

29 April 2026

S2G-DI: A Deep Learning Framework for Generating High-Resolution Wind Gust Fields from Sparse Mesonet Observations

Harish Baki¹, Maximilian Pierzyna², Sukanta Basu^{1,3}

1. Atmospheric Sciences Research Center University at Albany

2. Department of Geoscience and Remote Sensing Delft University of Technology

3. Department of Environmental & Sustainable Engineering University at Albany

Abstract

The growing deployment of surface observational networks (so-called Mesonets) has the potential to transform many sectors, such as agriculture, water resource management, renewable energy assessment, disaster risk reduction, power outage prediction, and high-impact weather monitoring. However, network observations are sparse, whereas these downstream applications typically require gridded data. Traditional interpolation methods, such as Barnes interpolation, do not perform well in complex terrain. Therefore, in this study, we propose a deep learning-based framework, called Sparse-to-Gridded Deep Interpolation (S2G-DI), to generate high-resolution spatially continuous fields from sparse surface observations. We focus on wind gusts, which are challenging to interpolate due to their highly localized and transient nature, but are critical for assessing weather-related hazards. We experimented with three deep learning architectures, employing convolutional neural networks and transformers with increasing complexity, and our sensitivity analysis shows that a U-Net architecture with meteorological and topographic inputs significantly outperforms Barnes interpolation, reducing RMSE by over 30% and better preserving fine-scale features. The model remains robust to missing data and generalizes well even with limited training data. Evaluations during extreme wind gust events confirm that the framework captures both spatial structure and intensity more reliably than traditional methods. However, for some of these cases, there is still room for improvement for the S2G-DI framework.

Keywords

atmospheric sciences, Barnes interpolation, deep learning, mesonet, meteorology, S2G-DI, wind gust

1 **S2G Deep Interpolation: A Deep Learning Framework for Generating**
2 **High-Resolution Wind Gust Fields from Sparse Mesonet Observations**

3 Harish Baki^a, Maximilian Pierzyna^b, Sukanta Basu^{a,c}

4 ^a*Atmospheric Sciences Research Center, University at Albany, New York, USA*

5 ^b*Department of Geoscience and Remote Sensing, Delft University of Technology, Delft,*
6 *Netherlands*

7 ^c*Department of Environmental & Sustainable Engineering, University at Albany, New York, USA*

8 *Corresponding author: Harish Baki, hbaki@albany.edu*

9 ABSTRACT: The growing deployment of surface observational networks (so-called Mesonets) has
10 the potential to transform many sectors, such as agriculture, water resource management, renewable
11 energy assessment, disaster risk reduction, power outage prediction, and high-impact weather
12 monitoring. However, network observations are sparse, whereas these downstream applications
13 typically require gridded data. Traditional interpolation methods, such as Barnes interpolation, do
14 not perform well in complex terrain. Therefore, in this study, we propose a deep learning-based
15 framework, called Sparse-to-Gridded Deep Interpolation (S2G-DI), to generate high-resolution
16 spatially continuous fields from sparse surface observations. We focus on wind gusts, which
17 are challenging to interpolate due to their highly localized and transient nature, but are critical
18 for assessing weather-related hazards. We experimented with three deep learning architectures,
19 employing convolutional neural networks and transformers with increasing complexity, and our
20 sensitivity analysis shows that a U-Net architecture with meteorological and topographic inputs
21 significantly outperforms Barnes interpolation, reducing RMSE by over 30% and better preserving
22 fine-scale features. The model remains robust to missing data and generalizes well even with
23 limited training data. Evaluations during extreme wind gust events confirm that the framework
24 captures both spatial structure and intensity more reliably than traditional methods. However, for
25 some of these cases, there is still room for improvement for the S2G-DI framework.

26 **1. Introduction**

27 High-resolution gridded meteorological data is essential for a wide range of applications across
28 science and industry. It supports land surface and agricultural studies, including agricultural
29 decision-making (Subedi et al. 2025), soil moisture modeling (Peraza et al. 2025), and snowpack
30 distribution modeling (Le Toumelin et al. 2023). It plays a critical role in natural hazard and disaster
31 management, such as forest fire monitoring and prediction (Chen et al. 2023), and in hydrological
32 extremes, including floods (Amponsah et al. 2018) and droughts (Peng et al. 2020). In atmospheric
33 science, it is used for monitoring and characterization of extreme weather events, such as cyclones
34 (Knapp and Kossin 2007), heatwaves (Hu et al. 2023), and thunderstorms (Houston et al. 2015).
35 Other applications include precipitation forecasting (An et al. 2025), evapotranspiration estimation
36 (Blankenau et al. 2020), and air pollution analysis (Venter et al. 2024). High-resolution data is
37 also crucial for wind resource assessment (Christiansen et al. 2006), as well as for water resource
38 management (Müller et al. 2021). The majority of these applications rely on key meteorological
39 variables, including precipitation, temperature, wind speed and direction, wind gust, humidity,
40 solar radiation, and soil moisture.

41 In the Continental United States (CONUS), the expansion of mesoscale observational networks
42 (mesonets) has been a key step toward addressing the demand for dense and continuous meteorolog-
43 ical measurements, with a total of 27 statewide mesonets established to date (Campbell Scientific,
44 Inc. 2023). The mesonets, including the New York State Mesonet (NYSM), typically have a
45 station spacing of about 30 km and are designed to capture mesoscale atmospheric phenomena
46 (Mahmood et al. 2017). These networks collect raw measurements of variables at 2-3 s intervals,
47 with quality-controlled products made available to the public at 5-10 min resolution. However,
48 despite their high temporal fidelity, the observations are inherently point-based (ungridded) and
49 spatially sparse. In particular, the NYSM network provides measurements at discrete locations,
50 whereas the target analysis fields aim to resolve variability at kilometer-scale resolutions (e.g.,
51 2–3 km) in modern high-resolution systems such as Real-Time Mesoscale Analysis (RTMA) and
52 High-Resolution Rapid Refresh (HRRR). Addressing this spatial sparsity requires transforming
53 point observations into gridded fields, also known as *objective analysis*, which is a longstanding
54 challenge in atmospheric sciences.

55 The objective analysis involves interpolating irregularly spaced measurements onto a structured
56 grid. Many empirical and statistical techniques have been developed for this purpose, including the
57 Cressman scheme (Cressman 1959), Barnes interpolation (Barnes 1964), kriging (Matheron 1967),
58 multiquadric interpolation (Hardy 1971), Shepard's method (Shepard 1968), and Gandin-based
59 statistical approaches (Gandin 1963). Although most of these techniques were originally proposed
60 in the mid 20th century, they are still employed in several recent studies, such as the Barnes
61 interpolation in the MicroMet framework (Liston and Elder 2006), kriging in the creation of the E-
62 OBS dataset (Haylock et al. 2008), inverse distance-weighted interpolation in deep learning-based
63 wind gust nowcasting (Xiao et al. 2023), and Gandin-based techniques for precipitation analysis
64 (Chen et al. 2008). While effective for some applications, these techniques make simplifying
65 assumptions and often fail to capture the complex geographical and topographical effects present
66 in real-world data.

