1
2 **Forecasting and Identifying the Meteorological and Hydrological Conditions Favoring the**
3 **Occurrence of Severe Hazes in Beijing and Shanghai using Deep Learning**
4
5
6 Chien Wang
7

8 Laboratoire d'Aerologie, CNRS and University Paul Sabatier
9 14 Avenue Edouard Belin, 31400 Toulouse, France

10
11
12
13
14
15
16
17
18
19
20
21

22 *Correspondence to*: Chien Wang (chien.wang@aero.obs-mip.fr)

23
24

25 **Abstract.** Severe haze or low visibility event caused by abundant atmospheric aerosols has
26 become a serious environmental issue in many countries. A framework based on deep
27 convolutional neural networks has been developed to forecast the occurrence of such events in
28 two Asian megacities: Beijing and Shanghai. Trained using time sequential regional maps of
29 meteorological and hydrological variables alongside surface visibility data over the past 41
30 years, the machine has achieved a good overall accuracy in associating the haze events with
31 favorite meteorological and hydrological conditions. Furthermore, an unsupervised cluster
32 analysis using features with a greatly reduced dimensionality produced by the trained machine
33 has, arguably for the first time, successfully categorized typical regional meteorological-
34 hydrological regimes alongside local quantities associated with haze and non-haze events in the
35 two targeted cities, providing substantial insights to advance our understandings of this
36 environmental extreme.

## 1 Introduction

38    Frequent low visibility or haze event caused by elevated abundance of atmospheric aerosols
39 due to fossil fuel and biomass burning has become a serious environmental issue in many Asian
40 countries in recent decades, interrupting economic and societal activities and causing human
41 health issues (*e.g*., Chan and Yao, 2008; Silva et al., 2013; Lee *et al*., 2017). For example, rapid
42 economic development and urbanization in China have caused various pollution-related health
43 issues particularly in populated metropolitans such as Beijing-Tianjin region and Yangtze river
44 delta centered in Shanghai (*e.g*., Liu *et al*., 2017). In Singapore, the total economic cost of severe
45 haze events in 2015 is estimated to be $510 million (0.17% of GDP), or $643.5 million based on
46 a wiling-to-pay analysis (Lin *et al*., 2016). To ultimately prevent this detrimental environmental
47 extreme from happening requires rigid emission control measures in place through significant
48 changes in energy consumption as well as land and plantation management. Before all these
49 measures could finally take place, it would be more practical to develop skills to accurately
50 predict its occurrence hence to allow mitigation measures to be implemented ahead of time.
51    Severe haze events arise from the solar radiation extinction by aerosols in the atmosphere,
52 this mechanism can be enhanced with the increase of relative humidity that enlarges the size of
53 particles (*e.g*., Kiehl and Briegleb 1993). Aerosols also need favorite atmospheric transport and
54 mixing conditions to reach places away from their immediate source locations, while their
55 lifetime in the atmosphere can be significantly reduced by rainfall removal. In addition, soil
56 moisture is also a key to dust emissions. Therefore, meteorological and hydrological conditions
57 are critical to the occurrence of haze events besides particulate emissions. To forecast the
58 occurrence of such events using existing atmospheric numerical models developed based on fluid
59 dynamics and explicit or parameterized representations of physical and chemical processes, the
60 actual task is to accurately predict the concentration of aerosols at a given geographic location
61 and a given time in order to correctly derive surface visibility (*e.g*., Lee *et al*. 2017 & 2018).
62 However, the propagation of numerical or parameterization errors through the model integration
63 could easily drift the model away from the original track, not mentioning that lack of real-time
64 emission data alone would simply handicap such an attempt. Therefore, a more fundamental
65 issue in practice is whether these models could reproduce the *a posteriori* distribution of the
66 possible outcomes of the targeted low-probability extreme events. Ultimately, lack of knowledge
67 about the extreme event would, in turn, hinder the effort to improve the forecasting skills.

68      Differing from the deterministic models, an alternative statistical prediction approach could
69   be adopted should the predictors of a targeted event could be identified and a statistical
70   correlation between them could be established with confidence. However, this is a rather difficult
71   task for the traditional approaches because it requires an analysis dealing with a very large
72   quantity of high-dimensional data in order to establish a likely multi-variate and nonlinear
73   correlation of generalization. Nevertheless, such attempts can obviously benefit now from the
74   fast-growing machine learning (ML) and deep learning (DL) algorithm development (*e.g.*,
75   LeCun *et al*., 2015). In addition, technological advancement and continuous investment from
76   governments and other sectors across the world have led to a rapid increase of quantity alongside
77   substantially improved quality of meteorological, oceanic, hydrological, land, and atmospheric
78   composition data. These data might still not be sufficient for evaluating and improving certain
79   detailed aspects of the deterministic forecasting models. However, rich information contained in
80   these data about favorite environmental conditions for the occurrence of extreme events such as
81   hazes could already have a great value for developing alternative forecasting skills.
82      Many Earth science applications dealing with meteorological or hydrological data need a
83   trained machine to not only forecast values but also recognize patterns or images. However, this
84   can easily lead to a curse of dimensionality of many traditional ML algorithms. Fortunately, deep
85   learning that directly links large quantity of raw data with targeted outcomes through deep
86   convolutional neural networks or CNNs (Goodfellow *et al*., 2016) offers a clear advantage in
87   sufficiently training deep networks suitable for solving highly nonlinear issues. In doing so, DL
88   can also eliminate the possible mistakes in data derivation or selection introduced by subjective
89   human opinion regarding a poorly understood phenomenon. Recently, DL algorithms have been
90   explored in various applications in atmospheric, climate, and environmental sciences, ranging
91   from recognizing specific weather patterns (*e.g*., Liu *et al*., 2016; Kurth *et al*., 2018; Lagerquist
92   *et al*., 2019; Chattopadhyay *et al*., 2020), weather forecasting including hailstorm detection (*e.g*.,
93   Grover *et al*., 2015; Shi *et al*., 2015; Gagne *et al*., 2019), to deriving model parameterizations
94   (*e.g*., Jiang *et al*., 2018), and beyond.
95      When weather patterns associated with targeted outcome are known or irrelevant to the task,
96   the forecasting can be normally proceeded to recognize a given pattern by using pattern-to-
97   pattern correlation from sequential training data with spatial-information-preserving full CNNs
98   such as U-net (Ronneberger *et al*., 2015; Weyn *et al*., 2020). However, this is certainly not the
99   case for the applications where the environmental conditions associated with targeted outcome
100  are yet known. For such applications, a possible solution is to utilize a large quantity of raw data
101  with minimized human intervention in data selection to train a deep CNN in order to associate
102  targeted outcomes with favorite environmental conditions. This study represents such an attempt,
103  where a DL forecast framework is trained to identify the meteorological and hydrological
104  conditions associated with the occurrences of severe hazes. The DL framework has been
105  developed initially with the severe hazes in Singapore (Wang, 2020), and now hazes in two
106  megacities of China, Beijing and Shanghai. In terms of particulate pollutant emissions, all these
107  cities share certain sources including fossil fuel combustions from transportation, domestic, and
108  industries. On the other hand, each city also has its own unique sources, for instance, desert and
109  perhaps anthropogenic dust for Beijing, and massive biomass burning in Singapore (Chen *et al*.,
110  2013; Liu *et al*., 2017; Lee *et al*., 2017, 2018, & 2019). It is obvious that besides meteorological
111  and hydrological conditions, dynamical patterns of anthropogenic activities leading to the
112  emissions of particulate matters are also important factors behind the occurrence of severe hazes.
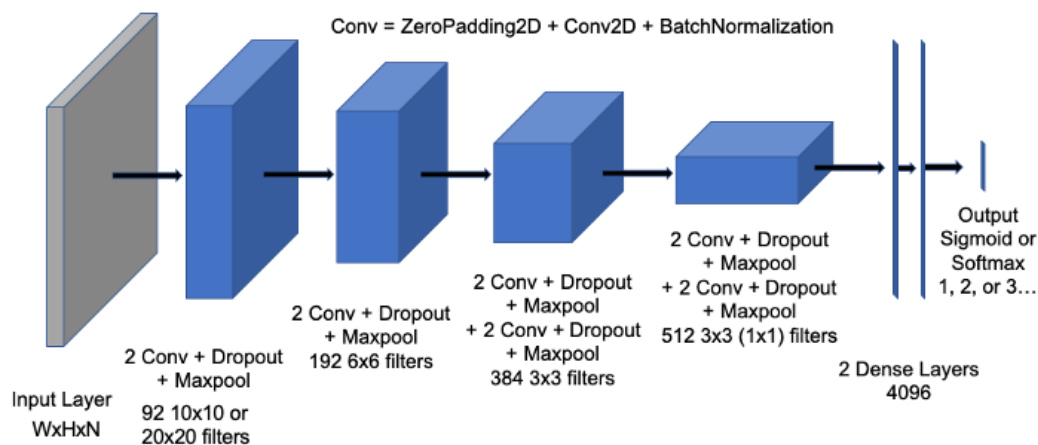113  Nevertheless, the major purpose of this study is to advance our fundamental knowledge about the

114    weather conditions favoring the occurrence of hazes and, through an in-depth analysis on the
115    forecasting results to identify the limit of such a machine and thus to provide useful information
116    for establishing a more complete forecasting platform for the task.
117        In the paper, the architecture alongside method and data for training are firstly described after
118    this Introduction, followed by a discussion of training and validation results. Then, an
119    unsupervised cluster analysis benefited from the trained machine is introduced along with the
120    results that furthers the understanding of the CNN's performance and summarizes, for the first
121    time, the various typical meteorological and hydrological regimes associated with haze versus
122    non-haze situations in the two cities. The last section concludes the major efforts and findings.

