a representation of the average distance  
 350 between the prototype vectors of nodes in the SOM grid. This figure provides a measure of how  
 351 smoothly the input space is represented by the SOM and it can visualize the average separation  
 352 between clusters. Nodes with a high value within the mean scaled distance matrix represent nodes  
 353 that are on average further apart from their surrounding neighbors. Higher values in this figure  
 354 indicate nodes pattern that are significantly different in structure (in our case meteorological  
 355 conditions) than its surrounding neighbors.

356

357 The middle sections of the LADCO SOM distance matrix have a higher total distance.  
 358 Unfortunately, the high dimensionality of our data appears to affect the summed distance matrix  
 359 quite a bit. The middle sections of the SOM (remember the hexagonal topology) with the more  
 360 pinkish and purple values inform us that our central nodes appear to sit along a boundary between  
 361 clusters (a result that is also apparent throughout the project). While we can read this result from the  
 362 summed distance matrix, given the relatively small size of our SOM in general, it could have been  
 363 inferred that nodes that are closer in weight values to one another would occur in a region of the  
 364 SOM where we expect to see more intermediate meteorological regimes.

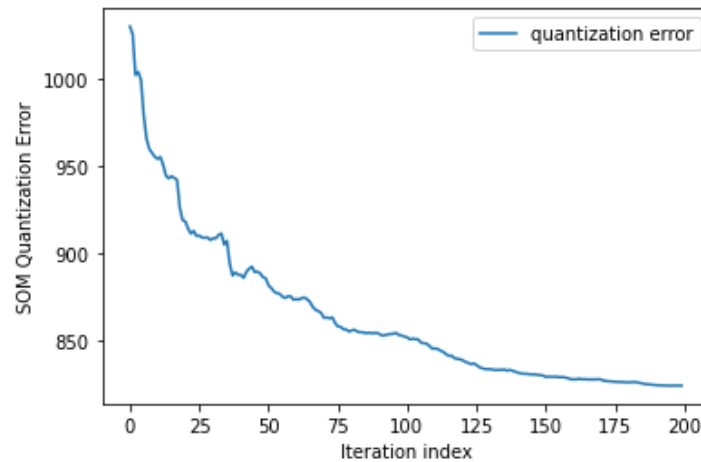


365

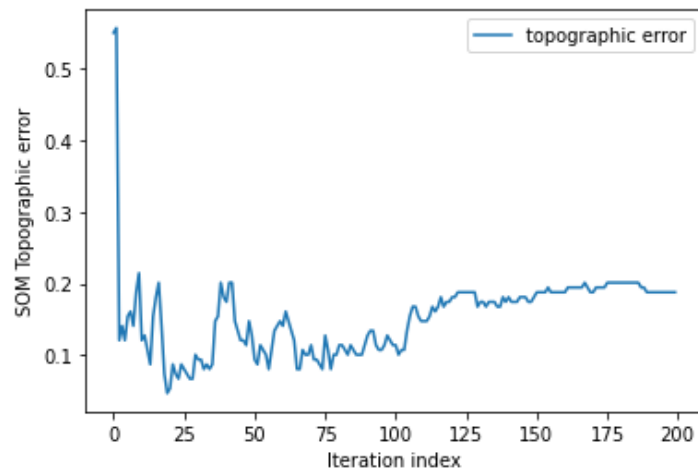
366

**Figure 7** LADCO SOM distance matrix (mean scaled)

367 As much as the sum scaled distance matrix is informative (although predictable), the mean scaled  
 368 matrix is even more so, especially in light of our highly dimensional data. **Figure 7**'s purpose is  
 369 threefold, it displays nodes that are comparatively different in relation to their neighbors (higher  
 370 average separation), it displays the relative smoothness of our SOM, and it illustrates how well the  
 371 topology of the original data is being preserved within the SOM. Interpretable from **Figure 7** is a  
 372 primary cluster of nodes with higher average separation in the lower left hand corner and  
 373 throughout the middle of the SOM. With these nodes all having values in the upper ranges for mean  
 374 scaled distance we can understand that (A) this local region of the SOM contains a wide range of  
 375 types of cases (meteorological setups representing the placement and orientation of high pressure in  
 376 this case) that are distinct from one another, and (B) looking globally at our entire mean scaled  
 377 distance matrix, we can see that although topological relationships are being preserved, there still  
 378 exist regions of the SOM that are better than others. Visualizations of the quantization error (QE)  
 379 and topographic error (TE) learning curve for the LADCO SOM in **Figure 8** and **Figure 9** will  
 380 examine possible reasons and justifications for this observation.



381  
 382 **Figure 8** Learning curve showing quantization error decreasing with 200 iterations



383  
 384 **Figure 9** Learning curve showing topographic error decreasing but then increasing slightly with  
 385 number of iterations



386 A very simple observation that can be reached from briefly examining the learning curves for the  
 387 LADCO SOM is that both the QE and TE learning curves display atypical behavior. For the QE  
 388 curve our final QE error is 824, which is much greater than the close to zero value normally  
 389 observed for typical SOM applications. This is explainable keeping in mind how QE is calculated in  
 390 the first place and considering the dimensionality of our data.

391 Quantization error is calculated by the formula:

$$392 \quad QE = \frac{1}{N} \sum_{i=1}^N \| \mathbf{x}_i - \mathbf{w}_{BMU(i)} \|$$

393 Where:

- 394 •  $N$  is the number of samples.
- 395 •  $\mathbf{x}_i$  is the  $i$ -th sample in the dataset.
- 396 •  $\mathbf{w}_{BMU(i)}$  is the weight vector of the Best Matching Unit (BMU) for the  $i$ -th sample.
- 397 •  $\|\cdot\|$  denotes the Euclidean distance.

398 As noted in the quantization error formula, QE is primarily a Euclidean distance measure. For the  
 399 LADCO SOM the QE converges to 824 as a side effect of our data’s dimensionality where our  
 400 SOM is actually performing quite well, however because each of our input variables has a vectorized  
 401 length of 186,480, combined with the fact that one sample in the dataset has 6 variables, each  
 402 sample has a vectorized length of 1,118,880. Considering our high dimensionality, this means that  
 403 even small residuals between an input vector and its matching BMU are propagated in relation to the  
 404 data’s dimensionality, when in reality, a QE of 824 means that  $824/1,118,880 = 0.00073\dots$  the  
 405 average error of individual elements within the input array and its BMU is comparatively very small.

406 Turning our attention to the topographic error, we notice an unusual trend by iteration 25, in that  
 407 our TE begins to increase and then level out with increasing iterations. Topographic error is  
 408 calculated by finding the first BMU and the second BMU, and a sample for which these two nodes  
 409 are not adjacent counts as an error. The topographic error is given by the total number of errors  
 410 divided by the total number of samples. A similar trend was observed in Forest et al. (2021) where  
 411 the phenomenon of increasing TE is correctly explained: “Topographic error shows the trade-off  
 412 between self-organization ... and the resulting clustering quality” who further went on to mention  
 413 how “A practitioner could thus choose to use an early stopping strategy ... but it would harm the  
 414 quality of the clustering.” Essentially, the increase in TE of the LADCO SOM is related to an  
 415 increase in clustering quality, where, by to some extent ignoring the data’s topological relationships,  
 416 better clustering can be achieved.

417 While this result may initially be concerning, given the day-to-day variability in mesoscale  
 418 meteorological patterns, it is expected that any given sample may not be similar enough to its second  
 419 BMU to count as a topographic error, the outcome of which can be observed in both the mean  
 420 scaled distance matrix in **Figure 7** and the explicit TE learning curve in **Figure 9**. Furthermore, the  
 421 phenomenon of increasing TE may also be in part due to the high dimensionality of the input data  
 422 diminishing the overall utility and representation of Euclidian distances in our data space, as seen by  
 423 the QE learning curve in **Figure 8**.

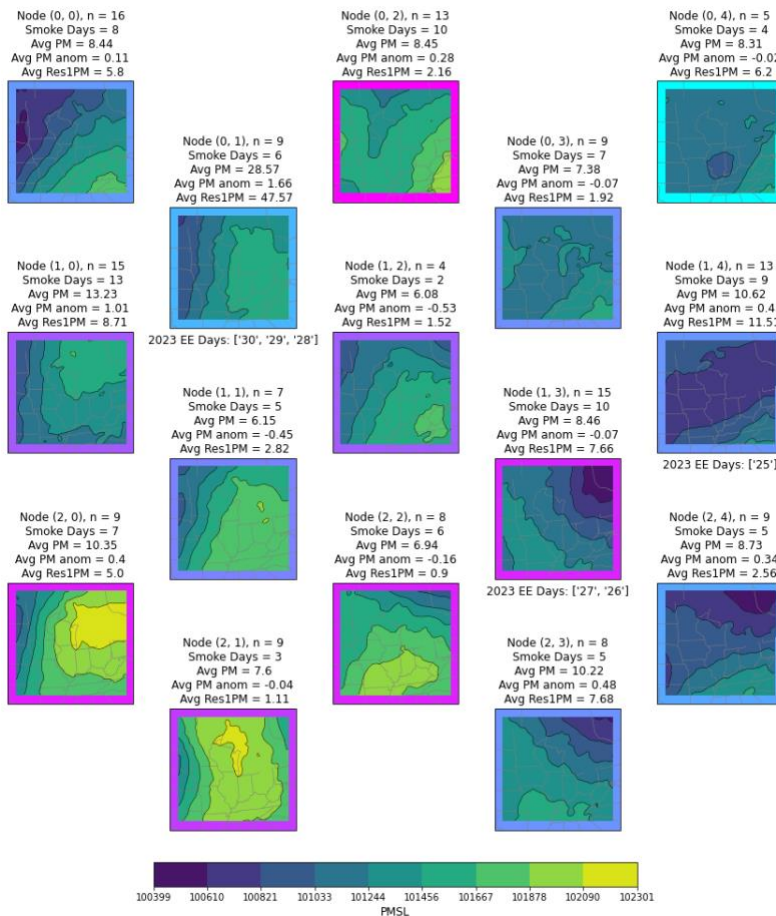
424

### 425 3.2 LADCO SOM WEIGHT VISUALIZATIONS

426 The primary results from the LADCO SOM present what a classification of meteorological regimes  
 427 looks like to a self-organizing map. The resulting clusters (nodes) are then compared using the  
 428 metrics presented in **section 2.4**. This section will serve primarily to introduce the weight  
 429 visualizations, where further discussion is prompted based on these results in **section 4**. Results will  
 430 be presented in the order they were introduced in **section 2.2**. Sections 3.2.1 through 3.2.5 will  
 431 cover the results of the purely meteorological LADCO SOM, in which other variations of this  
 432 primary LADCO SOM being introduced later.