67 These early objective analysis techniques laid the foundation for sophisticated data assimilation
68 (DA) frameworks, which integrate observations with numerical weather prediction (NWP) models
69 to produce physically and dynamically consistent gridded fields of meteorological variables (Kalnay
70 2003). Modern DA frameworks, such as the two/three/four-dimensional variational assimilation
71 (2D/3D/4D-Var), the ensemble Kalman filter (EnKF), and hybrid approaches can handle a wide
72 range of heterogeneous observations, including conventional surface and upper-air observations
73 from weather balloons, radiosondes, surface stations, aircraft, and buoys, as well as remotely
74 sensed data from radar, satellite radiances, and GPS-based retrievals (Barker et al. 2012). These
75 frameworks operate in close association with specific NWP models, such as the Weather Research
76 and Forecasting (WRF), Integrated Forecast System (IFS), Global Forecast System (GFS), since
77 they require model-based background fields, ensemble forecasts, and error covariance information
78 for continuous operation. Advances in NWP-DA systems have enabled the development of state-
79 of-the-art gridded meteorological products (Mankin et al. 2025), including the 5th generation
80 European Center for Medium-range Weather Forecast (ECMWF) ReAnalysis (ERA5) (Hersbach
81 et al. 2020), Copernicus European Regional ReAnalysis (CERRA) (Ridal et al. 2024), CONUS404
82 (Rasmussen et al. 2023), and Real-Time Mesoscale Analysis (RTMA) (De Pondecia et al. 2011).
83 Despite their demonstrated accuracy, however, these systems are computationally expensive due
84 to their reliance on large-scale integrations and ensemble techniques, which restrict operational

85 update frequencies to hourly or longer intervals. This trade-off between accuracy and computational
86 expense motivates the exploration of alternative, data-driven approaches, without the reliance on
87 complex DA and NWP model integrations to transform station-based observations into gridded
88 fields.

89 Deep learning (DL) methods have recently emerged as a promising alternative for generating
90 gridded meteorological fields from sparse observations. A detailed review of relevant DL-based
91 approaches is presented in Appendix A. In this work, we introduce the Sparse to Gridded Deep
92 Learning-based Interpolation (S2G-DI) framework for directly reconstructing high-resolution grid-
93 ded fields from sparse NYSM observations, covering the entire state of New York. Among the
94 key observed meteorological variables, we demonstrate the capabilities of the S2G-DI framework
95 with wind gusts as an illustrative example. Wind gusts are chosen because they are high-impact
96 weather events leading to severe societal and economic consequences (Kahl 2020) and are complex
97 to represent in NWP-DA systems (Sheridan 2018).

98 The proposed S2G-DI framework is designed to produce gridded wind gust fields directly
99 from sparse mesonet observations in a computationally efficient manner. It follows a decoupled
100 training–inference strategy: during training, the model learns the mapping between sparse and
101 gridded fields using a gridded NWP-based product (RTMA in this study) as reference, while
102 inference requires only real-time station observations. This design eliminates amplitude and
103 displacement errors (Gilleland et al. 2009) between NWP products and observations, enabling
104 gridded fields to be generated at the native frequency of the mesonet (5 min for NYSM) with fine
105 spatial resolution (2.5 km) and minimal computational cost. This approach eliminates the need
106 for NWP-DA-based products during inference, supporting near real-time analysis. Such analysis
107 fields enable rapid forecast updates, which are crucial for wind gusts. This study has three main
108 objectives: (i) to develop the S2G-DI framework for producing gridded wind gust fields from sparse
109 mesonet observations, (ii) to assess its sensitivity to model architecture, input configurations, data
110 availability, and applicability to other meteorological variables, and (iii) to evaluate its ability
111 to capture extreme gust events in real-world conditions. Performance is benchmarked against the
112 Barnes interpolation method (Barnes 1964), a widely used operational objective analysis technique.

113 The remainder of the paper is organized as follows. Section 2 describes the S2G-DI framework.
114 Section 3 presents the study area, observational networks, and simulated and observed data used

115 for model training and inference. Section 4 presents deep learning methodology and model
 116 performance evaluation methods. Section 5 present results for the sensitivity analysis on identifying
 117 optimized framework, while Sections 6, 7 and 8 present the evaluation of the optimized framework
 118 with real observations as well as an independent source data, and framework’s generalization across
 119 meteorological variables, respectively. Section 9 concludes the paper and discusses directions for
 120 future work.

121 2. S2G-DI framework

122 The S2G-DI framework is illustrated in Fig. 1, which follows a decoupled training–inference
 123 approach through two phases. The first phase is the *training/validation/testing phase* (Fig. 1(a–d)),
 124 where a deep learning (DL) model is trained on a reference high-resolution gridded meteorological
 125 dataset (acting as a proxy, RTMA in this study), to learn the reconstruction of continuous fields from
 126 sparsified versions of themselves, assuming spatial similarity between the gridded and observational
 127 data. This phase involves four main steps:

- 128 (a) A high-fidelity 2D gridded field of the target variable of interest at a time instance is selected
 129 from the reference proxy dataset $X_{ref}^{(hf)} \in \mathbb{R}^{n \times m}$ (Fig. 1(a)).
- 130 (b) Using a set of known station coordinates, values from the $X_{ref}^{(hf)}$ are extracted at those station
 131 locations, resulting in a sparse 1D vector $X^{(s)} \in \mathbb{R}^s$ (Fig. 1(b)). This mimics the observed
 132 values from sparse stations.
- 133 (c) The $X^{(s)}$ is then nearest-neighbor interpolated onto the full latitude–longitude grid of the
 134 $X_{ref}^{(hf)}$. This produces a low-fidelity 2D gridded field $X^{(lf)} \in \mathbb{R}^{n \times m}$ with the same shape and
 135 coordinates as the $X_{ref}^{(hf)}$, but introduces discontinuities and artifacts that degrade fidelity
 136 (Fig. 1(c)). This intermediate gridding step enables the formulation of the problem within an
 137 image-to-image learning framework. Although the original observations are spatially sparse
 138 point measurements, representing them on a structured grid allows the use of convolutional
 139 architectures, which are well suited for learning spatial dependencies.
- 140 (d) This $X^{(lf)}$ is provided as input to a deep learning model, which outputs a high-fidelity 2D
 141 prediction $X_{DL}^{(hf)} \in \mathbb{R}^{n \times m}$, trained by minimizing the loss relative to the reference field $X_{ref}^{(hf)}$
 142 (Fig. 1(d)).

143 For clarity, the framework is described in a variable-agnostic manner so that X can represent any
144 variable of interest. Here, m and n represent the grid points along latitude-longitude directions,
145 and s represents the number of stations. The pipeline is $X_{ref}^{(hf)} \rightarrow X^{(s)} \rightarrow X^{(lf)} \rightarrow X_{DL}^{(hf)} \Leftrightarrow X_{ref}^{(hf)}$.
146 This procedure enables the DL model to learn a mapping from sparse, discontinuous inputs to
147 spatially continuous and physically consistent outputs.

148 The second phase is the *inference phase*, during which the real sparse observations from mesonet
149 (e.g., NYSM) stations are used to obtain a continuous gridded field of the same, utilizing the
150 trained DL model from *training phase* (Fig. 1(b*–d*)). This involves three steps (b*–d*), which
151 are consistent with the steps (b–d) from the *training phase*:

152 (b*) The observed values of the selected variable are obtained from the NYSM stations, resulting
153 in a 1D vector $X^{(s)}$ (Fig. 1(b*)).

154 (c*) The $X^{(s)}$ is then nearest-neighbor interpolated onto the full latitude–longitude grid, resulting
155 in a low-fidelity 2D gridded field $X^{(lf)}$ (Fig. 1(c)).

156 (d*) Now, this $X^{(lf)}$ is provided as input to the trained DL model from *training phase*, which
157 outputs a high-fidelity gridded field $X_{DL}^{(hf)}$.