123    **2 Network Architecture, Training Methodology and Data**

124        The convolutional neural network used in this study, the HazeNet (Wang 2020), has been
125    developed by adopting the general architecture of the CNN developed by the Oxford
126    University's Visual Geometry Group or VGG-Net (Simonyan and Zisserman, 2015). The actual
127    structure alongside hyper-parameters of HazeNet have been adjusted and fine-tuned based on
128    numerous test trainings. In addition, certain techniques that were not available when the original
129    VGG net was developed, *e.g.*, batch normalization (Ioffe and Szegedy, 2015), have been
130    included as well. The current version for haze applications of Beijing and Shanghai contains
131    20,507,161 parameters (11,376 non-trainable). Figure 1 shows the general architecture of a
132    HazeNet version with 12 convolutional and 4 dense layers (in total 57 layers).
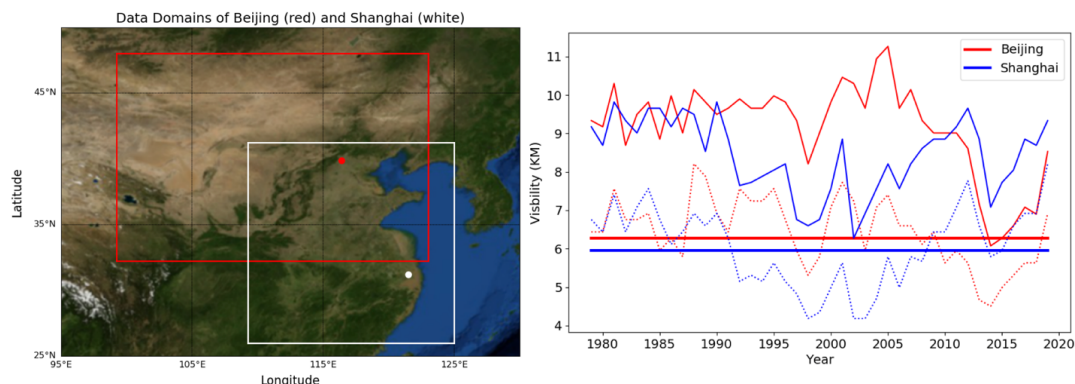


133
134    **Figure 1**. Architecture of the 12 convolutional plus 4 dense layer HazeNet. Here "Conv" represents a unit
135    containing a zero-padding then a 2D convolutional layer, followed by a batch normalization layer. There
136    is a flatten layer before the 2 dense layers. W = width, H = height, and N = number of features of the
137    input fields, they are 64, 96, and 16 for Beijing, and 64, 64, and 16 for Shanghai case, respectively.

138        The network has been trained in a standard supervised learning procedure for classification.
139    In this procedure, the network takes input features to produce classification output that are then
140    compared with known results or labels based on observations. The coefficients of the network
141    are thereafter optimized in order to minimize the error between the prediction and the

142   observation or label. The loss function used in optimization is cross-entropy (*e.g.*, Goodfellow *et*
143   *al.*, 2017). Such a procedure is repeated until the performance of the network can no longer be
144   improved. In practice, the trainings usually last about 2000 epochs (each epoch is a training cycle
145   that uses up the entire training dataset). This procedure in nature is to train a deep CNN to
146   recognize then associate input features (bundled meteorological and hydrological conditions in
147   this case) with corresponding class, *i.e.*, severe haze events or non-haze events. As a result, the
148   knowledge specifically about the favorite meteorological and hydrological conditions of severe
149   hazes could be advanced.
150       The labels for the training are derived using the observed daily surface visibility (*vis.*
151   thereafter), obtained from the Global Surface Summary Of the Day or GSOD dataset consisting
152   of daily observations of meteorological conditions from tens of thousands of airports around the
153   globe (Smith *et al.*, 2011). In the cases of Beijing and Shanghai, data are from the time period
154   from 1979 to 2019, containing 14975 samples. For simplicity, the discussions will be mainly on
155   the 2-class training, where events with *vis.* ≤ the long-term mean value of the 25$^{th}$ percentile or
156   p25 of *vis.* (6.27 km in Beijing, 5.95 km in Shanghai; Fig. 2, right panel; also Fig. S1 in
157   Supplementary) are defined as class 1 or severe hazes, otherwise the class 0 or non-haze cases.
158   The p25 values actually represent a substantial reduction of *vis.* due to high particulate pollution
159   (*e.g.*, Lee *et al.*, 2017). Note that unlike in the case of Singapore (Wang 2020), fog and mist are
160   more common low visibility events in Beijing and Shanghai and thus have been excluded from
161   the labels of severe hazes by following GSOD fog marks. The number of severe haze events
162   occurred during 1979-2019 defined in the above procedure is 2999 and 3099 for Beijing and
163   Shanghai, or in a frequency of 20.0% and 20.7%, respectively.



164
165   **Figure 2**. (Left) The input-feature defining domains for Beijing (red box and dot, 99.25 - 123E, 32.25-
166   48N; 96x64 grids with ERA5 data) and Shanghai (white box and dot, 109.25-125E, 26-41.25N; 64x64
167   grids), made using Basemap library, a matplotlib extension. (Right) Annual means (solid curves), 25$^{th}$
168   percentiles (dash curves), and 25$^{th}$ percentile means (solid straight lines) of surface visibility in Beijing
169   (red) and Shanghai (blue) between 1979 and 2019.

170       The training and validation of HazeNet also need the input features with the same sample
171   dimension of the labels. These input data are derived from hourly longitude-latitude maps of
172   meteorological and hydrological variables covering the data collection domain (Fig. 2, Left),
173   obtained from ERA5 reanalysis data produced by the European Centre for Medium-range
174   Weather Forecasts or ECMWF (Hersbach *et al.*, 2020). These data are distributed in a grid
175   system with a horizontal spatial interval of 0.25 degree. Up to 16 features are derived from the
176   original hourly data fields covering the analysis domain respectively for Beijing (64x96 grids)

Atmospheric
Chemistry
and Physics
Discussions
Open Access
EGU

177  and Shanghai (64x64 grids), including: daily mean of surface relative humidity (REL thereafter);
178  diurnal change as well as daily standard deviation of 2-meter temperature or DT2M and T2MS,
179  respectively; daily mean of 10-meter zonal and meridional wind speed or U10 and V10,
180  respectively; daily mean of total column water (TCW); daily mean (TCV) and diurnal change
181  (DTCV) of total column water vapor; daily mean of planetary boundary layer height (BLH);
182  daily mean soil water volume in soil layer 1 and 2 or SW1 and SW2, respectively; daily mean of
183  total cloud cover (TCC); daily mean geopotential heights at 500 (Z500) and 850 (Z850) hPa
184  pressure levels along with their diurnal changes D500 and D850, respectively. All input features
185  have been normalized into a range of [-1, +1] (Fig. S2 in Supplementary).
186      Before the training, the entire samples of labels alongside corresponding input features were
187  randomly shuffled first then split as: 2/3 of the samples went to training set and 1/3 to validation
188  set, each is used duly for its designated purpose throughout the entire training process without
189  switch. The above procedure treats each of the events as an independent one. For the
190  convenience in comparing performance or restarting training based on a saved machine, a pair of
191  saved training and validation datasets produced following the above procedure was used.
192      The number of samples used in training HazeNet is rather limited in deep learning standard.
193  However, to associate 16 joint two-dimensional maps with targeted labels even with the current
194  number of samples is still a demanding task, requiring a deep CCN to accomplish. Furthermore,
195  targeted severe hazes are a low probability event. Its frequency of appearance is about 20.0% in
196  Beijing and Shanghai cases. Therefore, trained machine would easily bias toward the
197  overwhelming non-haze events. To resolve these issues, a combination of class-weight and batch
198  normalization has been implemented in HazeNet. This approach has effectively reduced the
199  overfitting while overcome the data imbalance issue, making the long training of a deep CNN
200  become possible (Wang, 2020).