#### 433 3.2.1 Variable: Mean Sea Level Pressure (surface level)

#### Weights for PMSL Level None



434

435 **Figure 10** The weights for Mean Sea Level Pressure within LADCO SOM in Pascals.

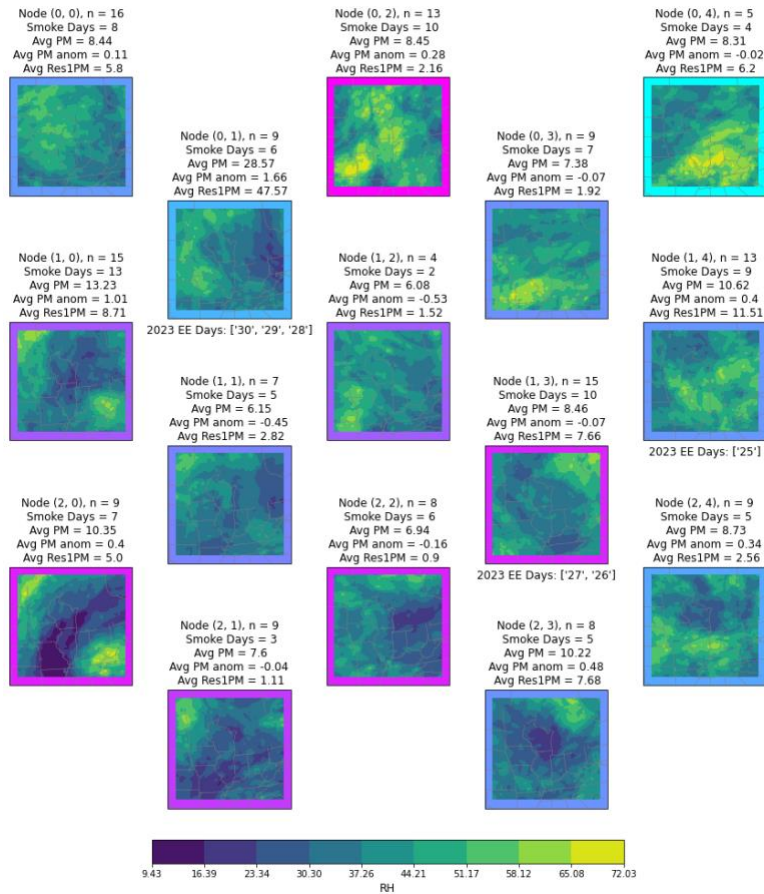
436 Mean Sea Level Pressure (MSLP) is primarily characterized by either high pressure (in the lower left  
 437 corner) or low pressure (in the lower right corner). In between transition states with no dominant  
 438 pressure pattern occurring within the middle of the SOM and in the upper right and left most  
 439 corners. The MSLP weights present an overall view of conditions at the surface and will be referred

440 to frequently through the remainder of this report. Node (0,1) is dominated by a weak high-pressure  
 441 pattern and node (1,3) is dominated by a strong northeastern low-pressure system. As evident by the  
 442 caption below, each node presents two very different meteorological pressure patterns that occurred  
 443 within the 2023 Canadian wildfire event introduced in **section 1**, which will hereby be referred to as  
 444 the 2023 EE (exceptional event).

445 The MSLP weights for within LADCO SOM suggest that higher PM<sub>2.5</sub> anomalies can be expected  
 446 during high-pressure patterns less than 1016hPa. These nodes are associated with stagnation  
 447 conditions and less dynamic motion within the atmosphere, and diminished advective processes  
 448 transporting emissions or wildfire smoke out of the LADCO region. The MSLP field in nodes (2,3)  
 449 and (2,4) describe a situation where the LADCO region is in between a high pressure system to the  
 450 south (with presumably anticyclonic motion) and a low pressure system to the north (with  
 451 presumably cyclonic motion). These two flows will enhance transport within the region which  
 452 indicates that in low pressure dominated nodes, environments conducive to westerly transport are  
 453 associated with stronger transport into the region and elevated PM<sub>2.5</sub> impacts occur as a result.  
 454 MSLP was chosen as an input variable as MSLP is one of the most recognizable patterns in  
 455 forecasting, and it is a surface variable that is normalized to account for the Appalachian Mountains.

456 3.2.2 Variable: 500 hPa level Relative Humidity

Weights for RH Level 22



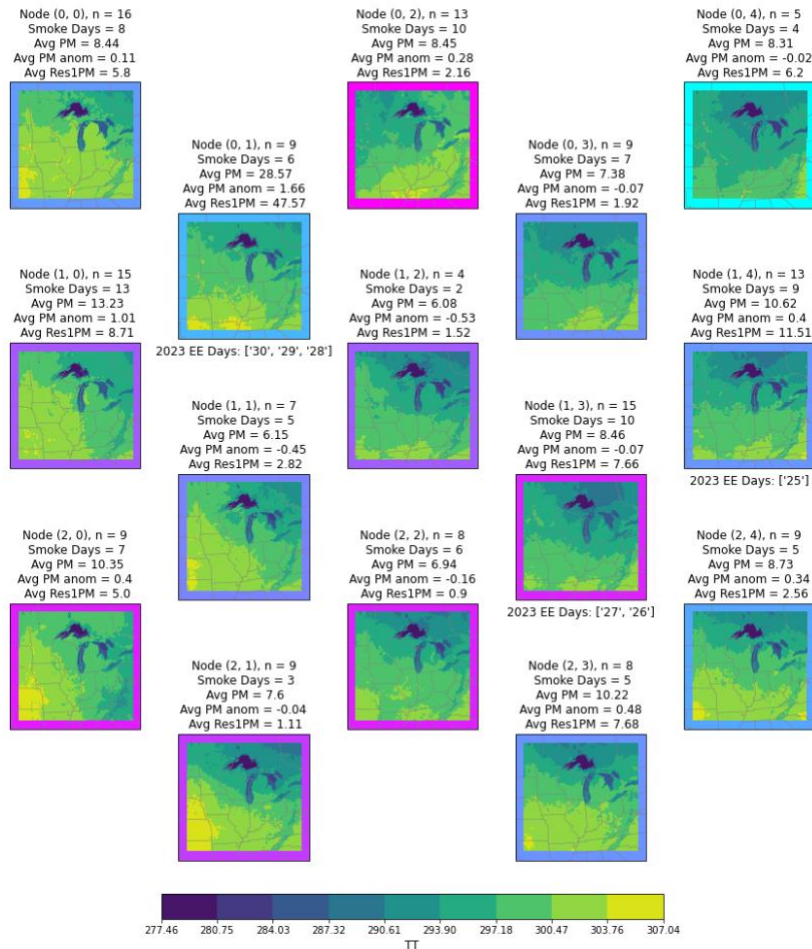
458 **Figure 11** The weights for 500 hPa relative humidity within LADCO SOM in percentage (%).

459

460 The 500hPa relative humidity (RH) is a weak predictor for the LADCO SOM for June. Taking the  
 461 average RH value for a node and comparing it to variation in the PM<sub>2.5</sub> field yields a Spearman  
 462 Correlation of -0.067 and a P-value of 0.81. Despite RH being a weak predictor, some notable  
 463 trends are still interpretable from the RH variable. Mainly within our high pressure dominated nodes  
 464 we see definitive dry streaks at the mid-levels, and in the opposite corner we also see sharp moisture  
 465 gradients within nodes that have near 0 average PM<sub>2.5</sub> anomaly. This may have not been noticeable  
 466 from the MSLP weights as there is most likely some dilution of the field due to averaging within the  
 467 pressure field, but these sharp RH gradients may be indicative of higher cloud cover over the  
 468 LADCO region which would be associated in this case with frontal passage and storms, in turn  
 469 leading to increase wet deposition and lower PM<sub>2.5</sub> anomaly. The choice to include relative humidity  
 470 at the 500hPa level is motivated by variations in the mid-level moisture profile, variations that can  
 471 become plainly visible in the vertical profile plots presented in **section 3.3**.