158 The pipeline is $X^{(s)} \rightarrow X^{(lf)} \rightarrow X_{DL}^{(hf)}$. Since the DL model learned to interpolate the sparse field
159 (c) to the full spatially consistent field (d), we assume that the same interpolation applies when using
160 the real sparse observations as input (c* to d*). In other words, based on spatial correlations and
161 features learned from proxy data during training, the model can interpolate real sparse observations
162 to a full spatially consistent analysis field.

163 The key strength of the framework is its decoupled design, which reduces mismatches between
164 the temporal resolution of gridded mesoscale datasets and station observations. Each timestamp
165 is treated as a stand-alone sample, so the method does not depend on aligned time steps across
166 datasets. This flexibility allows gridded fields to be generated at the observation frequency (5
167 min for the NYSM) even when the training dataset has a coarser resolution (1 hr for the RTMA).
168 The method can be trained with alternative mesoscale products, such as HRRR (Dowell et al.
169 2022) or CONUS404 (Rasmussen et al. 2023), and can be adapted to other regions, observational
170 networks, and variables. Additionally, the NWP-based proxy dataset is only required during
171 training, whereas the inference phase relies solely on station observations. This design removes

172 dependence on operational NWP models, enabling gridded fields to be generated directly from
 173 sparse observations with minimal computational cost. As a result, the framework supports near
 174 real-time updates, a significant advantage over conventional DA–NWP systems.

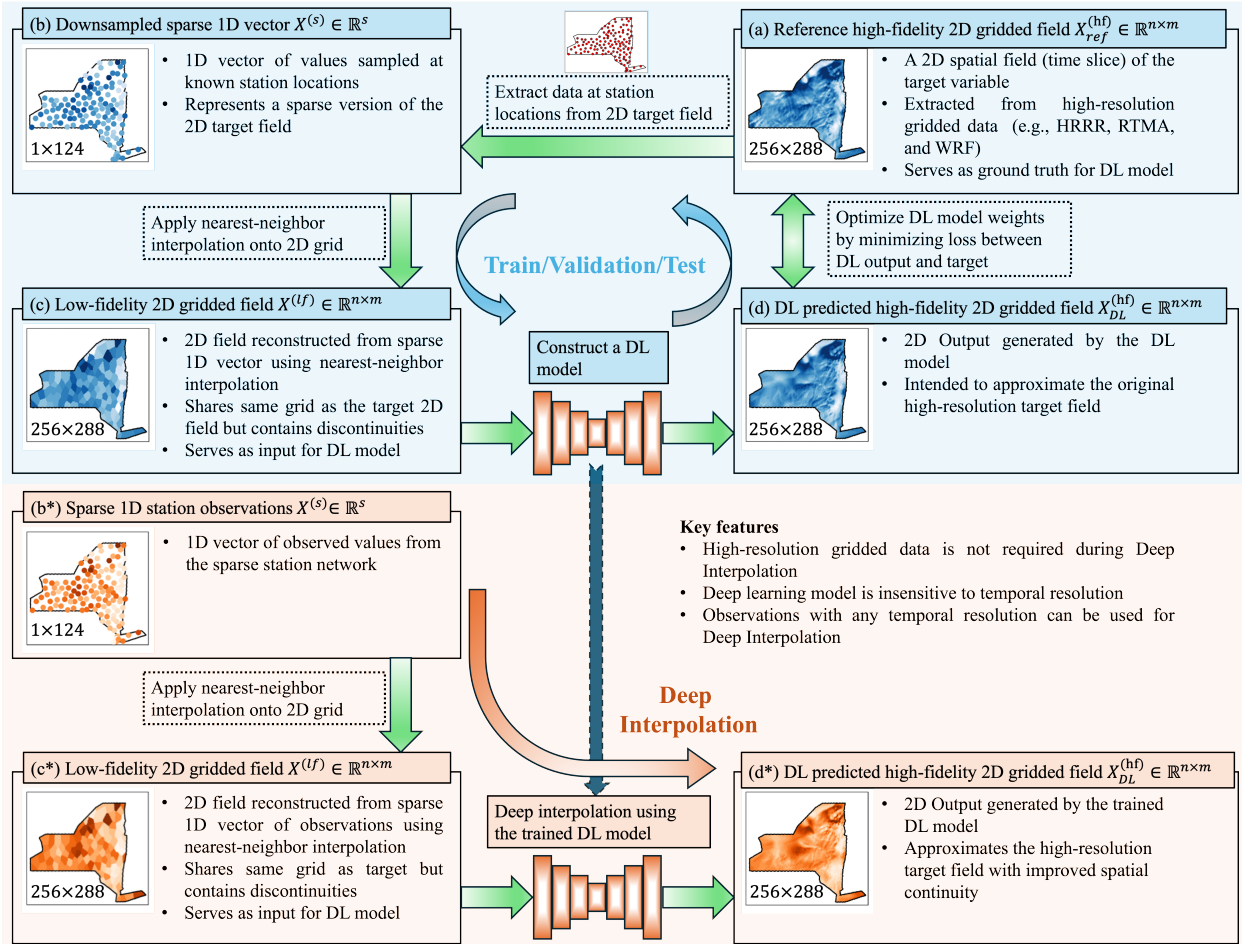


FIG. 1: Flowchart of the S2G-DI framework developed in this study.

175 3. Data

176 a. Study area and observation network

177 Our study focuses on the State of New York (NYS), located between 40.5° – 45° N latitude and
 178 72° – 80° W longitude (Fig. 2). NYS exhibits a wide range of topographical features, including
 179 elevated terrains such as the Adirondack Mountains in the northeast and the Catskill Mountains in
 180 the southeast-central region; major water bodies such as Lake Ontario (northwestern NY), Lake
 181 Erie (far western NY), the Hudson River, and the Atlantic Ocean; and lower terrain regions,

182 including the lake plains surrounding Lake Ontario and Lake Erie, the St. Lawrence River Valley
183 (north of the Adirondacks), the Hudson Valley (east of the Catskills), the Mohawk Valley (extending
184 west–east between the Adirondacks and Catskills), Long Island, and the coastal plain.

185 This topographic diversity results in substantial spatiotemporal variability in meteorological
186 conditions, contributing to a wide range of climatological phenomena across the state. For instance,
187 the lake plains frequently experience snowstorms driven by lake-effect snow processes (Niziol et al.
188 1995; Niziol 1987). Mountain-valley interactions influence the distribution and intensity of severe
189 weather events such as wind, hail, and tornadoes (Wasula et al. 2002). Additionally, Long Island is
190 subject to coastal meteorological effects from the Atlantic Ocean, including sea-breeze circulations
191 and low-level jet formations (Colle and Novak 2010).

192 Considering these diverse topographical features and associated climatology in the state of New
193 York, it is of paramount importance to have a denser observational network. For this purpose, the
194 NYSM has been in operation since 2018, with the main goal of monitoring high-impact weather
195 events across the state (Brotzge et al. 2020). It consists of 126 weather stations, out of which 124
196 are used in this study (Fig. 2) located across nearly all 62 counties, with additional stations in
197 larger or geographically diverse areas to capture local conditions more accurately. The network
198 has an average spacing of 30 km to observe mesoscale weather features. It plays a key role
199 in detecting and tracking events such as hurricanes, tornadoes, severe thunderstorms, flooding,
200 blizzards, snow squalls, ice storms, and extreme temperatures. Each station records a variety
201 of near-surface and subsurface variables, including air temperature (at two heights), dew point
202 temperature, relative humidity, wind speed and direction, pressure, precipitation, solar radiation,
203 snow depth, soil temperature and moisture (at three depths), and visible images.