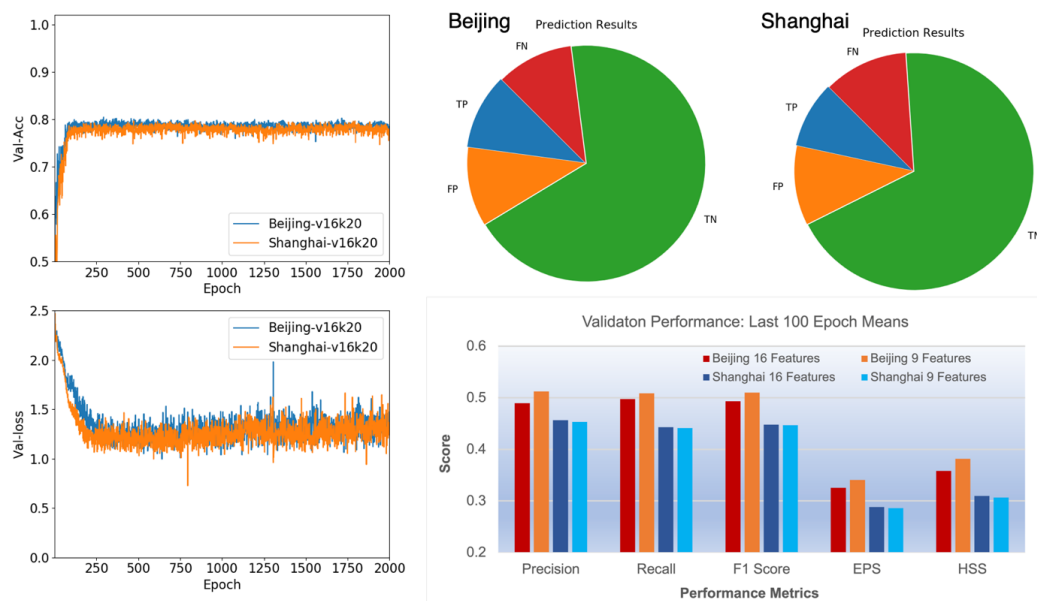201  **3 Training and Validation Results of Haze Forecasting**

202      Currently, it is still difficult to find any practical score in forecasting the occurrence of severe
203  hazes for comparison. Therefore, the performance of HazeNet has been mainly measured by
204  using certain commonly adopted metrics for classification largely derived from the concept of
205  the so-called confusion matrix (*e.g.*, Swets, 1988; Table A), including *accuracy*, *precision*,
206  *recall*, *F1 score*, *equitable threat score* or *ETS*, and *Heidke skill score* or *HSS* (Appendix A).
207  Unless otherwise indicated, the discussions on the performance scores are hereafter referring to
208  the severe haze class, or class 1, and obtained from validation rather than training. In all the
209  cases, the performance metrics referring to non-haze or class 0 has much better scores.
210      In order to train a stable machine, trainings with 2000 epochs or longer have been conducted
211  instead of using certain commonly used skills such as early stop. As a result, the validation
212  performance metrics of the trained machines all appeared to be stabilized by approaching the end
213  of training (Fig. 3). These scores were consistent with the results of ensemble training with the
214  same configuration but different randomly selected training and validation datasets, and also
215  comparable among trainings with different configurations. Overfitting has been clearly overcome
216  due to such a long training procedure alongside the adoption of class0weight and batch
217  normalization. In a 2-class classification (haze vs. non-haze), trained deep HazeNet can always
218  reach an almost perfect training accuracy (e.g., 0.9956 for Beijing cases) and a validation
219  accuracy of 80% in both Beijing and Shanghai cases, or the no-skill forecast accuracy for no-
220  haze (Fig. 3, left). At the same time, the performance scores in predicting specifically severe

221    hazes are also very reasonable, *e.g*., for Beijing cases either precision or recall exceeds 0.5 (they
222    normally evolve in opposite direction), leading to a nearly 0.5 *F1 Score* (Fig.3, right). The
223    corresponding scores in training are obviously much higher, *e.g*., with precision, recall, and F1 as
224    0.9804, 0.9980, and 0.9880, respectively for Beijing cases, owing to the deep and thus powerful
225    CNNs. HazeNet performed slightly better than several known deep CNNs such as Inception Net
226    V3 (Szegedy *et al*., 2015), ResNet50 (He *et al*., 2015), and VGG-19 (Simonyan and Zisserman,
227    2015) in the same haze forecasting task (Wang, 2020). Nevertheless, as indicated previously that
228    a nearly perfect validation performance is not realistic since meteorological and hydrological
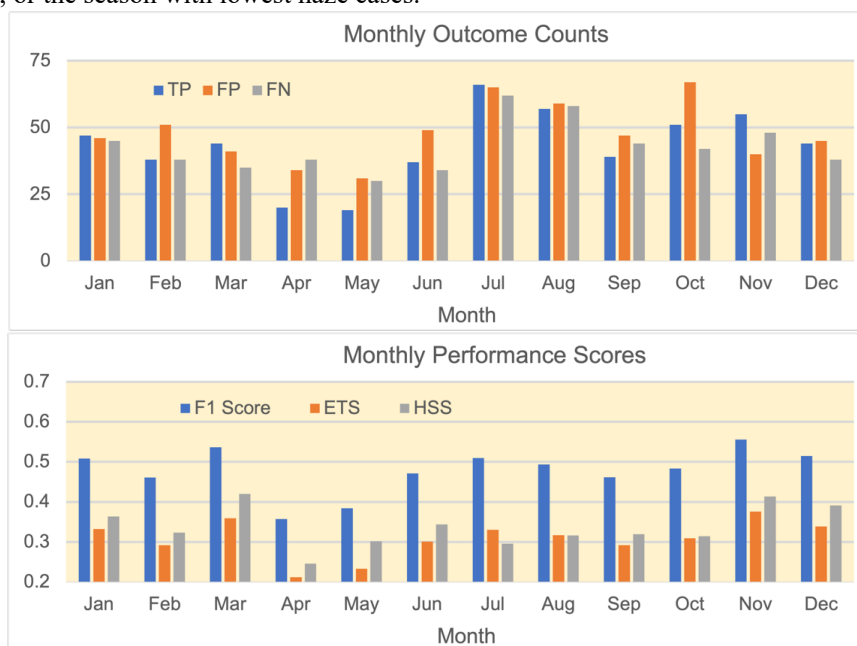229    conditions are not the only factors behind the occurrence of haze events.

230



231
232    **Figure 3**. (Left) Validation accuracy (top panel) and loss (lower panel) of HazeNet with 16 features for
233    Beijing and Shanghai cases, kernel size for the first filter is 20x20. (Right Top) Prediction outcomes in
234    reference to haze events or class-1 of Beijing and Shanghai. Here TP = true positive, TN = true negative,
235    FP = false positive, and FN = false negative prediction outcomes. (Right Bottom) Scores of performance
236    metrics as last 100 epoch means for Beijing and Shanghai with 16 and 9 features, respectively.

237    Looking into the specific prediction outcomes in referring to severe haze, the trained machine
238    has produced considerably higher ratio of true positive or TP outcomes than in the Southeast
239    Asia cases (Wang, 2020) despite a number of outcomes of false positive or FP (*i.e*., false alarm)
240    and false negative or FN (*i.e*., missing forecast). In forecasting the severe hazes in Beijing, the
241    trained machine performs reasonably well throughout all months except for April and May or the
242    major dusty season there, producing F1 score, ETS, and HSS all exceed or near 0.5 as well as the
243    number of TP outcomes is higher than that of FN (Fig. 4). The performance of HazeNet actually
244    improves in months with higher observed haze events. For Beijing, the lowest haze season is
245    during the dusty April and May when all the major performance metrics are lower than 0.4, and
246    the machine produces more missing forecasts than true positive outcomes. The relatively poor
247    performance in spring suggests that the weather and hydrological features associated with dust-

248 dominated haze events during this period might differ from the situations in the other seasons
249 when hazes are mainly caused by local particulate pollution. For Shanghai cases, HazeNet
250 performs better during late autumn and entire winter (from November to February) when haze
251 occurs most frequently. The worst performance comes from the monsoon season (July to
252 October), or the season with lowest haze cases.

253



254
255 **Figure 4.** (Top) Predicted TP, FP, and FN outcomes and (Bottom) performance scores for each month.
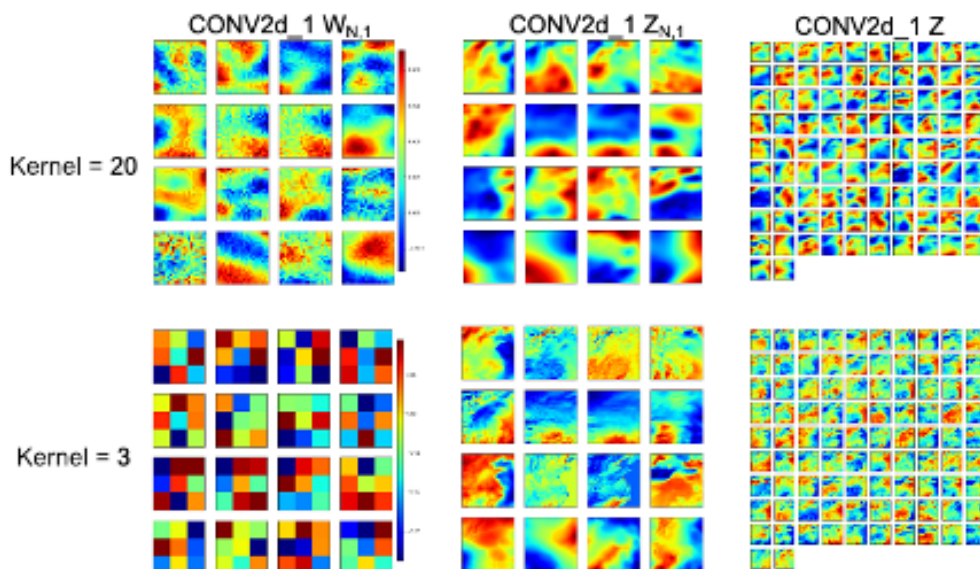256 All from validation of Beijing cases with 16 features.