472 3.2.3 Variable: Surface Level Temperature

Weights for TT Level 1



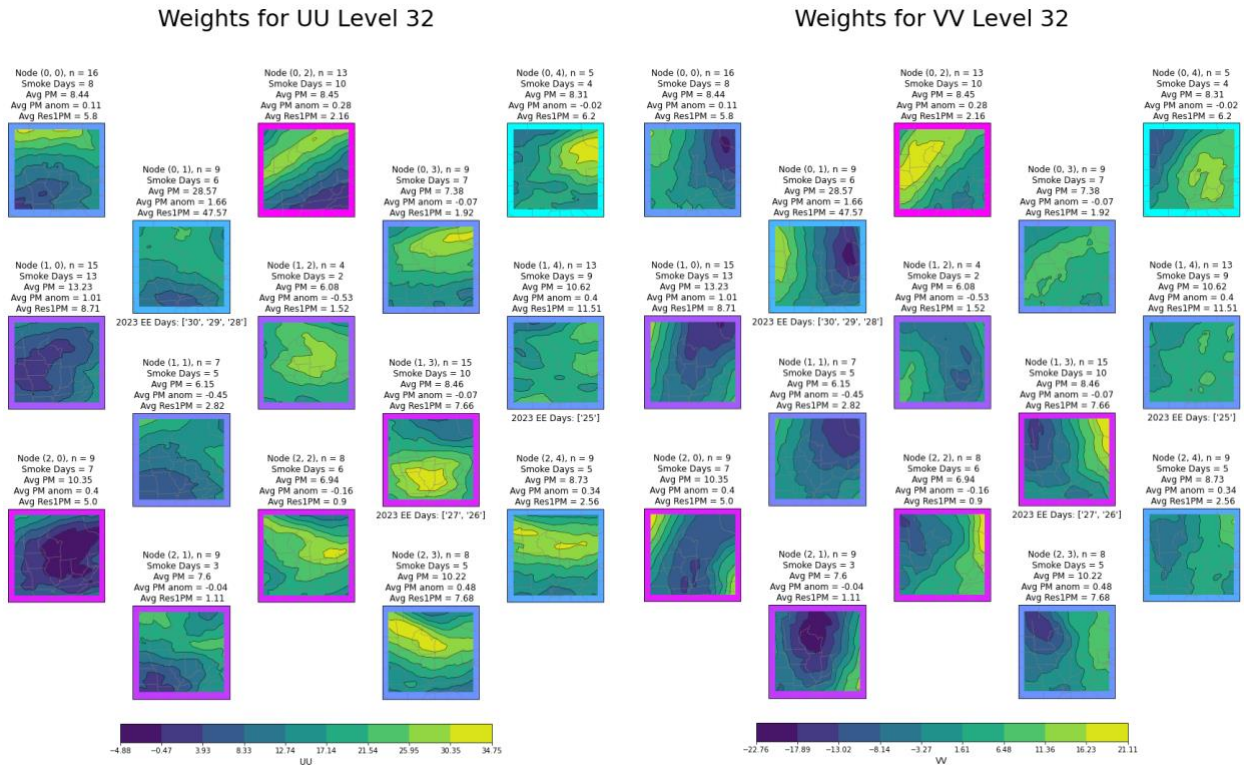
473



474 **Figure 12** The weights for surface temperature within LADCO SOM in Kelvin (K).

475  
 476 Two primary surface temperature patterns emerge based on **Figure 12**. The left side of the SOM is  
 477 characterized by warmer southern temperatures extending northward, and the right side is  
 478 characterized by cooler northern temperatures extending southward. The statistical relationship  
 479 between node averaged surface temperatures and PM<sub>2.5</sub> concentration is slightly stronger with a  
 480 Spearman correlation 0.44 and a p-value of 0.099 indicating the relationship is slightly positive  
 481 (higher temperatures correlate with higher PM<sub>2.5</sub> anomaly) and it is statistically significant at the 10%  
 482 level. In addition, standard temperature tends where colder temperatures appear northward, and  
 483 warmer temperatures occur southward. Nodes such as (0,2), which present a more unique  
 484 temperature setup with a conveyor belt of warm air extending as far north as southern Michigan,  
 485 associated with anticyclonic motion from the south, and node (0,4), which has a protrusion of colder  
 486 air extending into Missouri, a trend consistent with previous associations of node (0,4) with the  
 487 fronts passing over the LADCO region. Although low variability surface temperature makes trends  
 488 within LADCO SOM harder to visualize, it was selected as input variable because surface  
 489 temperature is one of the most commonly measured parameters in meteorology, and both PM<sub>2.5</sub>  
 490 impacts, and human impacts can be better understood considering it.

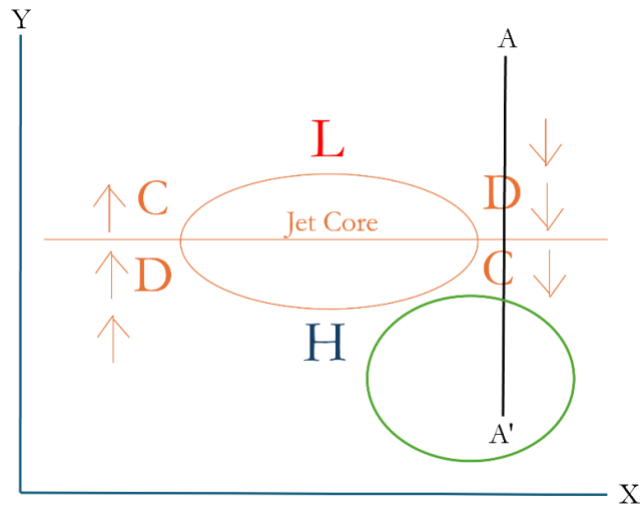
491 **3.2.4 Variable: 250hPa level U and V wind components**



492  
 493 **Figures 13a-b** The weights for U and V wind components within LADCO SOM in meters per  
 494 second (m/s).

495 The U **Figure 13a** and V **Figure13b** wind vectors are incorporated into the SOM as separate  
 496 variables, however, to improve readability these are commonly combined into the total wind speed  
 497 magnitude as seen in **Figure 16**. The 250hPa level is informative when diagnosing warm season jet  
 498 streak patterns (as opposed to the more traditional 300hPa level in the cool season). **Figure 14** and  
 499 **Figure 15** present a very abbreviated summary of some primary concepts of jet streak motion and  
 500 dynamics from figures adapted from (Keyser and Shapiro 1986).

501

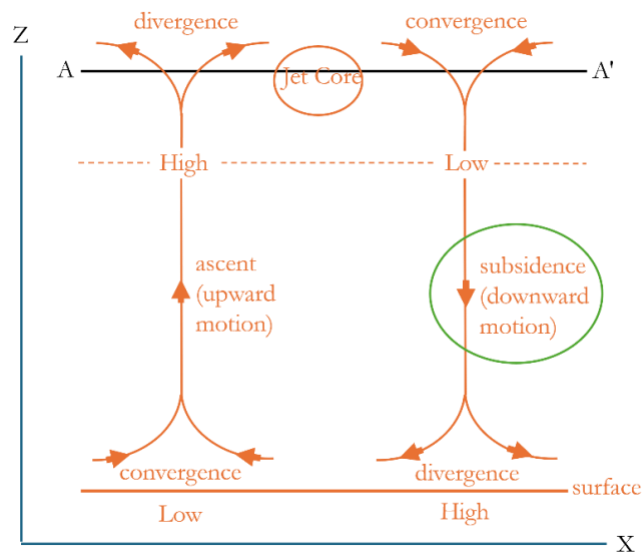


502

503

**Figure 14** A schematic of jet core dynamics in horizontal plane.

504



505

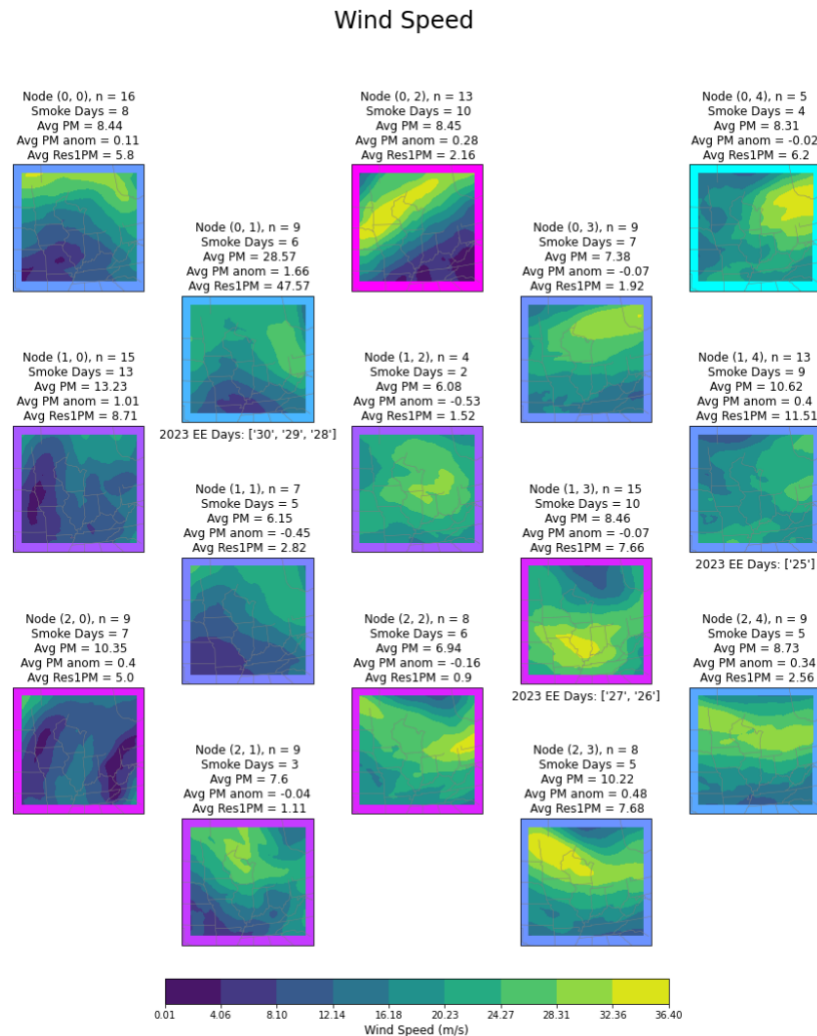
506

**Figure 15** A schematic of jet core dynamics in a vertical cross-section view.

507



508 To summarize, quadrants of the jet core correspond with either upper-level convergence or upper-  
 509 level divergence, which themselves are associated with vertical ascent or subsidence within the  
 510 atmosphere. Our interpretation of **Figure 14** will rely on knowledge of these concepts.