204 *b. Data for Phase I: Training*

205 In this study, we used the National Centers for Environmental Prediction’s RTMA (De Pondeva
206 et al. 2011) as a gridded meteorological reference dataset for the deep learning training. The
207 RTMA is considered a proxy for the observations, since we expect this dataset to capture spatial
208 correlations reliably. The RTMA provides hourly analyses at a horizontal resolution of 2.5 km,
209 has been available since 2006, and is being continuously developed. It incorporates data from both
210 satellite and ground-based observations through data assimilation to attain near-real-time analysis.

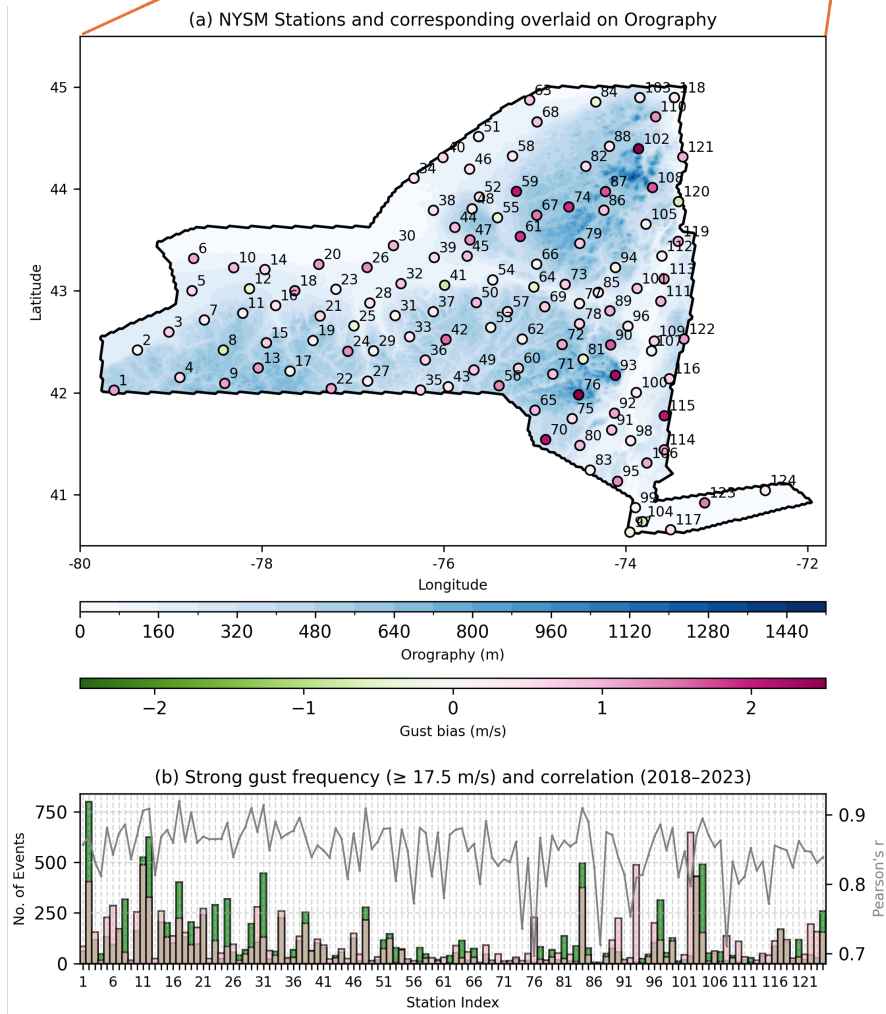
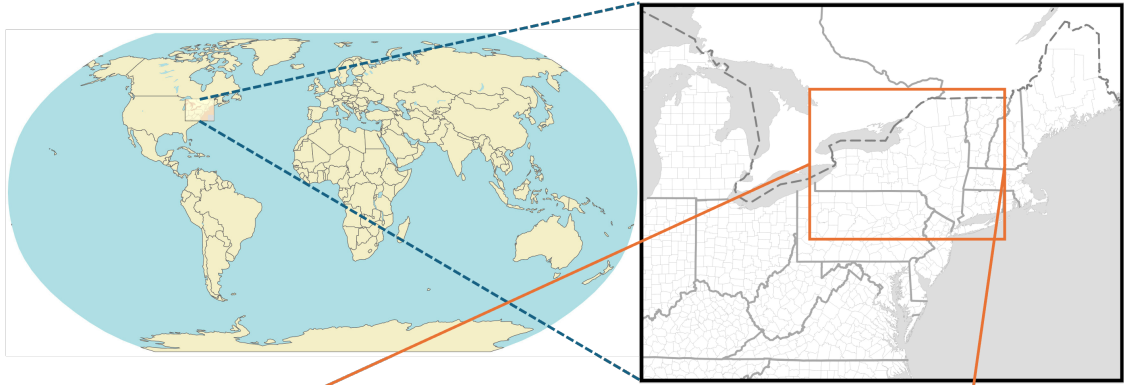


FIG. 2: An illustration of the study area and the NYSM observational network. (a) Spatial distribution of the NYSM stations, numbered in ascending order), overlaid on the orography. The station markers are colored according to the overall mean bias in 10 m wind gust between the RTMA analysis and NYSM observations, computed during 2018–2023 across individual stations. (b) The bar plot (left y-axis) shows the stationwise frequency of strong wind gust events (≥ 17.5 m s^{-1}) seen in NYSM and RTMA datasets; the line plot (right y-axis) indicates the stationwise Pearson’s correlation coefficient of wind gust between RTMA and NYSM datasets.

211 From the RTMA, we used 10 m wind speed (M_{10}), 10 m wind gust (G_{10}), 2 m air temperature (T_2),
212 and 2 m specific humidity (Q_2). These meteorological variables are time-varying. In contrast, the
213 geopotential height at the lowest model level is also used as static topographic information. Further
214 details of how the RTMA analysis is generated, the full list of variables, and their spatial coverage
215 can be found at De Pondeca et al. (2011). We used RTMA data for the period 2018 to 2023, at
216 hourly frequency. To generate spatially continuous fields within New York State, where the NYSM
217 stations are located, a rectangular domain of 256×288 was extracted from the RTMA dataset that
218 covers the entire state. This domain was then masked using the New York State boundary, keeping
219 only the data within the state and setting values outside the boundary to zero. This masking
220 approach was implemented to improve the accuracy of the deep learning models during training.

221 *c. Data for Phase II: Inference*

222 From the 124 NYSM stations, we obtained the same variables as the RTMA, which are M_{10} ,
223 G_{10} , T_2 , and Q_2 , at a frequency of 5 minutes. Since the NYSM records observations at a frequency
224 of 3 s, the wind gust is computed as the maximum within a 5-minute time window.

225 *d. Comparison of observed and reference data*

226 To evaluate the RTMA data with the observations, we selected wind gusts, and the NYSM wind
227 gusts were resampled by taking the maximum within a ± 10 -minute window centered on each
228 hour to match RTMA's hourly temporal resolution. The RTMA wind gusts exhibit a positive
229 mean bias relative to NYSM observations (Fig. 2(a)), indicating a general tendency to overestimate
230 gust magnitudes. Figure 2(b, left y-axis) presents the frequency of strong wind gust events
231 ($\geq 17.5 \text{ ms}^{-1}$), showing a mixed pattern: while several stations report more frequent events in
232 NYSM data, the majority display a higher frequency in RTMA. The Pearson correlation coefficients
233 between RTMA and NYSM gusts (Fig. 2(b), right y-axis) are above 0.85 at most stations, though
234 only a few exceed 0.90. Collectively, these results suggest that RTMA generally overestimates both
235 the magnitude and frequency of wind gusts compared to NYSM observations. Nonetheless, our
236 S2G-DI framework is not vulnerable to this bias, due to the decoupled design.