257     **Kernel size and CNN performance.** The deep architecture of HazeNet and the long training
258 procedure have actually made the performance less sensitive to many hyperparameters of the
259 network. One hyperparameter, however, is specifically interesting to explore for an application
260 using large quantity of meteorological maps, that is the kernel size of the first convolutional
261 layer, where the input data, *i.e.*, meteorological and hydrological maps are convoluted then
262 propagated into the subsequent layers. Meteorological maps or images often contain
263 characteristic patterns with different spatial scales. Intuitively, preserving these patterns could be
264 important in predicting the targeted extremes. Apparently, a larger kernel size produces smoother
265 output images from the first convolutional layer, while a smaller kernel size can preserve many
266 spatial details of the meteorological maps as demonstrated from the layer output shown in Fig. 5.
267 In practice, however, the patterns produced by the latter configuration might be too complicated
268 for the networks to recognize and to perform classification, whereas patterns resulted from a
269 relatively larger kernel size for the first convolutional layer might be more characteristic for the
270 task. The actual result suggests that HazeNet configured with a first-layer kernel size of 20 to 26
271 or close to 5 – 6 degree in spatial 'resolution', consistently produces a better performance (about
272 a 10% improvement in *F1 score*) than that by a smaller kernel size of 3 or 6. As a result, a kernel
273 size of 20 has been adopted as the default configuration for the first 2 convolutional layers in this
274 study.

275
276  **Figure 5**. (Left column) Weight coefficients of the first filter set ($W_{N,1}$), (Middle column) partial output
277  for each feature ($Z_{N,1}$), and (Right column) the output ($Z$) of the first convolution layer (CONV2d_1) with
278  two selected kernel sizes or ks: (upper panels) 20x20 and (lower panels) 3x3. Here $W$ represents the filters
279  and $Z$ the output of convolution, the subsets of $Z$ before the feature dimension is merged can be expressed
280  as: $Z_{N,i} = W_{N,i}(ks, ks) \cdot f_N^T(ks, ks)$, with the order of input features $N = 1,\ldots 16$ and $i$ represents the
281  convolutional layer index, *i.e.*, 1 is the first layer or CONV2d_1. For the first layer, input feature size is
282  $(h,w) = (64, 64)$, the sets of filters is 92, thus the final output $Z$ has a dimension of $(h\text{-}ks\text{+}1, w\text{-}ks\text{+}1, 92)$.
283  Shown are results from the trainings for Shanghai haze cases.

284  **Reducing the number of input features**. One recognized advantage of deep CNN in
285  practice is its capacity to directly link the targeted outcome with a large quantity of raw data to
286  avoid human misjudgment in selecting and abstracting input features due to a lack of knowledge
287  about the application task. Nevertheless, for an application such as this one that uses a large
288  number of meteorological and hydrological variables (or channels in machine learning term),
289  reducing the number of input features with minimized influence on the performance can still
290  benefit the efforts of establishing physical or dynamical causal relations and beyond.
291      There are certain available methods to rank feature then reduce some unimportant ones.
292  These do not work straightforwardly for deep CNNs (*e.g.*, McGovern *et al.*, 2019). In the
293  previous effort, this has been done by testing the sensitivity of the full network performance in
294  real training with either a single feature or all but one features (Wang, 2020), which apparently is
295  also a demanding task. Here, another attempt has been made to use a trained then saved machine
296  to examine the sensitivity of the network to various features (Appendix B).
297      The sensitivity analyses for Beijing and Shanghai cases have obtained largely consistent
298  results, indicating that the network is more sensitive to the same 9 features than the other 7 (Fig.
299  S3). The highest-ranking features though differ, with diurnal change of column vapor (DTCV)
300  and soil water content in the second soil layer (SW2) as the most sensitive features for Beijing,
301  while relative humidity (REL) and planetary boundary layer height (BLH) for Shanghai. Most
302  importantly, trainings using only the top 9 most sensitive features have produced a performance
303  equivalent to or even better than the same training but with 16 features (Fig. 3). With reduced
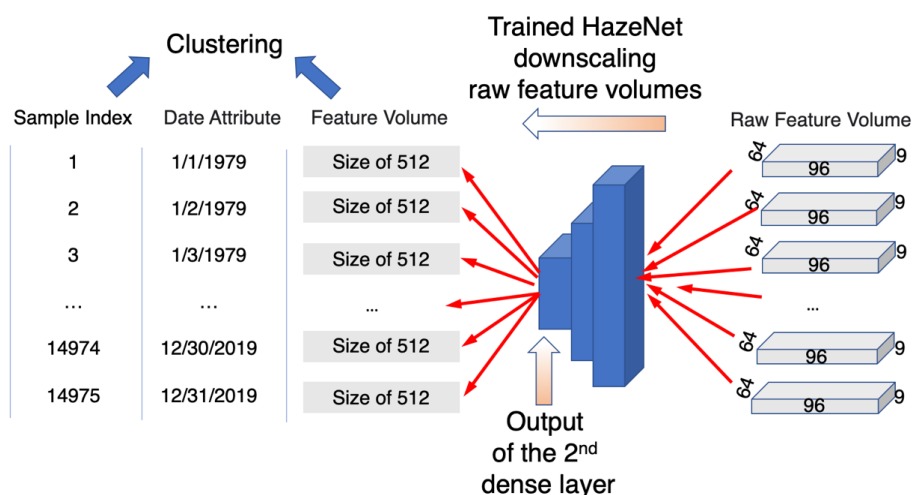
304    number of features, many further analyses can be conducted with less workload and produce
305    results that are easily understood.

**4 Identifying and Categorizing the Typical Regional Meteorological and Hydrological**
**Regimes Associated with Haze and Non-Haze Events**

308        A major purpose of this study is to identify the meteorological and hydrological conditions
309    favoring the occurrence of severe hazes in the targeted cities. When using a dataset with a large
310    number of samples, this type of analyses could be better accomplished by applying, *e.g.*, cluster
311    analysis (*e.g.*, Steinhaus, 1957), a standard unsupervised ML algorithm that groups data samples
312    into various clusters in such a way that samples in the same cluster are more similar to each other
313    than to those in other clusters. Specifically for this study, the derived clusters would likely
314    represent various regimes in terms of combined meteorological and hydrological conditions for
315    associated events. However, applying cluster analysis directly to a large number of samples, each
316    with a feature volume of ~50000 is an uneasy task. A dimensionality reduction is apparently
317    needed to reduce the feature volume of data.
318        In practice, a trained CNN is actually an excellent tool for this purpose. It encodes
319    (downscales) the input with large feature volume into data with a much smaller size in the so-
320    called latent space (*i.e.*, the output of the layer before the output layer) while equal predictability
321    for the targeted events. This feature has been used in developing various generative DL
322    algorithms from variational autoencoder or VAE to different generative adversarial networks or
323    GANs (*e.g.*, Forest, 2019). Therefore, the trained HazeNet for Beijing and Shanghai have been
324    used in this study to produce data with reduced size suitable for clustering (Fig. 6; see also
325    Appendix C). The new sample-feature set with a size of 14975×512 produced from this
326    procedure was then used in cluster analysis.



Figure 6. A diagram of the cluster analysis procedure. Here 96, 64, and 9 represent the number of
longitudinal, latitudinal grids, and number of features (variables), or the size of the input feature volume
of a trained HazeNet for Beijing cases, while 512 is the size of the output from the second dense layer of
HazeNet or the new feature volume.

333     In order to provide useful information for understanding the performance of the trained
334     networks, the clustering has been performed for each of the prediction outcomes rather than just
335     haze versus non-haze events (Appendix C). In this configuration, haze associated regimes are
336     represented by derived clusters of TP plus FN outcomes, while non-haze regimes by those of TN
337     plus FP. Since the clusters were actually derived using the indices of samples as the record for
338     members, the actual feature maps of the members in any cluster thus can be conveniently
339     retrieved then used to identify the representative regimes in terms of combined 9 meteorological
340     and hydrological features of various prediction outcomes or haze versus non-haze events. Here
341     the clustering results have been analyzed using the feature maps in both normalized (machine
342     native) and unnormalized (original reanalysis data) format. The characteristics of various
343     regimes can be easily identified from the former as they represent anomalies to climatological
344     means. An added benefit is to advance the understanding of the performance of the trained
345     networks. The analysis using the latter maps aims to better appreciate the conventional regional
346     and local meteorological and hydrological patterns associated with various regimes. The feature
347     maps used in both analyses have been averaged across each cluster for clarity.