511  
 512 **Figure 16** 250hPa level wind speed field derived from LADCO SOM in meters per second (m/s).  
 513

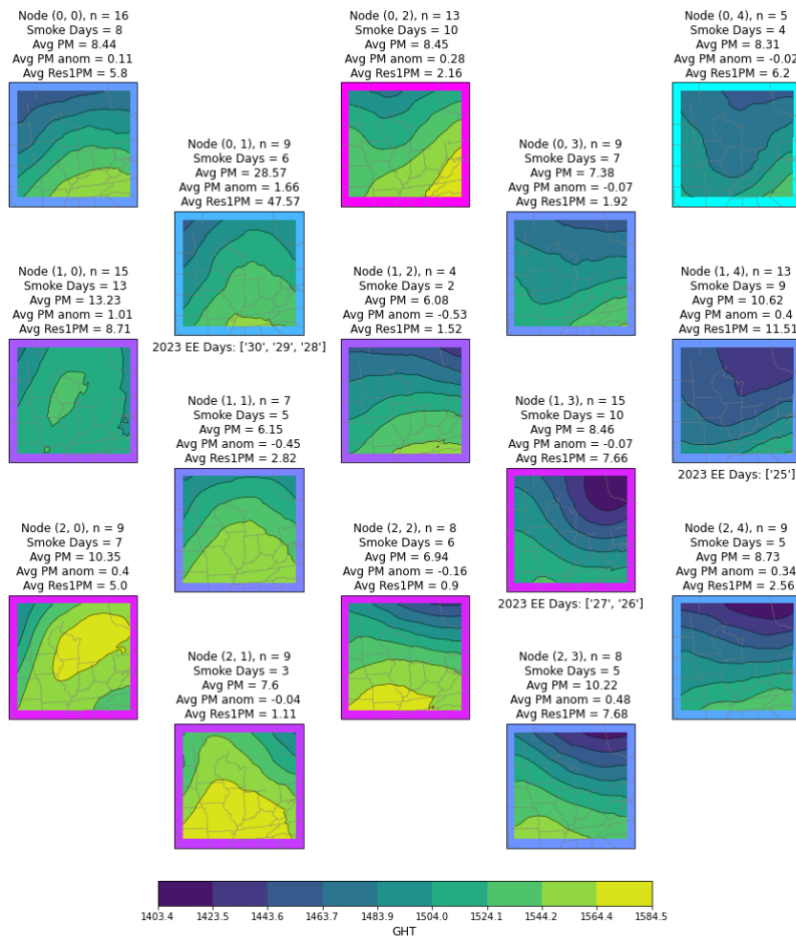
514 Analyzing the wind speed field from the LADCO SOM yields several conclusions. Among them, is  
 515 evidence that upper-level convergence (atmospheric subsidence) and downward motion within the  
 516 atmosphere is associated with higher PM<sub>2.5</sub> anomaly; however, the inverse is not explicitly true.  
 517 Upper-level divergence is generally a feature found within the LADCO region in conjunction with  
 518 severe storms and deepening mid-latitude low pressure systems. While not necessarily for severe  
 519 storms, these general quadrants of the jet streak are often used as forecasting tools that can inform  
 520 where areas of more severe weather are expected. LADCO SOM adds an additional layer of  
 521 understanding and importance to analyzing the jet streak layer as areas where upper-level  
 522 convergence is expected can be seen as an indicator for higher PM<sub>2.5</sub> concentrations at the surface.  
 523 Given the conclusions interpretable by these fields, as well as the valuable information the jet streak

524 layer presents to forecasters in terms of synoptic level transport, the U and V wind fields presented a  
 525 natural choice in terms of inclusion into LADCO SOM.

526

527 3.2.5 Variable: 850hPa Geopotential Height

Weights for GHT Level 8



528

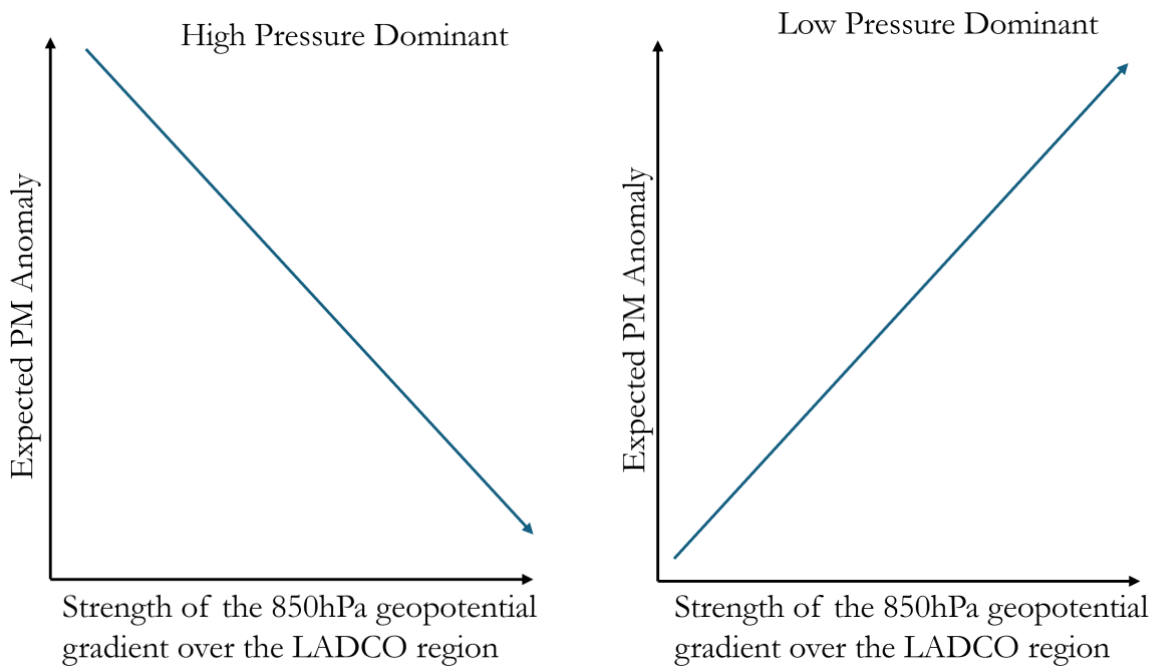
529 **Figure 17** The weights for geopotential height within LADCO SOM in meters (m).

530

531 The geopotential height field was chosen as an input into the SOM in an attempt to give LADCO  
 532 SOM a variable that can act as a classification basis for considering mid-tropospheric flow and the  
 533 vertical orientation of fronts or pressure systems. The primary geopotential height trend observable  
 534 from the SOM is in relation to the geopotential gradient. For high pressure dominated nodes, when  
 535 the distance between isohypses (lines of constant geopotential height) is large, this is associated with  
 536 less dynamical motions and stagnation conditions, and consequently, higher PM<sub>2.5</sub> anomaly. When  
 537 the LADCO region falls under a tighter geopotential gradient the result is a lower PM<sub>2.5</sub> anomaly as  
 538 seen in nodes (1,1) and (1,2). This tighter geopotential height gradient indicates stronger advection

539 within the region that clears out pollutants. However, this trend does not seem entirely robust for  
 540 lower pressure dominated nodes as the LADCO region appears to be within a loose geopotential  
 541 gradient in nodes (0,3) and (0,4) whilst the  $PM_{2.5}$  anomaly is near zero and even slightly negative.  
 542 However, it has yet to be proven if this is an artifact of our averaging methodology, where these  
 543 nodes may contain samples where both positive and negative extremes lead to a near zero average  
 544 and should be examined in future work. In any case, considering the position of the jet streak, and  
 545 the established westerly flow in nodes (2,3) and (2,4), our gradient trend appears to inverse for low  
 546 pressure dominated nodes, with tighter gradients leading to higher  $PM_{2.5}$  anomaly. Although this  
 547 trend may prove to not be linear in nature, the following mental model is provided for air quality  
 548 forecasters.

549



550

551 **Figures 18a-b** Mental model to aid in operational forecasting of air quality given the 850hPa  
 552 geopotential height field.