237 *e. Independent evaluation data*

238 Apart from observational datasets, we also utilize meteorological data from the High-Resolution
239 Rapid Refresh (HRRR) model (Dowell et al. 2022) as an independent validation source. HRRR
240 is a convection-allowing implementation of the Weather Research and Forecasting (WRF) model,
241 operating at a horizontal resolution of 3 km with an hourly data assimilation cycle over the CONUS.
242 From HRRR, we obtained the same variables as in RTMA, namely M_{10} , G_{10} , T_2 , and Q_2 , at hourly
243 frequency for the year 2023. HRRR data are used exclusively during independent evaluation
244 experiments and are not involved in any stage of model training.

245 **4. Deep Learning methodology**

246 *a. Model architectures and training strategy*

247 We chose three types of deep learning (DL) model architectures that vary in parameter count and
248 complexity (Appendices B(a–c)). First, we selected a simple deep convolutional neural network
249 (DCNN) architecture (Fig. B1) with minimal complexity and a low parameter count. Next,
250 we implemented a U-Net architecture (Fig. B2), which is widely used in image processing and
251 related tasks due to its encoder–decoder structure with skip connections. Finally, we increased
252 the model complexity by adopting a U-Net-style architecture with a Swin Transformer backbone
253 (SwinT2UNet; Fig. B3). All three architectures are well established and have been used in a variety
254 of related studies (Fukami et al. 2021; Sunderhaft et al. 2024; Wang et al. 2022). To maintain
255 focus and avoid cluttering the main text with standard model descriptions, we provide detailed
256 explanations of these architectures in Appendices B(a–c).

257 Following the framework in Fig. 1, the target tensor corresponds to the high-fidelity gridded field
258 of the variable of interest, taken from the reference RTMA dataset. For the 10 m wind gust as
259 the target variable ($X = G_{10}$), the reference field is denoted as $X_{ref}^{(hf)} = G_{10,ref}^{(hf)}$, and is represented
260 as a single-channel tensor of shape $\text{batch} (B) \times 1 \times 256 \times 288$. On the other hand, the DL models
261 take input tensors of shape $B \times C_{in} \times 256 \times 288$, where C_{in} denotes the number of input channels.
262 The first input channel corresponds to the low-fidelity (lf) interpolated field $G_{10}^{(lf)}$, as shown in
263 Fig. 1(c). The last input channel is a station mask, which has the same spatial dimensions as the
264 domain (256×288). This binary mask assigns a value of 1 to grid points corresponding to NYSM
265 station locations and 0 elsewhere, explicitly informing the model of the station locations within

266 the domain. The DL models' output is a single-channel tensor of the same shape as the target,
 267 representing the predicted high-fidelity (*lf*) gridded field $G_{10,DL}^{(hf)}$.

268 The RTMA data from 2018 to 2021, at an hourly frequency, are used for model training. The
 269 years 2022 and 2023 are used for validation and testing, respectively. While Fig. 1 illustrates the
 270 framework using a single time instance for clarity, in practice training is performed with mini-
 271 batches of samples. The models were trained using the Charbonnier loss (Charbonnier et al.
 272 2002), computed between the predicted and target fields:

$$\mathcal{L}(y, \hat{y}) = \sqrt{\|y - \hat{y}\|^2 + \epsilon^2}, \quad (1)$$

273 where $y = G_{10,ref}^{(hf)}$ denotes the reference field, $\hat{y} = G_{10,DL}^{(hf)}$ the model prediction, and $\epsilon = 0.001$.
 274 When $\epsilon = 0$, the loss reduces to the mean absolute error. Among conventional loss functions, mean
 275 squared error (MSE) tends to over-penalize large errors, while mean absolute error (MAE) is more
 276 robust but not differentiable at zero. The Charbonnier loss combines the advantages of both: it
 277 does not excessively penalize outliers like MSE, yet provides a smooth, continuously differentiable
 278 approximation to MAE. It has also been shown to achieve better performance than both MSE and
 279 MAE in super-resolution applications (Anagun et al. 2019; Cai et al. 2024).

280 We aim to predict the optimal continuous spatial field within the NYS. However, the selected
 281 domain covers the vast majority of land outside NYS, where the NYSM stations will have little
 282 observability. Thus, to improve model predictability and avoid nonphysical predictions, we only
 283 considered points lying inside the NY State boundary by masking outside points with zeros in the
 284 input and target tensors. In general, one would not expect to predict the target variable values at
 285 the station locations, since they are in fact already known by the DL models through the first input
 286 channel (interpolated field). Thus, the target variable values at the station locations are preserved
 287 by copying them directly from the first input tensor to the output tensor. To reflect this setup in the
 288 loss function, a masking strategy was employed: the loss was computed only over the grid points
 289 inside the New York State boundary, excluding those corresponding to NYSM station locations.
 290 This ensured that model updates were applied only within the masked boundary region, while
 291 preserving the NYSM station values unchanged in the output tensor.

292 The models are optimized using the Adam optimizer with the following hyperparameters: an
 293 initial learning rate of 0.003, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$, and no weight decay. Training

294 is carried out with an early stopping criterion: if the validation loss does not decrease for 20
 295 consecutive epochs, training is halted. The training continues for a maximum of 200 epochs if
 296 early stopping is not triggered. Additionally, an exponential learning rate decay with a decay factor
 297 of $\gamma = 0.9$ is applied during training. Each run is trained using two NVIDIA A100 GPUs, each
 298 with 80 GB of memory, with a batch size of 16. The models converge after approximately 60–70
 299 epochs, corresponding to a total training time of about 3 hours. Once trained, the inference time
 300 for generating a single analysis field is approximately 5 ms on a single A100 GPU, and less than
 301 400 ms when executed on 32 CPU threads on an AMD EPYC 7742 processor using PyTorch. This
 302 enables near real-time application even in CPU-only environments.

303 To assess how various design choices of our framework affect the quality of the S2G-DI analysis
 304 fields in comparison to Barnes baseline, we present several sensitivity studies in section 5, with
 305 RTMA acting as training data. Further, the validity of the framework is evaluated with true
 306 observations from NYSM (section 6) and with HRRR acting as an independent evaluation dataset
 307 (section 7).