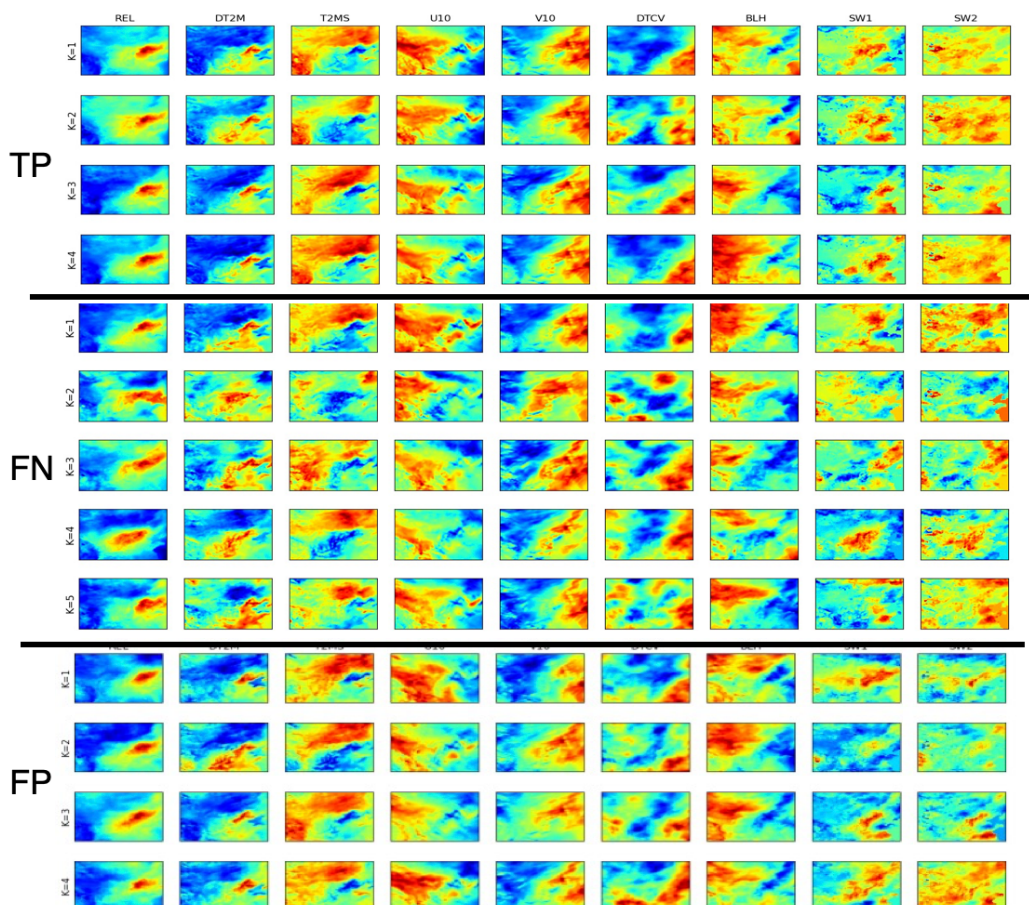### 4.1 Results based on normalized feature maps

349     As shown in Figure 7, the 4 clusters of true positive or TP in Beijing cases exhibit a clear
350     similarity in general feature patterns among themselves, differing only in rather minor details.
351     The differences between clusters are more evident in the daily change of column water vapor or
352     DTCV and in two soil water contents (SW1 and SW2). On the other hand, FN clusters (also
353     associated with haze events but missed in prediction) also display a clear similarity to the
354     patterns of TP clusters across most features except DTCV, SW1, and SW2.
355     Generally speaking, the common patterns in normalized feature maps shared by most clusters
356     associated with observed haze events (*i.e.*, TP plus FN outcomes) include an isolated positive
357     relative humidity (REL) center in the southeast region covering Beijing associated with mild
358     temperature variations (DT2M and T2MS) as well as zonal wind (U10) and lower boundary
359     layer height (BLH). Note that the mild daily temperature variation alongside lower BLH
360     indicates that the haze region is not experiencing drastic weather system change such as fronts
361     and likely covered by low cloud, hence the high REL can be easily formed. All these characters
362     reflect a stable regional weather conditions over the southeastern half of the domain where
363     targeted hazes occurred. They are also in a sharp contrast to the conditions in the northwestern
364     half of the domain as well as the conditions associated with non-haze events represented by TN
365     outcomes (Fig. S4).
366     Interestingly, the 4 FP (false alarm) clusters actually display a similarity in normalized
367     feature patterns to those of TP as those of FN (Fig. 7). In addition, despite an anticipated
368     diversity in feature patterns across TN clusters (Fig. S4), four of its clusters (i.e., 2, 5, 12, and 13)
369     exhibit a certain level of similarity to those of TP clusters. All these could offer an explanation
370     for the forecast errors made by the machine, *i.e.*, the machine could have simply been confused
371     by such similarities between certain FN and TN members, or between certain TP and FP
372     members. Nevertheless, these could also suggest an alternative reason behind the incorrect
373     forecasts. It is worth indicating again that meteorological or hydrological conditions are not the
374     only factors determining the occurrence of hazes. Other factors such as abnormal energy
375     consumption events or long-range transport of aerosols could all cause haze to occur even under
376     unfavorable weather and hydrological conditions. This could well be the reason for some of the
377     missing forecasts (FN outcomes) when haze occurred under unfavorable conditions, or for false

378    alarms (FP outcomes) when low aerosol events occurred even under a weather condition

379    favorable to haze. Future improvement of the skill could benefit from this knowledge.

380      The results of Shanghai are largely the same as in Beijing case (Fig S5 & S6).



381

382 **Figure 7**. Maps of 9 features in normalized format for 4 clusters of true positive or TP outcome, 5 clusters

383 of false negative or FN outcome, and 4 clusters of false positive or FP outcome. Here TP plus FN = haze

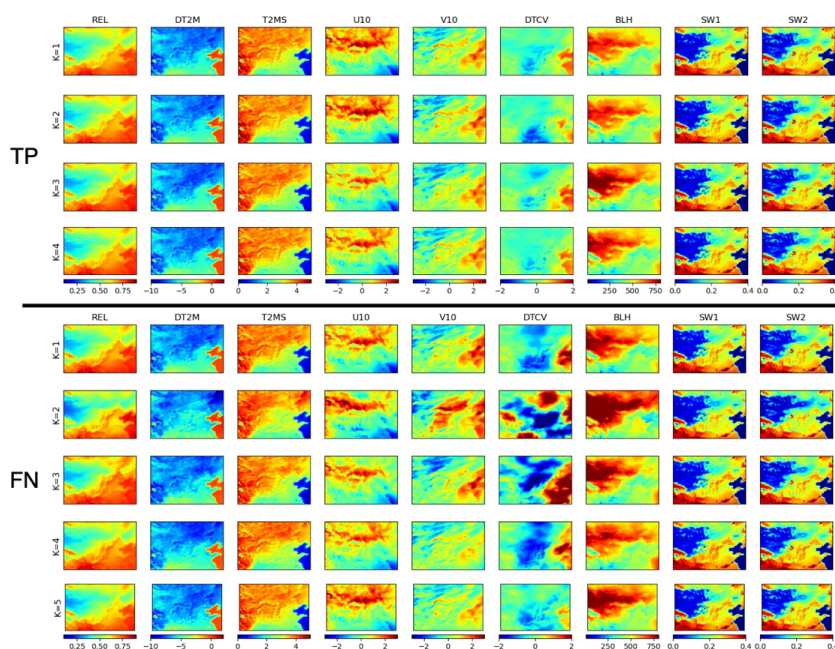384 events. Results shown are cluster averages for Beijing cases.

385 **4.2 Results based on original unnormalized feature maps**

386      Utilizing feature maps in their original unnormalized format represented by actual physical

387 quantities could provide a convenience to appreciate the conventional regional and local

388 meteorological and hydrological patterns associated with various events. Note that the visual

389 differences between unnormalized feature maps particularly in cluster-mean format might be

390 subtle for bare eyes to recognize.

391      For haze events in Beijing (*i.e.*, TP and FN outcomes; Fig. 8), the associated cluster-mean

392 regional meteorological and hydrological patterns of most features except DTCV contain two

393 regions with sharply contrasting quantities, roughly separated by a line linking the southwest and

394 northeast corner of the domain, likely due to the nature of weather system besides meridional