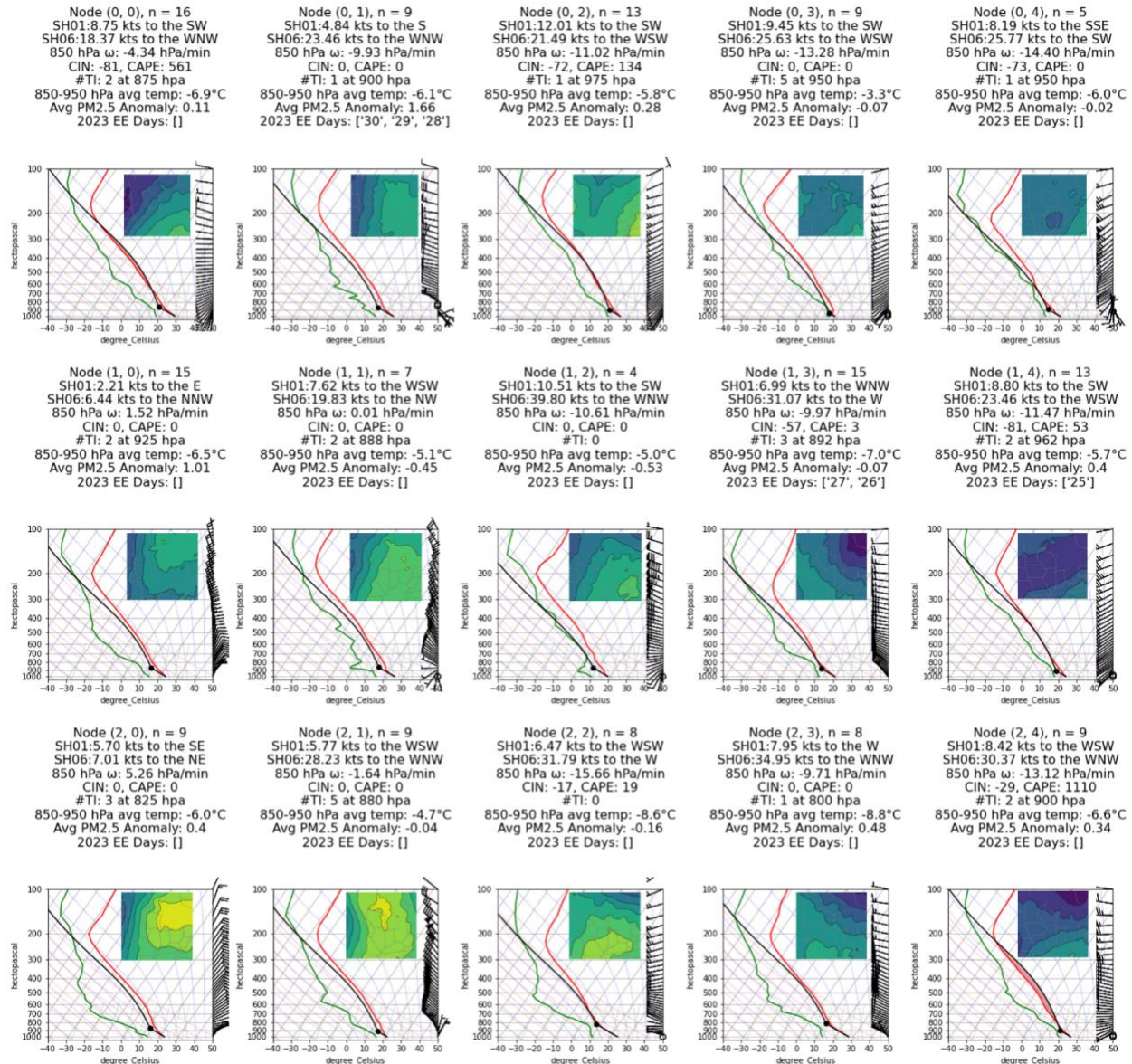
553

### 554 3.3 LADCO SOM VERTICAL PROFILE ANALYSIS

555 Following the procedure described in section 2.5, vertical profiles for each node were generated.

556 **Figure 19** visualizes the results from this procedure.

557



558

559 **Figure 19** Averaged vertical profiles for each node within LADCO SOM with MSLP pattern in  
560 upper right corner.

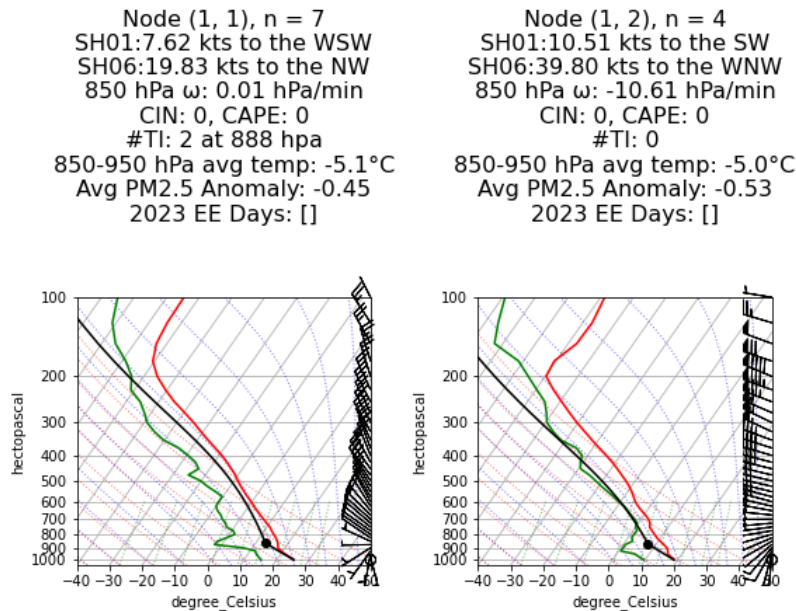
561

562 **Figure 19** presents several notable trends, the first of which can be seen by looking at the moisture  
563 profiles across all SOM nodes. Although there are exceptions, generally mid-level moisture around  
564 the 500hPa level varies greatly, and across the SOM, entire moisture profiles get more saturated  
565 starting from the lowest left node (2,0) (driest) to the uppermost right node (0,4) (most saturated).  
566 Notable exceptions to this rule are nodes (1,1) and (1,2) which appear to be surrounded by vertical  
567 profiles that are drier. However, these two nodes (that both have very negative PM<sub>2.5</sub> anomaly)  
568 uniquely seem to have a low-level dry layer near the 900hPa level. **Figure 20** displays a zoomed in  
569 version of **Figure 19** with these two nodes highlighted. Although nodes (1,1) and (1,2) do not  
570 contain an above average number of identified temperature inversions, these dry conditions at the



571 surface could indicate stronger mixing at the surface which acts to disperse pollutants upward, or the  
 572 drier conditions at the surface could work to reduce the rate chemical reactions that lead to the  
 573 formation of secondary PM<sub>2.5</sub> production such as interaction with sulfate and nitrate aerosols (which  
 574 may be particularly applicable given these sounding originate over Chicago, IL). However, more  
 575 work is necessary to see if this is indeed the case. Another similarity these nodes seem to share is  
 576 comparatively less steep 850-950hPa environmental lapse rates, a trend shared by nodes (0,3) and  
 577 (2,1) which are both accompanied by slightly negative PM<sub>2.5</sub> anomaly.

578



579

580 **Figure 20** Averaged vertical profiles for nodes (1,1) and (1,2) with low level dry layer near 900hPa.

581

582 Another trend that is apparent from **Figure 19** corresponds to the direction and strength of winds at  
 583 and near the surface level. Nodes (2,0), (1,0), and (0,1) all have slow surface winds that blow  
 584 eastward, and all have elevated PM<sub>2.5</sub> anomaly, with nodes (1,4) and (2,4) also having significantly  
 585 positive PM<sub>2.5</sub> anomaly however calm winds at the surface as opposed to eastward.

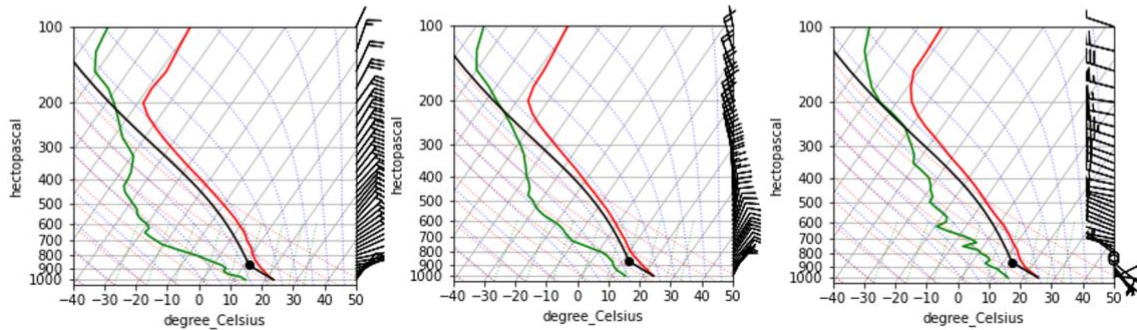
586 Also shared among high pressure dominated nodes that boast significantly high PM<sub>2.5</sub> anomaly is  
 587 low 0-1km shear. Nodes (2,0), (1,0), and (0,1) again all seem to have this in common. With all three  
 588 nodes having a high pressure dominated MSLP condition, low 0-1km shear points plainly to a  
 589 correlation of PM<sub>2.5</sub> impacts with stagnation at the surface. Nodes (2,0) and (1,0) share the additional  
 590 similarity of positive 850hPa  $\omega$  values. An indicator for downward vertical motion in the  
 591 atmosphere.

592 **Figure 21** presents a zoomed in view with a grouping of the high PM<sub>2.5</sub> anomaly nodes mentioned  
 593 above for easier reference.

Node (2, 0), n = 9  
 SH01:5.70 kts to the SE  
 SH06:7.01 kts to the NE  
 850 hPa  $\omega$ : 5.26 hPa/min  
 CIN: 0, CAPE: 0  
 #TI: 3 at 825 hpa  
 850-950 hPa avg temp: -6.0°C  
 Avg PM2.5 Anomaly: 0.4  
 2023 EE Days: []

Node (1, 0), n = 15  
 SH01:2.21 kts to the E  
 SH06:6.44 kts to the NNW  
 850 hPa  $\omega$ : 1.52 hPa/min  
 CIN: 0, CAPE: 0  
 #TI: 2 at 925 hpa  
 850-950 hPa avg temp: -6.5°C  
 Avg PM2.5 Anomaly: 1.01  
 2023 EE Days: []

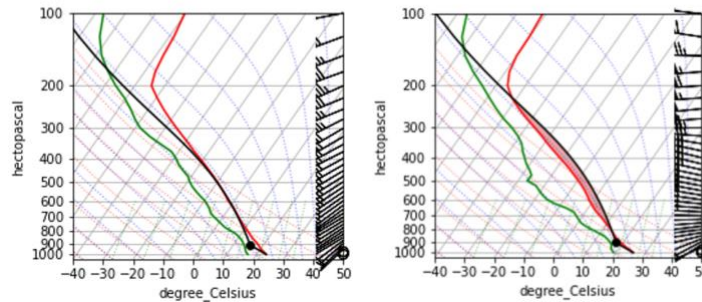
Node (0, 1), n = 9  
 SH01:4.84 kts to the S  
 SH06:23.46 kts to the WNW  
 850 hPa  $\omega$ : -9.93 hPa/min  
 CIN: 0, CAPE: 0  
 #TI: 1 at 900 hpa  
 850-950 hPa avg temp: -6.1°C  
 Avg PM2.5 Anomaly: 1.66  
 2023 EE Days: ['30', '29', '28']



594

Node (1, 4), n = 13  
 SH01:8.80 kts to the SW  
 SH06:23.46 kts to the WSW  
 850 hPa  $\omega$ : -11.47 hPa/min  
 CIN: -81, CAPE: 53  
 #TI: 2 at 962 hpa  
 850-950 hPa avg temp: -5.7°C  
 Avg PM2.5 Anomaly: 0.4  
 2023 EE Days: ['25']

Node (2, 4), n = 9  
 SH01:8.42 kts to the WSW  
 SH06:30.37 kts to the WNW  
 850 hPa  $\omega$ : -13.12 hPa/min  
 CIN: -29, CAPE: 1110  
 #TI: 2 at 900 hpa  
 850-950 hPa avg temp: -6.6°C  
 Avg PM2.5 Anomaly: 0.34  
 2023 EE Days: []



595

596 **Figure 21** Averaged vertical profiles for nodes (2,0), (1,0), (0,1), (1,4) and (2,4). The top row is  
 597 characterized by slow and eastward surface winds, bottom row with calm conditions at the surface.