308 *b. Quantitative evaluation metrics*

309 For qualitative evaluation, spatial maps of predicted and observed fields were visually compared.
 310 To quantitatively assess model accuracy, we adopted three widely used image-based metrics: Root
 311 Mean Square Error (RMSE) (Eq. 2), Peak Signal-to-Noise Ratio (PSNR) (Eq. 3), and Structural
 312 Similarity Index Measure (SSIM) (Eq. 4). These metrics treat both the predictions and targets
 313 as 2D spatial fields (images), allowing a pixel-wise comparison. RMSE measures the average
 314 magnitude of error, with larger errors penalized more heavily. PSNR quantifies the ratio between
 315 the maximum possible signal power and the power of the error, indicating reconstruction fidelity.
 316 SSIM evaluates the similarity in luminance, contrast, and structure between the prediction and
 317 target fields, and is particularly effective for assessing perceptual similarity.

$$\begin{aligned}
 \text{RMSE}(y, \hat{y}) &= \sqrt{\text{MSE}(y, \hat{y})} \\
 &= \sqrt{\frac{1}{|\Omega|} \sum_{(i,j,t) \in \Omega} (y_{i,j,t} - \hat{y}_{i,j,t})^2} \quad (2)
 \end{aligned}$$

318

$$\text{PSNR}(y, \hat{y}) = 10 \cdot \log_{10} \left(\frac{(y_{\max} - y_{\min})^2}{\text{MSE}(y, \hat{y})} \right) \quad (3)$$

319

$$\text{SSIM}(y, \hat{y}) = \frac{1}{|T|} \sum_{t \in T} \frac{(2\mu_{y_t} \mu_{\hat{y}_t} + C_1)(2\sigma_{y_t \hat{y}_t} + C_2)}{(\mu_{y_t}^2 + \mu_{\hat{y}_t}^2 + C_1)(\sigma_{y_t}^2 + \sigma_{\hat{y}_t}^2 + C_2)} \quad (4)$$

320 Here, $|\Omega|$ represents the total number of valid time instances and grid points lying inside the domain
 321 (excluding the station locations), $|T|$ is the total number of time instances (samples), t is individual
 322 time instance (sample), μ_y denotes pixel mean of the target tensor, $\mu_{\hat{y}}$ denotes mean of the predicted
 323 tensor, σ_y^2 denotes variance of the target tensor, $\sigma_{\hat{y}}^2$ denotes variance of the prediction tensor, $\sigma_{y\hat{y}}^2$
 324 is the covariance of target and predicted tensors, and $[C_1, C_2]$ are the two constants to stabilize the
 325 division with weak denominator.

326 The visual inspection or the metrics alone may fail to explain differences across spatial scales.
 327 Therefore, to provide a scale-aware quantitative assessment, we computed the power spectral
 328 density (PSD) of each prediction and target using the two-dimensional Fast Fourier Transform (2D
 329 FFT) (Skamarock 2004). The PSD characterizes how the variance of the spatial field is distributed
 330 across frequencies (scales). Since the resulting PSD is symmetric across quadrants, we extracted
 331 and analyzed only the diagonal of one quadrant to avoid redundancy.

332 *c. Cross-validation strategy*

333 We implemented a random-folding cross-validation (CV) strategy to evaluate model performance
 334 under a varying number of inference stations. This approach, illustrated in Fig. 3, simulates
 335 scenarios of random station dropout by selectively excluding groups of stations during inference.
 336 The 124 stations are divided into four spatially diverse groups, each with comparable intra-group
 337 distance characteristics. Two inference configurations are tested: one using three groups for
 338 inference and one for testing ($n_{\text{inference}} = 93$), and another using two groups for inference and two
 339 for testing ($n_{\text{inference}} = 62$). In the 93 inference stations setup, four folds (F1–F4) are generated by
 340 rotating which group is held out as unseen, while the remaining three groups are used for inference.
 341 The union of the held-out groups across all four folds ensures that every station is exactly unseen
 342 once, providing a single complete set of unseen stations. Similarly, in the 62 inference station
 343 configuration, six folds (F1–F6) are created by considering all possible combinations of two
 344 groups used for inference and the remaining two as unseen. This results in three disjoint sets of

345 entirely unseen stations, which can be considered as three ensemble members of the full set of
 346 unseen stations.

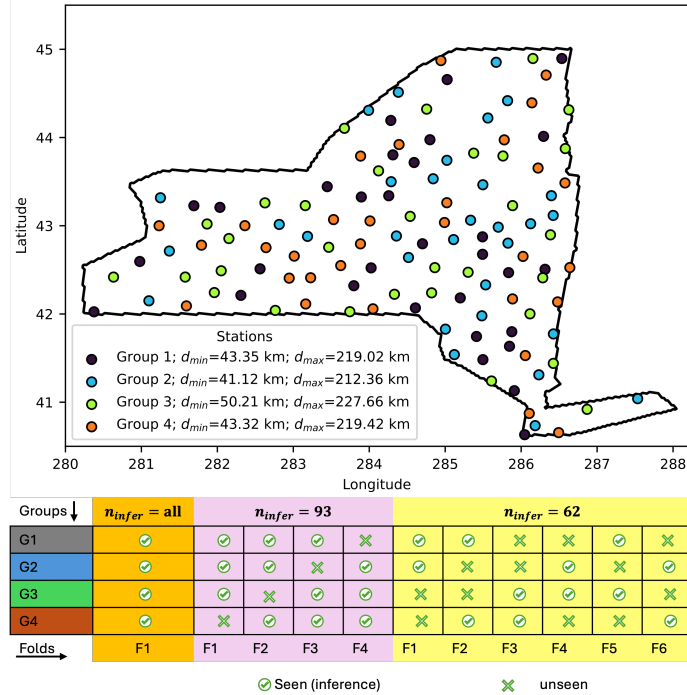


FIG. 3: Station grouping and fold-wise inference strategy for model evaluation. The map (top) shows the spatial distribution of 124 stations across New York, divided into four groups (G1–G4) with the same number of stations (31) and similar intra-group distance statistics. Each group maintains comparable minimum and maximum pairwise distances. The matrix (bottom) illustrates the folding strategies used during inference, with green checkmarks (✓) indicating groups used during inference (seen), and green crosses (✗) indicating groups excluded during inference (unseen).

346

347 5. Results: Sensitivity analysis with RTMA data

348 We conduct a sensitivity study to illustrate the influence of various design choices on the deep
 349 interpolation performance of our approach, i.e., to generate high-resolution gridded analysis fields
 350 from sparse inputs (Fig. 1(a)-(d)). The sensitivity study is conducted using gridded RTMA data only
 351 to determine a model configuration that achieves the best possible sparse-to-gridded reconstruction.
 352 We employ wind gust as an illustrative target variable for these sensitivity experiments because it
 353 is governed by complex processes and, thus, is challenging to model (Kwon and Kareem 2009).

354 *a. Sensitivity to model architecture and input variables*

355 We first aim to optimize the DL model architecture and the subset of available variables to
356 be used as inputs. More specifically, we compare a low-complexity deep convolutional neural
357 network (DCNN, ca. 0.7M parameters, cf. Appendix B(a)), a standard medium-complexity UNet
358 (ca. 10M parameters, cf. Appendix B(b)), and a high-complexity UNet with SWin transformer
359 (ca. 15M parameters, cf. Appendix B(c)). All three architectures are trained on different subsets of
360 input features, with the high-resolution wind gust G_{10} from RTMA as the target in all cases. The
361 different configurations are listed in Table 1. The simplest experiment takes only nearest-neighbor-
362 interpolated G_{10}^{lf} and the station masks as input. More input features are added successively,
363 beginning with the terrain height H_0 from the RTMA dataset, which is provided as a continuous
364 spatial field to inform the DL models about the underlying topography.

365 Other meteorological variables are the 10 m wind speed (M_{10}), 2 m temperature (T_2), and 2 m
366 specific humidity (Q_2), which are all provided as nearest-neighbor-interpolated fields as explained
367 in Fig. 1(c). Note that the same variables should be available from both RTMA and NYSM
368 observations, as generating analysis fields from only observations during inference is otherwise
369 not possible (cf. Fig. 1(b*-d*)).