395    variation of general climate. Beijing (at ~1/3 domain width from the east boundary and nearly
396    the north-south center) locates in the southeastern half of the domain. In comparison, as same as
397    shown in the previous analysis using normalized feature maps, the patterns of FN clusters share
398    many common patterns with those of TP clusters. Their differences are more evident in DTCV,
399    SW1, and SW2. In addition, cluster 5 of FN shows more diverse patterns than the rest. FP
400    clusters also display a similarity to those of TP clusters (Fig. S5), whereas TN clusters show
401    more visible differences particularly in patterns of meridional wind (V10) and daily change of
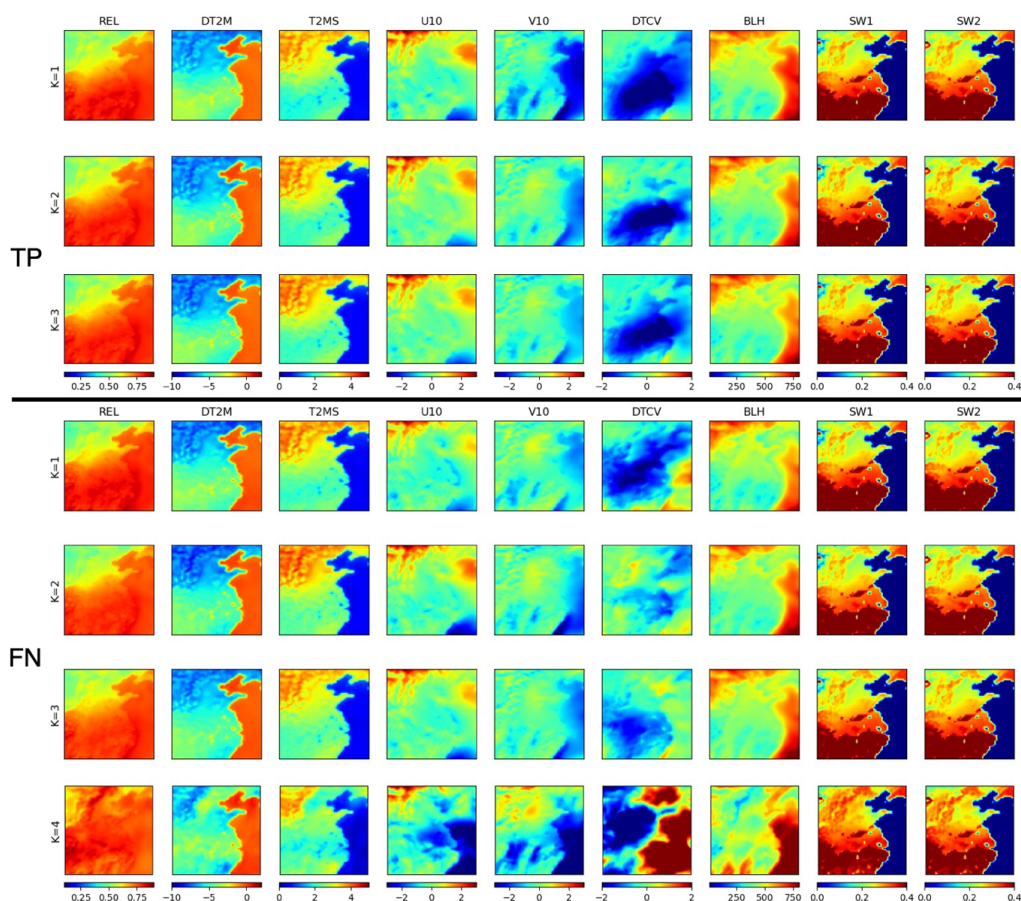402    column water vapor or DTCV (Fig. S6).



**Figure 8**. Feature maps associated with severe haze events in Beijing represented by 4 clusters of TP (4
top rows) and 5 clusters of FN (5 lower rows) predicted outcomes. Shown are cluster means of
unnormalized data of REL (ratio), DT2M and T2MS in degree, U10 and V10 in m/s, DTCV (kg/m$^2$),
BLH in meter, and SW1 and SW2 in kg/m$^2$.

408        The general regional meteorological and hydrological conditions during haze events in the
409    southeastern in comparison to the northwestern portion of the domain include a higher relative
410    humidity, lower variation of surface temperature, largely northward or northwestward wind,
411    lower planetary boundary layer height, and higher soil water content, and quantity wise these are
412    all in a sharp contrast to the situations in the other half of the domain. The visually recognized
413    cross-cluster differences of haze events mainly exist in DTCV patterns, represented by a strong
414    negative center in the middle of the domain with varying extent across different clusters. To a
415    less extent, patterns of surface wind V10 and U10 also offer some different characteristics
416    among various clusters particularly of FN clusters. Consistent to the analysis result using
417    normalized feature maps, all these indicate a stable weather condition over the southeastern half
418    of the domain for haze events in Beijing. It is known that dust can cause low visibility events in
419    Beijing. During dust seasons, the condition of the northwestern half of the domain, represented
420    by a dominant eastward wind and lower soil water content likely favors dust transport from

421  desert to Beijing. However, the details would need an in-depth analysis to examine since most
422  clusters having members rather well distributed through different months.
423      The cluster-means of 9 features for haze events (TP plus FN) versus non-haze (TN plus FP)
424  at the grid point of Beijing are also derived and listed in Table 1 for reference. Specifically, the
425  common local conditions associated with hazes in Beijing in comparison to those with non-haze
426  events include a higher humidity, less drastic variations in surface temperature, a northwestward
427  rather than southeastward wind, a lower planetary boundary layer height, and higher soil water
428  contents. Again, the most recognizable cross-cluster differences appear in DTCV, followed by
429  surface wind. In most of the local features, variabilities of FN clusters tend to be larger than
430  those of TP clusters. One interesting result of the local weather conditions shown in Table 1 is
431  that the cluster means of TN are sharply different than those of TP and FN, while the cluster
432  means of FP and those of TP+FN are likely to be statistically indifferent except for DTCV,
433  providing an evidence to support the assumption that FP outcomes might simply represent the
434  non-haze events caused by reasons other than weather and hydrological conditions.
435



436
437  **Figure 9.** The same as Figure 9 except for Shanghai with 3 clusters for TP and 4 for FN outcomes.
438

439    For the case of Shanghai, the general weather conditions associated haze events are likely
440    stable, with characters similar to the cases of Beijing (Fig. 9). Quantities of most feature patterns
441    display a sharply southeast versus northwest contrast. DTCV maps display a negative center over
442    a large area, its distribution and extent vary significantly among different clusters. The patterns
443    of soil water content in both soil layers exhibit a sharp meridional contrast, much higher in the
444    south part of the domain than in the north part, largely separated by the Yellow River. Local
445    quantities of all the features associated with haze events (TP plus FN) in Shanghai display clear
446    differences with those of non-haze prediction outcomes (TN) (Table 1). Similar to the cases of
447    Beijing, the cluster mean of the FP outcomes is statistically indifferent to those of haze (TP and
448    FN) than predicted non-haze (TN) events. Again, this result implies that even a weather pattern
449    favoring haze appeared and was correctly recognized by HazeNet, due to other factors such as
450    energy consumption variations, haze could still not to occur.

451    **Table 1**. Cluster means of features associated with haze events (TP and FN) in Beijing and Shanghai
452    versus means of all clusters of non-haze events of TN and FP, respectively. Number of cluster members
453    of each cluster are listed in bracket.

| Cluster | REL (0-1) | DT2 (ºC) | T2MS (ºC) | U10 (m/s) | V10 (m/s) | DTCV (kg/m²) | BLH (m) | SW1 (kg/m²) | SW2 (kg/m²) |
|---|---|---|---|---|---|---|---|---|---|
| Beijing | | | | | | | | | |
| TP1 (848) | 0.64 | -5.99 | 3.24 | -0.29 | 0.20 | 0.04 | 379.71 | 0.23 | 0.22 |
| TP2 (181) | 0.65 | -5.80 | 3.14 | -0.28 | 0.19 | 0.57 | 378.33 | 0.23 | 0.23 |
| TP3 (354) | 0.65 | -5.39 | 2.98 | -0.45 | 0.29 | 0.31 | 400.20 | 0.23 | 0.22 |
| TP4 (1208) | 0.64 | -5.82 | 3.18 | -0.34 | 0.28 | 0.27 | 381.28 | 0.23 | 0.22 |
| FN1 (157) | 0.66 | -5.83 | 3.16 | -0.43 | 0.34 | 0.15 | 379.91 | 0.23 | 0.21 |
| FN2 (13) | 0.65 | -5.05 | 2.98 | -0.52 | 0.48 | -1.88 | 422.35 | 0.23 | 0.22 |
| FN3 (29) | 0.69 | -5.90 | 3.05 | -0.41 | 0.36 | 0.99 | 393.52 | 0.24 | 0.23 |
| FN4 (86) | 0.64 | -5.64 | 3.02 | -0.19 | 0.11 | 0.10 | 420.49 | 0.23 | 0.22 |
| FN5 (223) | 0.60 | -6.56 | 3.45 | -0.14 | 0.11 | 0.01 | 449.48 | 0.23 | 0.22 |
| TN mean | 0.51 | -7.13 | 3.65 | 0.15 | -0.15 | 0.36 | 552.90 | 0.22 | 0.21 |
| FP mean | 0.65 | -5.84 | 3.15 | -0.35 | 0.25 | -0.26 | 386.27 | 0.24 | 0.23 |
| Shanghai | | | | | | | | | |
| TP1 (1228) | 0.81 | -3.44 | 1.79 | -0.16 | -0.55 | -2.25 | 415.59 | 0.35 | 0.35 |
| TP2 (135) | 0.81 | -3.10 | 1.71 | -0.12 | -0.66 | -2.08 | 422.04 | 0.36 | 0.36 |
| TP3 (689) | 0.81 | -2.95 | 1.59 | -0.17 | -1.28 | -2.29 | 472.74 | 0.36 | 0.35 |
| TP4 (355) | 0.81 | -3.52 | 1.82 | 0.03 | -0.57 | -2.74 | 411.96 | 0.35 | 0.35 |
| FN1 (102) | 0.82 | -3.33 | 1.80 | -0.67 | -0.36 | -0.14 | 409.55 | 0.35 | 0.35 |
| FN2 (113) | 0.80 | -3.64 | 1.84 | -0.34 | -0.51 | -1.21 | 423.09 | 0.35 | 0.34 |
| FN3 (370) | 0.80 | -3.47 | 1.80 | -0.41 | -0.42 | -0.84 | 421.36 | 0.35 | 0.35 |
| FN4 (7) | 0.80 | -2.82 | 1.39 | -1.19 | -2.18 | 3.63 | 596.53 | 0.36 | 0.36 |
| TN mean | 0.77 | -3.29 | 1.57 | -2.86 | 1.40 | 0.62 | 739.75 | 0.31 | 0.32 |
| FP mean | 0.82 | -3.26 | 1.71 | -0.48 | -0.85 | -2.26 | 438.55 | 0.35 | 0.35 |

## 5 Summary and Conclusions

Following an earlier preliminary attempt for hazes in Singapore, a deep convolutional neural network containing more than 20 million parameters, namely HazeNet, has been further developed to test forecasting the occurrence of severe haze events during 1979-2019 in two metropolitans of Asia, Beijing and Shanghai. By training the machine to recognize regional patterns of meteorological and hydrological features associated with haze events, the study would advance our knowledge about this still poorly known environmental extreme. The deep CNN has been trained in a supervised learning procedure using the time sequential maps of up to 16 meteorological and hydrological variables or features as inputs and surface visibility observations as the labels.