598

## 599 4. APPLICATIONS

600 Although pipelines and code for operational use have yet to be implemented, LADCO SOM (or an  
 601 improved version of the SOM in the future) has potential for operational use in the field of air  
 602 quality forecasting. By inputting the conditions of, for example, a forecast hour +48 HRRR run  
 603 initialized at 0z, a classification of the modeled atmospheric conditions can be outputted by LADCO  
 604 SOM. Doing so would allow decision makers an initial forecast of the expected air quality given the  
 605 atmospheric conditions modeled for the future. Currently LADCO SOM only has knowledge of  
 606 June PM<sub>2.5</sub> events across multiple years, an updated operational model would most likely require data



607 spanning multiple months and years, unless it was determined that a more specialized SOM for a  
608 specific time period performed better. Perhaps two versions of a SOM are given different data, one  
609 corresponding to cold season months and one corresponding to warm season months. Each SOM  
610 could be asked to classify the conditions for each day, and the two SOMs running in parallel could  
611 cover any days that might be during transitions between seasons, or any anomalously warm or cool  
612 days. Moreover, separate SOMs given data corresponding to ENSO patterns (El Niño & La Niña)  
613 may also provide additional insight into broader climatological trends specific to the LADCO  
614 region.

615 Air quality forecasters may also find use in inputting the current atmospheric conditions, or the  
616 atmospheric conditions from a past event to give an indication for whether or not that particular  
617 meteorological setup is synonymous with a certain type of PM<sub>2.5</sub> anomaly. In this way it is possible  
618 to examine the relative anomaly (anomaly detection) of certain events given knowledge of data  
619 within the same seasonal period.

620 Much like the insights gained from this study, given a larger and temporally comprehensive set of  
621 data, an updated version of LADCO SOM may be able to discover harder to detect meteorological  
622 relationships and classify these similarities into representative analogs, that may have transferable  
623 relationships between pressure dominance regardless of season, or seasonal relationships regardless  
624 of pressure dominance.

625 There also exists room for policy evaluation within LADCO SOM (or a future version with an  
626 expanded dataset), although perhaps it is not the most direct way of doing so. One way to  
627 accomplish a policy evaluation using a SOM would be to provide a SOM with historical  
628 meteorological data (and optionally air quality data) before policy implementation to get a baseline  
629 understanding of typical patterns and relationships identified by the SOM initially. Then after a  
630 policy has been implemented, either look for shifts in patterns that indicate changes in quality under  
631 similar meteorological conditions, or (if given air quality data) look to see if similar clusters exist  
632 whose primary distinction is on the basis of PM<sub>2.5</sub>. If possible, then examine further similar nodes to  
633 analyze what time period samples commonly mapped to each node are from. If two nearby nodes  
634 have a similar meteorological setup, but one node has a lower average PM<sub>2.5</sub> concentration, and the  
635 samples within that node come from a time after the policy was implemented, this could be a sign  
636 that a certain policy was effective. Although perhaps a more actionable response might be to then  
637 perform a more rigorous comparative analysis of each node, now informed by the SOM of the  
638 data's meteorological similarity.

639

## 640 5. FUTURE IMPROVEMENTS

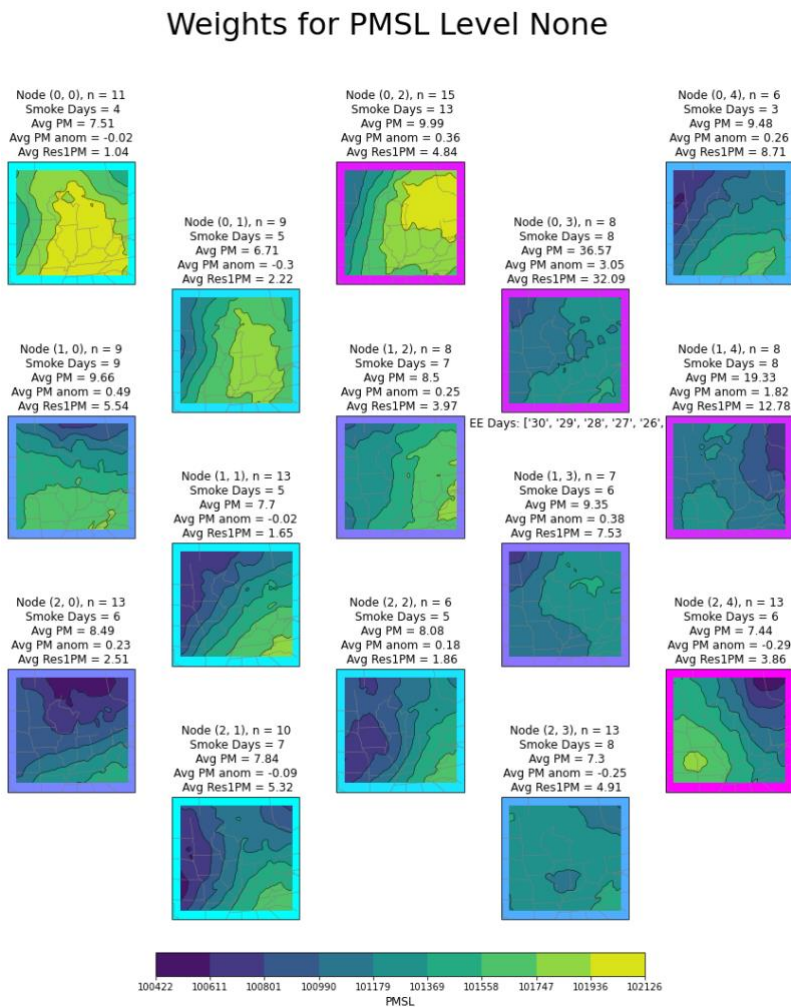
641 This section will outline aspects of LADCO SOM that should be targeted for improvement in the  
642 future. Many are simple enhancements, among them different (or just more) input data distributions  
643 as mentioned in **section 4**. However, of the improvements mentioned below, several would  
644 significantly alter (and improve) the current functionality of LADCO SOM. While worthwhile, the  
645 time required to perform such modifications is left to future work. This section will be broken down  
646 into subsections relating to each potential improvement.

647 **5.1 LADCO SOM + KRIGGED PM<sub>2.5</sub>**

648 In the later stages of the project a question occurred that was: would the results of the SOM  
 649 clustering significantly improve given knowledge of PM<sub>2.5</sub> concentrations on the ground? To answer  
 650 this question, we explored a few different options, although just imputing an average PM<sub>2.5</sub> variable  
 651 tacked on as a column on the end of the dataset was not successful, owed to the dataset’s high  
 652 dimensionality. For this a PM<sub>2.5</sub> variable with the same dimensions as a single meteorological  
 653 variable needed to be considered. This was accomplished via the creation of a “krigged” (spatially  
 654 interpolated) PM<sub>2.5</sub> dataset. A visualization of what this dataset looks like when plotted over the  
 655 study region is displayed in **Figure 5**.

656 The krigged PM<sub>2.5</sub> dataset was able to be successfully integrated into LADCO SOM’s code, however  
 657 the results were not necessarily informative to the research question. **Figure 22** presents the results  
 658 of LADCO SOM being run with all 6 meteorological variables + the krigged PM<sub>2.5</sub> dataset as input  
 659 variables.

660

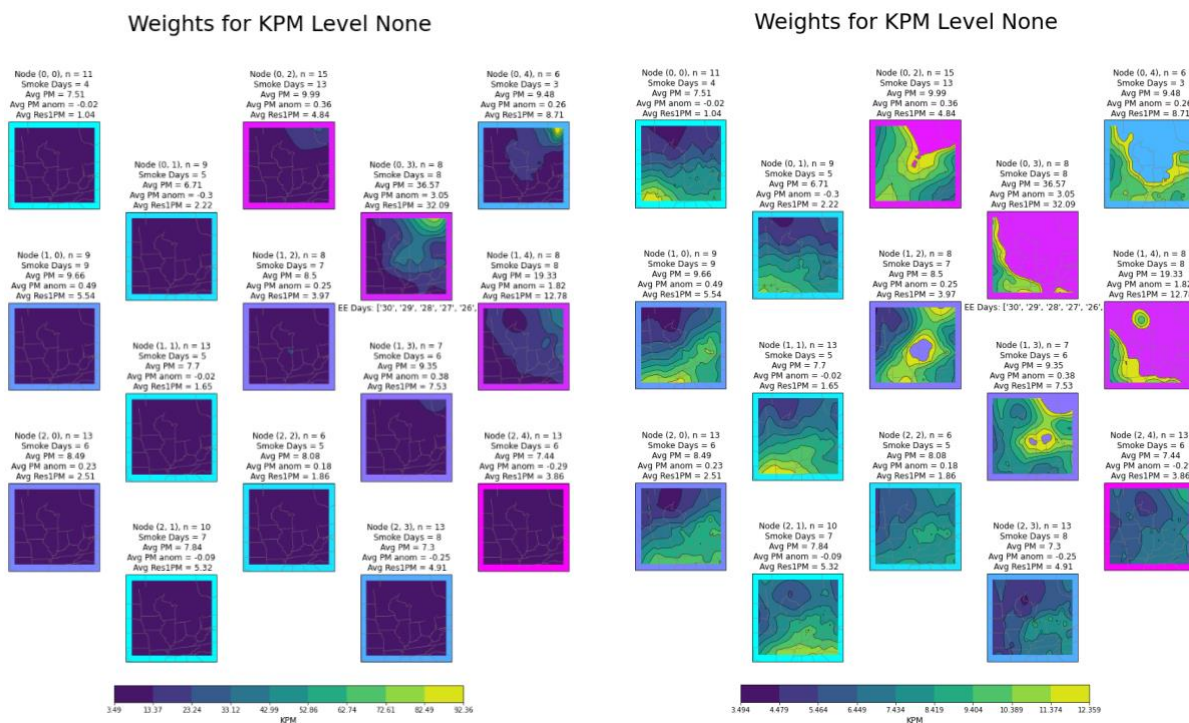


661

662 **Figure 22** The weights for MSLP within LADCO SOM with the krigged PM<sub>2.5</sub> variable.