370 The three different architectures and five different subsets of input variables yield 15 distinct
371 model configurations. The model weights of each combination are randomly initialized using
372 Xavier initialization (Glorot and Bengio 2010). To obtain an ensemble and examine the robustness
373 of the performance to different architectures and variables, we repeat each of the 15 experiments
374 five times with different random initializations, resulting in a total of 75 runs. Additionally, Barnes
375 interpolation is employed to generate baseline wind gust analysis fields from the same sparse
376 RTMA wind gust inputs used by the DL models.

377 The performance of the different configurations with respect to RMSE, PSNR, and SSMI is
378 presented in Fig. 4. All scores are normalized by the scores achieved by Barnes interpolation to
379 illustrate the benefit of our DL approach. To get a qualitative understanding of the performance,
380 examples of analysis fields generated by different models are presented in Fig. 5 as spatial maps. The
381 traditional Barnes interpolation fields are also displayed together with bias maps of the estimated
382 fields compared to the RTMA ground truth. Finally, the average power spectral density (PSD)

TABLE 1: List of experiments conducted to assess the sensitivity of deep interpolation performance to different input variable configurations. Here, lf indicates low-fidelity and hf indicates high-fidelity.

ID	Inputs						Target
	$G_{10}^{(lf)}$	$M_{10}^{(lf)}$	$T_2^{(lf)}$	$Q_2^{(lf)}$	Terrain height $H_0^{(hf)}$	Station mask	$G_{10}^{(hf)}$
G ₁₀	✓					✓	✓
G ₁₀ -H ₀	✓				✓	✓	✓
G ₁₀ -H ₀ -M ₁₀	✓	✓			✓	✓	✓
G ₁₀ -H ₀ -M ₁₀ -T ₂	✓	✓	✓		✓	✓	✓
G ₁₀ -H ₀ -M ₁₀ -T ₂ -Q ₂	✓	✓	✓	✓	✓	✓	✓

383 displayed in Fig. 6 presents how well the different model configurations can capture patterns of
 384 various sizes compared to the RTMA training data.

385 From Fig. 4 we see that the UNet performs overall best in all three metrics and all variable
 386 configurations, followed by SwinT2UNet and DCNN. Error bars indicate that UNet and DCNN
 387 are insensitive to random weight initialization, as their performance varies only slightly. This is
 388 not the case for SwinT2UNet, where different initializations results in different scores, as indicated
 389 by larger error bars. As the SWinT2UNet has 50% more parameters to tune than the UNet, the
 390 increased variance of the transformer UNet scores may also indicate too large a model capacity with
 391 respect to our training data. The spatial maps (first row) reveal significant qualitative differences
 392 between Barnes interpolation (b*) and the three DL methods (c* – e*). All three DL methods
 393 yield significantly better performance than Barnes interpolation, as is evident qualitatively from the
 394 maps and quantitatively from the statistics discussed before. The DL methods recover more spatial
 395 details related to topography compared with Barnes interpolation, as shown in the bias maps in the
 396 second row. That is also visible from the PSD in Fig. 6 where the curves of the three deep learning
 397 models overlap. All three DL models are close to the RTMA ground truth (black), indicating that
 398 spatial patterns of all scales are matched well. The PSD of Barnes interpolation (blue), on the
 399 other hand, deviates significantly from the ground truth and the DL models at wavelengths smaller
 400 than 40 km, indicating that finer structures, such as terrain influences, are missed. This difference
 401 is expected, as Barnes interpolation does not consider orography, whereas our framework can take
 402 it as input.

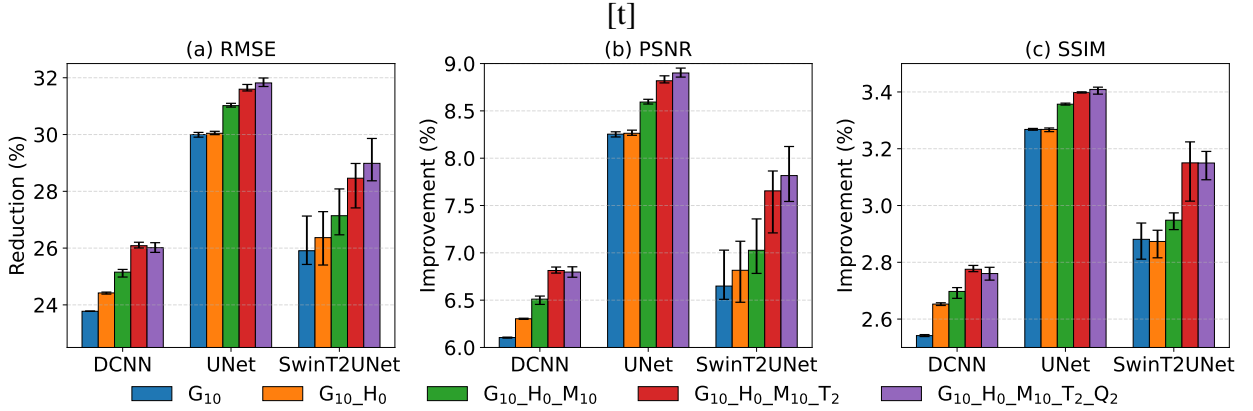


FIG. 4: Change of relative performance of our deep interpolation framework depending on the complexity of the deep learning model ($DCNN < UNet < SwinT2UNet$) and the number of input variables used (cf. Tab. 1). Statistics are given relative to traditional Barnes interpolation. Error bars indicate sensitivity to random weight initialization.

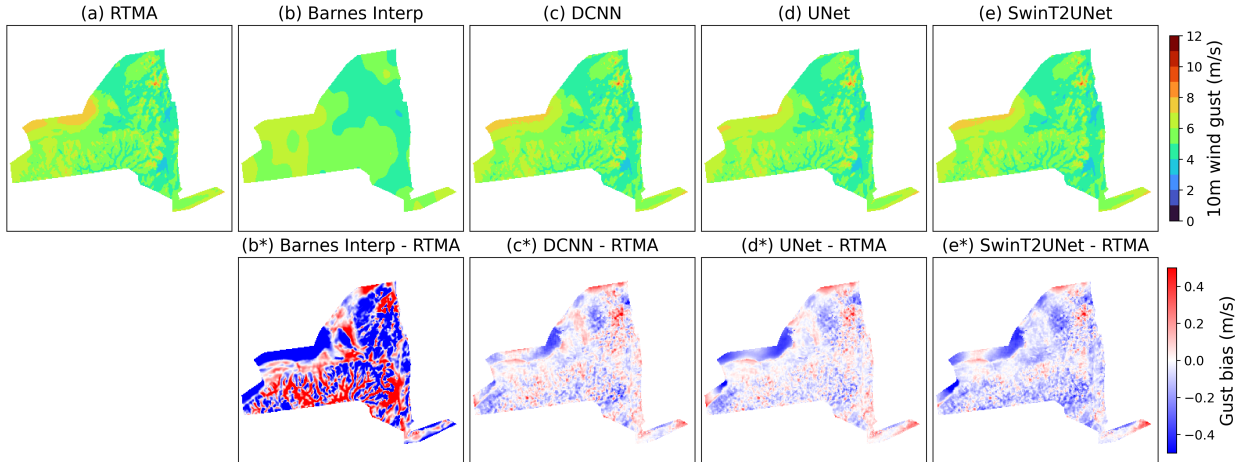


FIG. 5: (First row) Comparison of the annual spatial mean 10 m wind gust from (a) RTMA, (b) Barnes interpolation, (c) DCNN, (d) UNet, and (e) SwinT2UNet. The deep learning models are configured with the best input configuration ($G_{10_H_0_M_{10}_T_2_Q_2}$). (Second row) Spatial bias of the Barnes interpolation (b^*) and deep learning models ($c^* - e^*$) relative to the RTMA annual mean 10 m wind gust.