Even with a rather limited samples (14,975), the trained machine has displayed a reasonable performance measured by commonly adopted validation metrics. Its performance is clearly better during months with high haze frequency, *i.e.*, all months except dusty April and May in Beijing and from late autumn through entire winter in Shanghai. Relatively larger spatial patterns appear to be more effective than the smaller ones to influence the performance of forecasting. On the other hand, in-depth analysis on performance results has also indicated certain limitations of current approach of solely using meteorological and hydrological data in performing forecast.

The trained machine has also been used to examine the sensitivity of the CNN to various input features and thus to identify then remove features ineffective to the performance of the machine. In addition, in order to further categorize typical regional weather and hydrological patterns associated with severe haze versus non-haze events, an unsupervised cluster analysis has been subsequently conducted, benefited from using features with greatly reduced dimensionality produced by the trained machine.

The cluster analysis has, arguably for the first time, successfully categorized major regional meteorological and hydrological patterns associated with severe haze and non-haze events in Beijing and Shanghai into a limited number of representative groups, with the typical feature patterns of these clustered groups derived. It has found that the typical weather and hydrological regimes of haze events in Beijing and Shanghai are rather stable conditions, represented by increasing relative humidity, low planetary boundary layer, mild daily temperature change that likely associated with low cloud cover over the haze occurring regions, The result has further revealed a rather strong similarity between the meteorological and hydrological patterns associated with haze events and those with either false alarm or missing forecast prediction outcomes, implying that factors other than meteorological and hydrological ones such as energy consumption variations, long range transport of aerosols, or beyond, could cause haze events to occur even under unfavorite weather conditions.

Due to the exploratory nature of this specific effort, several aspects could be further optimized including the rather arbitrary though statistically meaningful labeling. Also, an in-depth analysis on weather regimes exceeds the extent of this paper. Nevertheless, this study has demonstrated the potential of applying deep CNNs with extensive multi-dimensional and time sequential environmental images to advance our understandings about poorly known environmental and weather extremes. The methodology, results alongside experience obtained from this study could benefit future improvement of the skills. Besides, the trained machines can be used in many other types of machine learning and deep learning applications as partially demonstrated here.

## Appendix A. Performance metrics

Several commonly used performance metrics have been used in this study. They are largely derived based on the so-called confusion matrix (e.g., Swets, 1988) as defined in the following Table A.

**Table A**. Confusion matrix for measuring the prediction outcomes of a given class.

| | | Observed | |
|---|---|---|---|
| | | *Positive* | *Negative* |
| *Predicted* | *Positive* | True Positive or TP | False Positive or FP |
| | *Negative* | False Negative or FN | True Negative or TN |

Here, *positive* or *negative* is referring to the outcome of a given event or class in the classification, *e.g.*, severe haze or non-haze events. Hence, the prediction outcome TP is a correct forecast of a severe haze while TN a correct forecast of a non-haze event, FP represents a false alarm, and FN a missing forecast. The context of outcomes changes when the designated class is switched. The major performance metrics used in this paper include:

$$accuracy = \frac{TP+TN}{N} \tag{A1}$$

$$precision = \frac{TP}{TP+FP} \tag{A2}$$

$$recall = \frac{TP}{TP+FN} \tag{A3}$$

$$F1\ score = 2 \cdot \frac{precision \cdot recall}{precision+recall} \tag{A4}$$

$$ETS = \frac{TP-Hit_{random}}{TP+FP+FN-Hit_{random}}; \tag{A5a}$$

$$where: \quad Hit_{random} = \frac{(TP+FN) \cdot (TP+FP)}{N} \tag{A5b}$$

$$HSS = \frac{2 \cdot (TP \cdot TN-FP \cdot FN)}{(TP+FP) \cdot (FP+TN)+(TP+FN) \cdot (TP+TN)} \tag{A6}$$

Note that *accuracy* has the same value for all the classes and thus is a good metrics for the overall classification. Values of all the other metrics differ depending on the referred specific class. Here, *F1 score* is the F-score with $\beta = 1$ (van Rijsbergen, 1974), *ETS* represents equitable threat score (or Gilbert skill score; Gilbert, 1884; range = [-1/3, 1]), *HSS* represents Heidke skill score (Heidke, 1926; range = [-∞,1]), and *N* is the number of total outcomes.

## Appendix B. Examining the network's sensitivity to features using trained machine

A method has been adopted in this study to use a trained machine from basic training to examine the sensitivity of the network to a random perturbation applied to the values of different features. The saved machine contains all the coefficients in different network layers and can be used to predict output from any of these layers using same input features for training or validation. The sensitivity of the network to a given feature is determined by comparing the prediction using input feature maps containing randomly perturbation applied to the map of this feature with the prediction using original input feature maps, and measured by the content loss between these two predictions, with *img1* with *MxN* pixels as the unperturbed and *img2* as perturbed network output:

$$Content\ Loss = \frac{1}{M \times N} \sum_{i,j}^{M,N} (img1_{i,j} - img2_{i,j})^2 \tag{B1}$$

The perturbation is applied as random patch with addition of -0.2 or 0.2 to 10% of the pixels of the input map of the targeted feature in each sample while maps of all the other features remain unperturbed. To reduce the workload, only validation input set corresponding to the class-1 events (about 1020 samples) are used. Therefore, the sensitivity tested here is actually the sensitivity of the network to a given feature in predicting class-1 events. To preserve the spatial information of the perturbation field, the output of the 9[th] layer, or the MaxPooling layer following the second convolutional layer (Fig. 1) is used as the prediction. It has a size of (15, 31, 92) for Beijing cases and (15, 15, 92) for Shanghai cases when a kernel size of 20x20 is adopted. A higher content loss represents that the performance of the network is more sensitive to the variations in value of this feature.

## Appendix C. Cluster analysis

535

536     The cluster analysis of this study was conducted in the following three steps (see also Fig. 6).

537     **(i)** Firstly, the trained and saved HazeNet for both Beijing and Shanghai cases with 9 input features have been
538 used to perform prediction using the entire 14975 input samples in original raw data format, *i.e.*, with a feature
539 volume size of 96x64x9 for Beijing and 64x64x9 for Shanghai for each sample. The prediction results were then
540 summarized into various outcomes, *e.g.,* as true positive (TP), true negative (TN), false positive (FP), or false
541 negative (FN) in referring to the haze class. In the meantime, the output of the second dense layer just before the
542 output layer or latent space (see Fig. 1 & Fig. 6) were further used to form the new data of each sample with reduced
543 feature volume of 512. This new dataset with 14075 samples and 512 feature volume were ready for clustering.

544     **(ii)** The second step is to actually perform clustering using the new data with reduced size resulted from the
545 previous step. For this purpose, it should be conducted separately for different types of samples or events, *e.g.*,
546 categorizing all the samples for haze into characteristic groups with similarity and same for non-haze events. In
547 order to provide additional information to further the understanding of the network's performance, the clustering
548 was actually conducted for different prediction outcomes, by taking corresponding samples from the new dataset. In
549 this case, TP plus FN would lead to haze events, and TN plus FP to non-haze events. The clustering calculations
550 were done by directly using the k-mean (Steinhaus, 1957) function of scikit-learn library (https://scikit-
551 learn.org/stable/modules/clustering.html#clustering). For Beijing cases, the trained machine with 9 features
552 produced 2591 TP, 11368 TN, 508 FP, and 508 FN outcomes, and 2407 TP, 11484 TN, 492 FP, and 592 FN for
553 Shanghai. The cluster analysis was performed separately for each of these outcomes in an unsupervised learning
554 procedure to let the machine to categorize corresponding samples into groups based on similarities among them. In
555 practice, similarity is judged by the so-called inertia for a cluster with members of $x_i$ and mean of $\mu$:

556 $$inertia = \sum_i^N (\|x_i - \mu\|)^2 \qquad \text{(C1)}$$

557 The clustering is to seek a grouping with minimized inertia within each cluster. The overall measure is the
558 summation inertia that decreases almost exponentially with the increase of number of clusters. In practice, the
559 cluster analysis was first tested with various given number of clusters ranging from 1 to 100, to examine the values
560 alongside decay of the inertia. This provided a base to identify the smallest possible number of cluster centers with
561 reasonably low inertia in actual cluster analysis. This has actually been decided by using square root of the inertia
562 weighted by the number of samples to put the varying number of samples across various outcomes in consideration.
563 An optimized number of clusters was chosen with a weighted inertia lower than 1/e of that of the single cluster case.
564 For TN, due to the large sample number, this criterion was set to be half of 1/e. As a result, the optimized numbers
565 of clusters for TP, FN, FP, and TN outcomes are 4, 5, 4, and 15 for Beijing and 4, 4, 3, and 10 for Shanghai,
566 respectively,

567     **(iii)** The members of each cluster derived from (ii) were recorded by the actual sample indices with date
568 attribute. Therefore, actual samples of input data grouped into various clusters can be thus conveniently identified
569 with corresponding feature maps retrieved, either in the format of normalized or unnormalized (*i.e.*, in original
570 quantity as in reanalysis dataset), and used for further analyses. In practice, cluster-averaged maps for various
571 features were performed beforehand.