663 The resulting SOM appears to classify every day from the 2023 EE into one node. Except, this is  
 664 questionable because we are certain that the meteorological conditions changed significantly  
 665 throughout the course of the event period. This version of LADCO SOM model built with krigged  
 666 PM<sub>2.5</sub> indicates that i) The 2023 EE days are being classified together into a node could be due to the  
 667 extremely anomalous PM<sub>2.5</sub> impacts, and ii) it calls for examination of surface input variables when  
 668 using the krigged PM<sub>2.5</sub> field. Instead of gleaning insight into the meteorological conditions  
 669 associated with these impacts, we get yet another indicator that is event is extreme, which is already  
 670 known. **Figure 23a** displays this undesirable behavior, pay particular attention to the scale in  
 671 ( $\mu\text{g}/\text{m}^3$ ), **Figure 23b** is generated with a high value mask to see PM<sub>2.5</sub> values within the first  
 672 interval of **Figure 23a**.

673



674

675 **Figure 23a-b** The weights for the krigged PM<sub>2.5</sub> variable within LADCO SOM + krigged PM<sub>2.5</sub>.

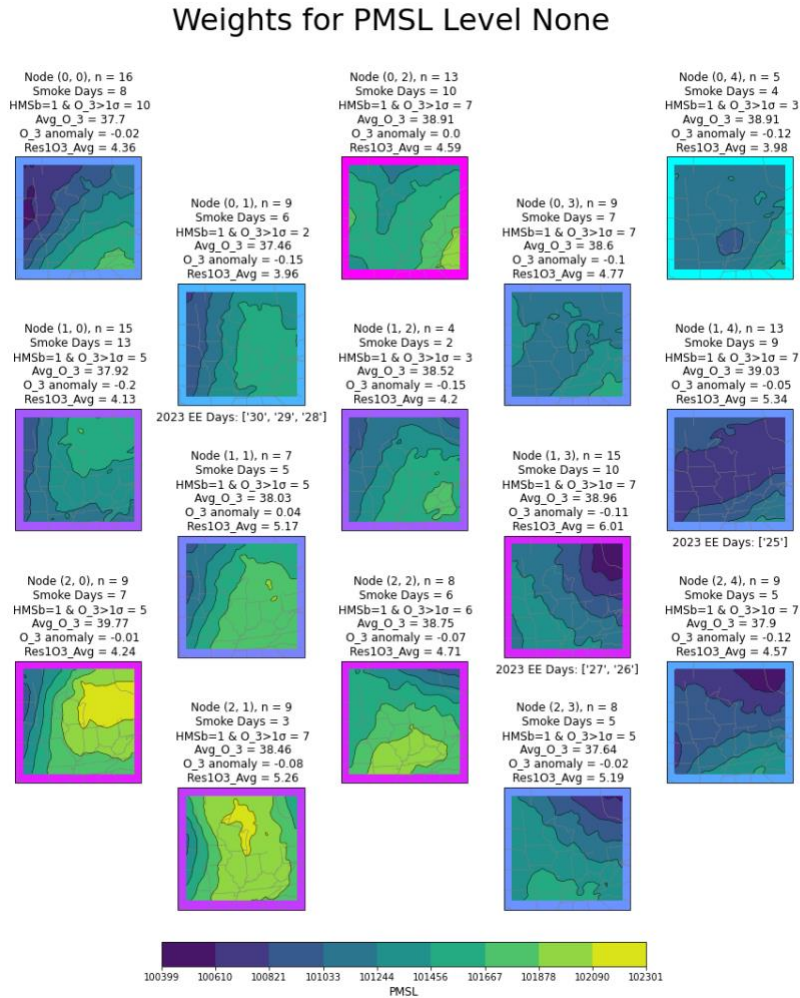
676 We can see clearly that there certainly exist interesting PM<sub>2.5</sub> relationships with meteorology data,  
 677 although we are also certain that this data is being significantly impacted by the 2023 EE. Hence,  
 678 improvements to the meteorological variables + krigged PM<sub>2.5</sub> SOM is left for future work, a  
 679 possible direction may be as simple as excluding the 2023 EE.

680

## 681 5.2 OZONE SOM ANALYSIS

682 Another SOM application that initially started off as a part of the project, but then fell off due to  
 683 time constraints is an analysis of ground level ozone. Given ozone's more predictable relationship  
 684 with meteorological conditions, a version of LADCO SOM that does a thorough analysis of ozone

685 would be extremely beneficial. In its current state LADCO SOM is capable of calculating statistics  
 686 of SOM-grouped ozone data, however interpretation of these results remains a challenge. **Figure 24**  
 687 displays LADCO SOM with ozone statistics visualized.



688  
 689 **Figure 24** The weights for MSLP within LADCO SOM with calculated ozone statistics.

690 As the data currently stands the highest ozone anomaly is 0.04 associated with node (1,1) which isn't  
 691 particularly anomalous, although interestingly is associated with a node whose vertical profile has a  
 692 dry layer near the surface. The shortfalls of LADCO SOM in ozone classification perhaps lay within  
 693 low ozone variability in the month of June for the LADCO region? Therefore, an increasingly  
 694 seasonal dataset may prove useful for further ozone analysis. Particularly for the ozone case,  
 695 incorporation of a krigged ozone dataset based off of ground monitors may provide the SOM with  
 696 extra knowledge of ozone concentrations to cluster nodes off of.

697  
 698 **5.3 LADCO SOM WITH ADDITIONAL METEOROLOGICAL VARIABLES**

699 A potential way to address the current issues with the SOM presented in 5.1, is to add more  
 700 meteorological variables to offset the proportion of the input space the krigged PM<sub>2.5</sub> takes up.



701 While this significantly adds to the dimensionality of the SOM (and substantially to the overall  
 702 runtime), considering that dimensionality already presents an issue within the study, a test was ran  
 703 considering an expanded array of meteorological variables. While the same six meteorological  
 704 variables were used within as introduced in **section 2.2** this expanded version of the SOM, LADCO  
 705 SOM was given the data for these variables at the “critical levels” within the atmosphere. The full  
 706 array of considered variables is as follows in the code:

```
707 variables = [  

  708     ('PMSL', None),  

  709     ('RH',1), ('RH',4), ('RH',8), ('RH',14), ('RH',22), ('RH',32),  

  710     ('TT',1), ('TT',4), ('TT',8), ('TT',14), ('TT',22), ('TT', 32),  

  711     ('UU',1), ('UU',4), ('UU',8), ('UU',14), ('UU', 22), ('UU',32),  

  712     ('VV',1), ('VV',4), ('VV',8), ('VV',14), ('VV',22), ('VV', 32),  

  713     ('GHT',1),('GHT',4),('GHT',8),('GHT',14),('GHT',22),('GHT',32),  

  714     ('KPM', None)  

  715 ]
```

716 Where variables only available at the surface have a level “None” and the associated pressure levels  
 717 with model levels are as such:

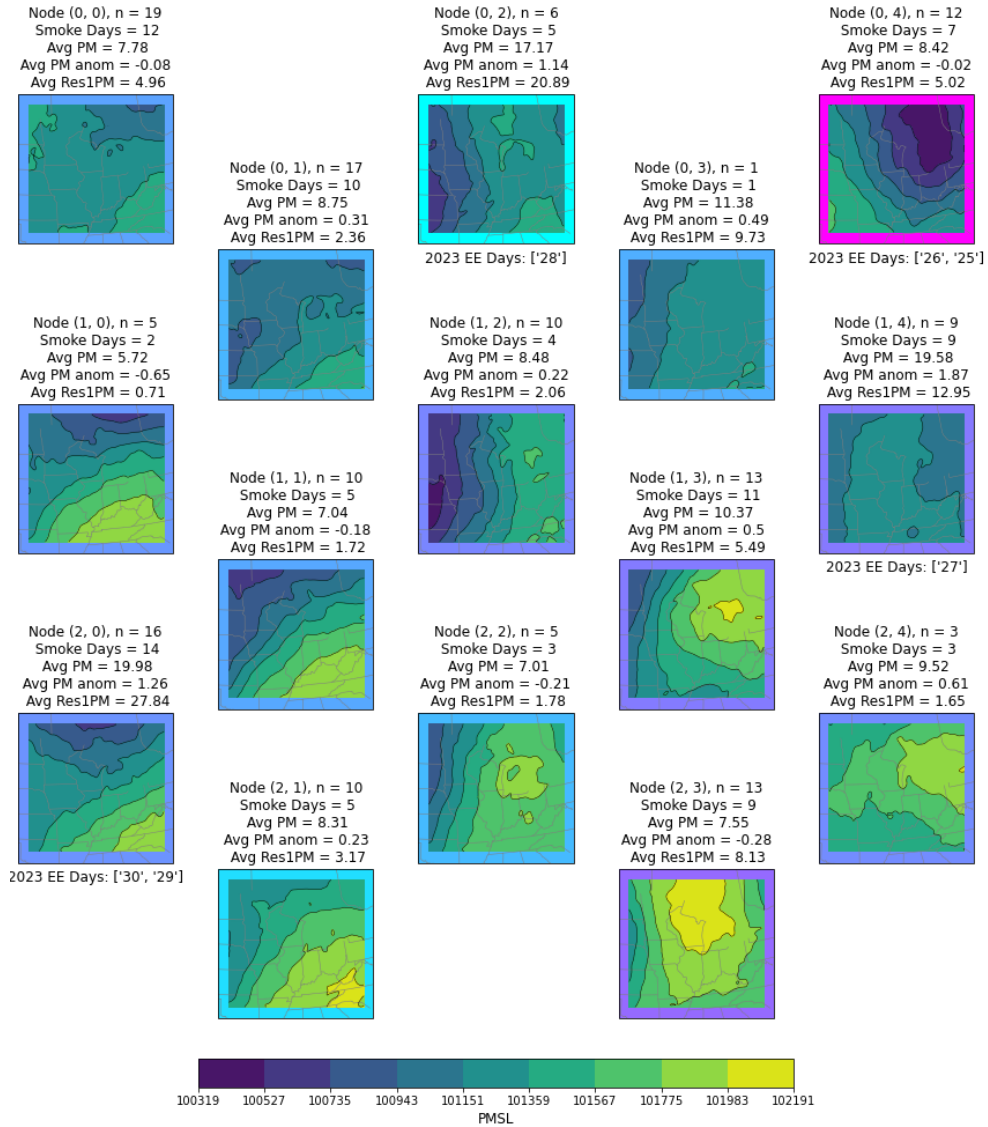
- 718 1 = Surface level (near 1000hPa)
- 719 4 = 950hPa
- 720 8 = 850hPa
- 721 14 = 700hPa
- 722 22 = 500hPa
- 723 32 = 250hPa