## Code and data availability

573     The Python script for network architecture, training and validation is rather straightforward and simple,
574 basically consisting of directly adopted function calls from Keras interface library (https://github.com/keras-
575 team/keras) with TensorFlow-GPU (https://www.tensorflow.org) as backend, or from scikit learn library
576 (https://scikit-learn.org/). All the data used for analyses are publicly available as indicated in the
577 Acknowledgements.

## Competing interests

579 The author declares that he has no conflict of interest.

580

## Acknowledgements

## References

Chan, C. K. and Yao, X.: Air pollution in mega cities in China, Atmos. Environ., 42, 1-42, 2008.

Chattopadhyay, A., Nabizadeh, E. and Hassanzadeh, P.: Analog forecasting of extreme-causing weather patterns using deep learning. *J. Adv. Modeling Earth Sys*., 12, e2019MS001958. Doi:/10.1029/2019MS001958, 2020.

Forest, D.: *Generative Deep Learning*, O'Reilly Media, Inc., 2019.

Gagne, D., Haupt, S. and Nychka, D.: Interpretable deep learning for spatial analysis of severe hailstorms. *Mon. Wea. Rev*., 147, 2827–2845, Doi:/10.1175/MWR-D-18-0316.1, 2019.

Gilbert, G. K.: Finley's tornado predictions, *Amer. Meteor. J*., 1, 166–172, 1884.

Goodfellow, I., Bengio, Y. and Courville, A.: *Deep Learning*, MIT Press, 800pp., 2017.

Grover, A. Kapoor, A. and Horvitz, E.: A deep hybrid model for weather forecasting, *Proc. 21st ACM SIGKDD Intern'l Conf. KDD*, p.379-386, Sydney, Australia, August 10, 2015. ACM. ISBN 978-1-4503-3664-2/15/08. Doi:10.1145/2783258.2783275, 2016.

He, K., Zhang, X., Ren, S. and Sun, J.: Deep residual learning for image recognition, arXiv:1512.03385, 2015.

Heidke, P.: Calculation of the success and goodness of strong wind forecasts in the storm warning service, *Geogr. Ann. Stockholm*, 8, 301–349, 1926.

Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., De Chiara, G., Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes, M., Geer, A., Haimberger, L., Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S., Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., de Rosnay, P., Rozum, I., Vamborg, F., Villaume, S. and Thépaut, J.-N.: The ERA5 global reanalysis, *Q. J. R. Meteorol. Soc*., 146, 1999-2049, 2020.

Ioffe, S. and Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift, arXiv:1502.03167, 2015.

Jiang, G.-Q., Xu, J. and Wei, J.: A deep learning algorithm of neural network for the parameterization o typhoon-ocean feedback in typhoon forecast models, *Geophys. Res. Lett*., 45, https://doi.org/10.1002/2018GL077004, 2018.

Kiehl, J. T. and Briegleb, B. P.: The relative roles of sulfate aerosols and greenhouse gases in climate forcing, *Science*, 260, 311-314, 1993.

Kurth, T., Treichler, S., Romero, J., Mudigonda, M., Luehr, N., Phillips, E., Mahesh, A., Matheson, M., Deslippe, J., Fatica, M., Prabhat and Houston, M.: Exascale deep learning for climate analytics, arXiv:1810.01993, 2018.

Lagerquist, R., McGovern, A. and Gagne II, D.: Deep learning for spatially explicit prediction of synoptic-scale fronts. *Wea. Forecasting*, 34, 1137–1160. Doi:10.1175/WAF-D-18-0183.1, 2019.

626    LeCun, Y., Bengio, Y. and Hinton, G.: Depp learning, *Nature*, 521, 436-444,
627        doi:10.1038/nature14539, 2015.
628    Lee, H.-H., Iraqui, O. and Wang, C.: The impacts of future fuel consumption on regional air
629        quality in Southeast Asia, *Sci. Rep*., 9:2648, doi:10.1038/s41598-019-39131-3, 2019.
630    Lee, H.-H., Iraqui, O., Gu, Y., Yim, H.-L. S., Chulakadabba. A., Tonks, A. Y. M., Yang, Z. and
631        Wang, C.: Impacts of air pollutants from fire and non-fire emissions on the regional air
632        quality in Southeast Asia, *Atmos. Chem. Phys*., 18, 6141–6156, doi:10.5194/acp-18-6141-
633        2018, 2018.
634    Lee, H.-H., Bar-Or, R. and Wang, C.: Biomass Burning Aerosols and the Low Visibility Events
635        in Southeast Asia, *Atmos. Chem. Phys*., 17, 965-980, doi:10.5194/acp-17-965-2017, 2017.
636    Lin, Y., Wijedasa, L. S. and Chisholm, R. A.: Singapore's willingness to pay for mitigation of
637        transboundary forest-fire haze from Indonesia, *Environ. Res. Lett*., 12, 024017,
638        doi:10.1088/1748-9326/aa5cf6, 2016.
639    Liu, M., Huang, Y., Ma, Z., Jin, Z., Liu, X., Wang, H., Liu, Y., Wang, J., Jantunen, M., Bi, J. and
640        Kinney, P. L.: Spatial and temporal trends in the mortality burden of air pollution in China:
641        2004-2012, Environ. Int., 98, 75-81, 2017.
642    Liu, Y., Racah, E., Prabhat, Correa, J., Khosrowshahi, A., Lavers, D., Kunkel, K., Wehner, M.
643        and Collins, W.: Application of deep convolutional neural networks for detecting extreme
644        weather in climate datasets. arXiv:1605.01156, 2016.
645    McGovern, A., Lagerquist, R., Gagne II, D. J., Jergensen, G. E., ElmLMore, K. L., Homeyer, C.
646        R. and Smith, T.: Making the black box more transparent: Understanding the physical
647        implications of machine learning, *Bull. Amer. Meteor. Soc*., 100, 2175-2199, 2019.
648    Ronneberger, O., Fischer, P. and Brox, T.: U-Net: Convolutional networks for biomedical image
649        segmentation, arXiv:1505.04597, 2015.
650    Shi, X., Chen, Z., Wang, H. and Yeung, D.-Y.: Convolutional LSTM network: A machine
651        learning approach for precipitation nowcasting, arXiv:1506.04214, 2015.
652    Silva, R. A., West, J. J., Zhang, Y., Anenberg, S. C., Lamarque, J.-F., Shindell, D. T., Collins,
653        W. J., Dalsoren, S., Faluvegi, G., Folberth, G., Horowitz, L. W., Nagashima, T., Naik, V.,
654        Rumbold, S., Skeie, R., Sudo, K., Takemura, T., Bergmann, D., Cameron-Smith, P., Cionni,
655        I., Doherty, R. M., Eyring, V., Josse, B., MacKenzie, I. A., Plummer, D., Righi, M.,
656        Stevenson, D. S., Strode, S. Szopa, S. and Zeng, G.: Global premature mortality due to
657        anthropogenic outdoor air pollution and the contribution of past climate change, *Environ.*
658        *Res. Lett*., 8, 034005, doi:10.1088/1748-9326/8/3/034005, 2013.
659    Simonyan, K. and Zisserman, A.: Very deep convolutional networks for large-scale image
660        recognition, arXiv:1409.1556, 2015.
661    Smith, A., Lott, N. and Vose, R.: The integrated surface database: Recent developments and
662        partnerships, *Bull. Ameri. Meteorol. Soc*., 92, 704-708, doi:10.1175/2011BAMS3015.1,
663        2011.
664    Steinhaus, H.: Sur la division des corps matériels en parties, *Bull. Acad. Polon. Sci*., 4, 801–804,
665        1957.
666    Swets, J.: Measuring the accuracy of diagnostic systems, *Science*, 240, 1285–1293, 1988.
667    Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. and Wojna, Z.: Rethinking the inception
668        architecture for computer vision, arXiv:1512.00567, 2015.
669    van Rijsbergen, C., 1974: Foundation of evaluation, *J. Documentation*, 30, 365–373.
670    Wang, C.: Exploiting deep learning in forecasting the occurrence of severe haze in Southeast
671        Asia, arXiv:2003.05763, 2020.

672  Weyn, J. A., Durran, D. R. and Caruana, R.: Improving data-driven global weather prediction
673      using deep convolutional neural networks on a cubed sphere, *J. Adv. Modeling Earth Sys.*,
674      e2020MS002109, https://doi.org/10.1029/2020MS002109, 2020.
675
676