724 **Figure 25** demonstrates how the 2023 EE days are no longer classified into a single node, and  
 725 meteorological variables are once again primarily used for distinctions between nodes. The most  
 726 apparent problem the expanded SOM presents is exponentially higher dimensionality where each  
 727 input vector has of length ~6Million “columns”. For this reason, it is hard to gauge whether  
 728 classifications made the expanded LADCO SOM are accurate, visualizations of other variables  
 729 within the expanded LADCO SOM present somewhat contradictory information to those discussed  
 730 within this report, although the legitimacy and verification of these results is questionable and hence  
 731 left to future work. An initial step for improving the LADCO SOM model is to conduct an  
 732 exploratory analysis on a representative set of input variables guided by Principal Component  
 733 Analysis or other dimensionality reduction techniques prior to building SOM models.

734



## Weights for PMSL Level None



735

736

**Figure 25** The weights for MSLP within the expanded LADCO SOM.

737

### 738 5.4 DIMENSIONALITY AND QUANTITATIVE ANALYSIS IMPROVEMENTS

739 Owing to the extremely high dimensionality of our input data, quantitative clustering metrics have a  
 740 harder time diagnosing proper hyperparameters for LADCO SOM. The SOM presented within this  
 741 paper has the following clustering metrics:

742 Calinski-Harabasz Score: 14.7473, Silhouette Score: -0.0242, Davies-Bouldin Index: 2.4361

743 While these metrics are generally used to evaluate different hyperparameter configurations, it is  
744 plainly noticeable that our Silhouette score is negative (when higher values for Silhouette score are  
745 supposed to represent better clustering); and while this might be cause for concern in other cases, it  
746 was not mentioned earlier in the report as there is an explainable reason for this. **Figure 8**  
747 displayed the best QE achieved by LADCO SOM is around 824, and this is due to our highly  
748 dimensional data. Our Silhouette score is negative for a similar reason. Silhouette score measures  
749 how similar a point is to its own cluster compared to other clusters, which in this case leads to all  
750 points being roughly equidistant from each other, which is resulting in a negative Silhouette score.  
751 Similar interpretations can be reached when considering Calinski-Harabasz Score and the Davies-  
752 Bouldin Index.

753 Calinski-Harabasz Score: This score evaluates the ratio of the sum of between-cluster dispersion and  
754 within-cluster dispersion. High dimensionality is leading to increased within-cluster dispersion due to  
755 the curse of dimensionality, where distances between points become less meaningful as  
756 dimensionality increases.

757 Davis-Bouldin Index: This index evaluates the average similarity ratio of each cluster with its most  
758 similar cluster, considering cluster centroids. Considering the dimensionality of input data, our  
759 centroid might not be a good representation of the cluster, leading to higher index values.

760 In all three cases, the results align with challenges posed by high dimensionality. How then is  
761 LADCO SOM being evaluated? To this we point to trial and error during hyperparameter  
762 adjustment, manual inspection, and using close to default settings, with the initial requirement that  
763 every step (including the results) is understood, explainable, or interpretable. We also know from  
764 (Hewitson and Crane 2002) that the results of the SOM are less dependent on the data conforming  
765 to a specific distribution or underlying model.

766 Since no methods exist for perfectly reconstructing high dimensional data spaces within non-linear  
767 manifold representation learning, and that LADCO SOM relies on data of the same shape to  
768 reconstruct any useful visuals, we are currently at a stalemate with this dimensionality.

769 However, potential improvements on the front of LADCO SOM's dimensionality could come in  
770 the form of applying additional dimensionality reduction techniques before consideration by the  
771 SOM such as t-Distributed Stochastic Neighbor Embedding (t-SNE) or Uniform Manifold  
772 Approximation and Projection. Although in effect we are not trying to explicitly reduce the  
773 dimensionality of our data (as in get rid of less informative columns) and instead are trying to reduce  
774 the amount of data it takes to represent them. In this way LADCO SOM can still be thought of as  
775 on par with these dimensionality reduction techniques in regards to the presentation of analogs for  
776 meteorological conditions, that describe (in much fewer cases) a representative map for classifying  
777 the meteorological conditions of future PM<sub>2.5</sub> events.

778

## 779 6. CONCLUSIONS

780 The report has demonstrated how Self-organizing maps (SOMs) can provide additional insight into  
781 classifying meteorological conditions and their associated PM<sub>2.5</sub> impacts and provided justification

782 for modes of vertical transport capable of carrying fire smoke to the surface. Moreover, based on a  
783 known extreme event such as the late June 2023 event we confirm that the SOM's behavior is both  
784 predictable and explainable. Finally, we present applications for this research in the field of air  
785 quality forecasting and analysis and demonstrate the need for future research on the topic.

## 786 **ACKNOWLEDGEMENTS**

787 The author acknowledges the work of Tsengel Nergui for her preparation of all input data for  
788 LADCO SOM including the krigged PM<sub>2.5</sub> data, as well as her support throughout the duration of  
789 the project and beyond. The author also acknowledges the support of all LADCO staff for their  
790 positive encouragement and productive criticism and conversations.

791

792 **REFERENCES**

- 793 Carrasco Kind, Matias, and Robert J. Brunner. 2014. "SOMz: Photometric Redshift PDFs with Self-  
794 Organizing Maps and Random Atlas." *Monthly Notices of the Royal Astronomical Society* 438 (4): 3409–21.  
795 <https://doi.org/10.1093/mnras/stt2456>.
- 796 Deboeck, Guido, and Teuvo Kohonen. 2013. *Visual Explorations in Finance: With Self-Organizing Maps*.  
797 Springer Science & Business Media.
- 798 Forest, Florent, Mustapha Lebbah, Hanene Azzag, and Jérôme Lacaille. 2021. "Deep Embedded  
799 Self-Organizing Maps for Joint Representation Learning and Topology-Preserving Clustering."  
800 *Neural Computing and Applications* 33 (24): 17439–69. <https://doi.org/10.1007/s00521-021-06331-w>.
- 801 Hewitson, B. C., and R. G. Crane. 2002. "Self-Organizing Maps: Applications to Synoptic  
802 Climatology." *Climate Research* 22 (1): 13–26.
- 803 Hrust, Lovro, Zvezdana Bencetić Klaić, Josip Križan, Oleg Antonić, and Predrag Hercog. 2009.  
804 "Neural Network Forecasting of Air Pollutants Hourly Concentrations Using Optimised Temporal  
805 Averages of Meteorological Variables and Pollutant Concentrations." *Atmospheric Environment* 43 (35):  
806 5588–96. <https://doi.org/10.1016/j.atmosenv.2009.07.048>.
- 807 Hulle, Van, and M Marc. 2012. "Self-Organizing Maps." *Handbook of Natural Computing* 1:585–622.
- 808 Keyser, Daniel, and M. A. Shapiro. 1986. "A Review of the Structure and Dynamics of Upper-Level  
809 Frontal Zones." *Monthly Weather Review* 114 (2): 452–99. [https://doi.org/10.1175/1520-0493\(1986\)114<0452:AROTSA>2.0.CO;2](https://doi.org/10.1175/1520-0493(1986)114<0452:AROTSA>2.0.CO;2).
- 811 Kohonen, Teuvo. 1982. "Self-Organized Formation of Topologically Correct Feature Maps."  
812 *Biological Cybernetics* 43 (1): 59–69. <https://doi.org/10.1007/BF00337288>.
- 813 Törönen, Petri, Mikko Kolehmainen, Garry Wong, and Eero Castrén. 1999. "Analysis of Gene  
814 Expression Data Using Self-organizing Maps." *FEBS Letters* 451 (2): 142–46.  
815 [https://doi.org/10.1016/S0014-5793\(99\)00524-4](https://doi.org/10.1016/S0014-5793(99)00524-4).
- 816 Vesanto, J., and E. Alhoniemi. 2000. "Clustering of the Self-Organizing Map." *IEEE Transactions on*  
817 *Neural Networks* 11 (3): 586–600. <https://doi.org/10.1109/72.846731>.

818