403 Considering also the different input configurations presented in Fig 4, it is clear that adding
 404 more variables improves the performance for all DL architectures. Adding orography does not
 405 seem to improve the performance of the two UNets much compared to DCNN, but providing M_{10}
 406 and T_2 yields better scores across all metrics and models. Including also the specific humidity
 407 Q_2 , only yields a minor performance improvement, suggesting that it does not contain much extra
 408 information for gust prediction compared to the previously added variables. Nevertheless, we keep
 409 Q_2 to achieve maximum performance. The overall performance increase due to extra variables is

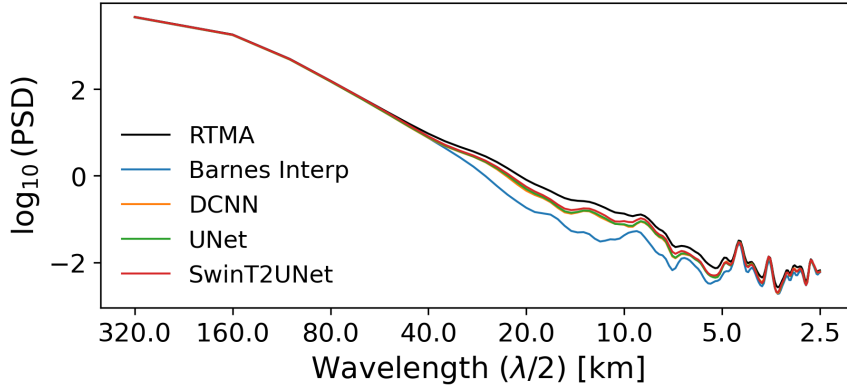


FIG. 6: Annual mean power spectral density (PSD) of RTMA (black) compared with interpolated fields using Barnes’s method (blue) and different deep learning architectures provided the best input configuration ($G_{10_H0_M10_T2_Q2}$).

410 attributed to physical dependencies between variables, which improve the predictability of wind
 411 gusts.

412 **Intermediate result:** While all DL architectures are close to each other in performance, we
 413 select the UNet because of its overall high skill, its insensitivity to random initialization, and its
 414 balanced complexity. As adding more variables improves performance, we continue the study with
 415 the UNet trained on all inputs.

416 *b. Sensitivity of Barnes interpolation: bias correction*

417 To further investigate the origin of the performance gains observed for the deep interpolation
 418 framework, we design a modified Barnes configuration that explicitly tests whether the improve-
 419 ments of the U-Net can be attributed to distributional bias correction alone.

420 We apply a grid-point quantile mapping (QM) correction (Themeßl et al. 2011) to the Barnes-
 421 interpolated fields (hereafter Barnes_QM), using RTMA as the reference dataset. The QM correc-
 422 tion is constructed during the training period (2018–2021) by first generating Barnes-interpolated
 423 fields from RTMA data using the same station sampling strategy as in the main experiments.
 424 For each grid point, empirical cumulative distribution functions are estimated for both the Barnes-
 425 interpolated values and the corresponding RTMA values using 1000 quantile intervals to adequately
 426 resolve the full distribution, including its upper tail.

427 The quantile mapping function is defined at each grid point by pairing corresponding quantiles
 428 of the Barnes and RTMA distributions, thereby establishing a transfer function that maps Barnes

429 values to RTMA-consistent values. During the evaluation period (2023), this precomputed mapping
 430 is applied to Barnes-interpolated fields generated from RTMA-based station inputs, without further
 431 adjustment. In this way, Barnes_QM preserves the spatial structure imposed by the interpolation
 432 while enforcing the marginal distribution of the reference field at each grid point.

433 This experiment isolates the role of distributional alignment in the observed performance gains
 434 and assesses whether correcting marginal distributions alone is sufficient to reproduce the spatial
 435 characteristics captured by the deep learning approach. As shown in Fig. 7, the annual mean gust
 436 field from Barnes_QM is visually closer to RTMA in terms of small-scale texture compared to
 437 the Barnes baseline, which is expected given the imposed RTMA-consistent distributions. The
 438 correction drastically reduced the overall bias compared to Barnes baseline, but not as better
 439 compared to the UNet. Further, the correction only aligns the mean characteristics, but fails at
 440 temporal consistency.

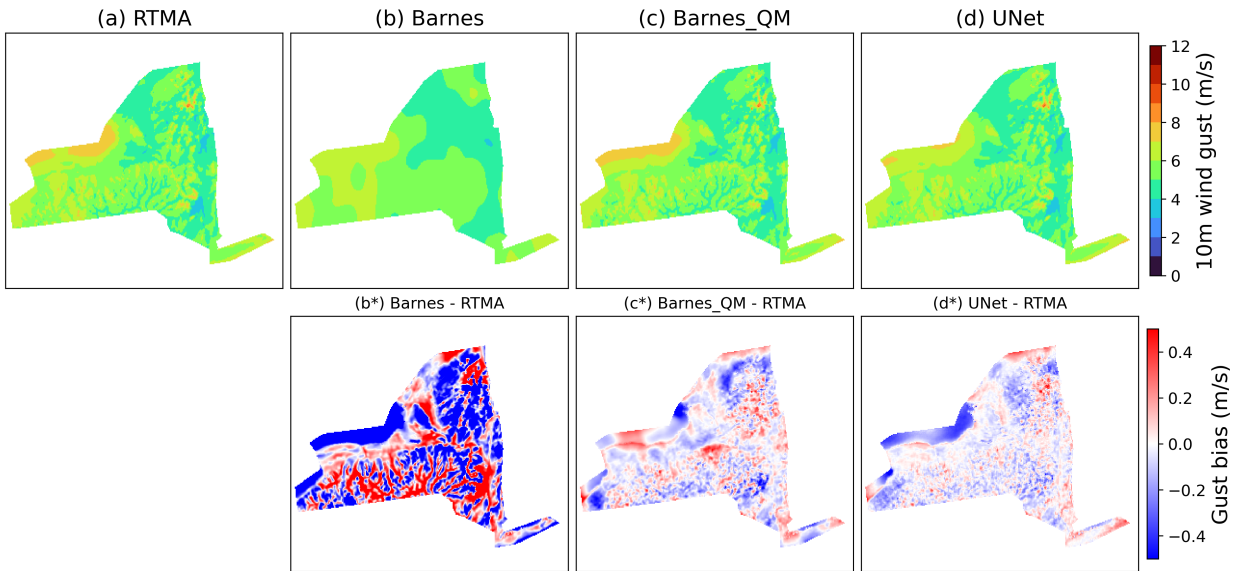


FIG. 7: Comparison of annual mean 10 m wind gust fields over New York State using (a) RTMA, (b) Barnes interpolation, (c) grid-point quantile-mapped Barnes, and (d) UNet. Panels (b*)–(d*) show the corresponding bias fields relative to RTMA. All fields are aggregated over the evaluation period 2023. Barnes and UNet are included for reference.

441 Table 2 supports these findings quantitatively. QM-based correction yields only modest improve-
 442 ments over Barnes (7.5% RMSE reduction, 1.1% PSNR gain, and 1.1% SSIM gain), indicating that
 443 aligning the marginal distribution alone does not substantially improve the reconstructed fields.

444 This reflects the locality of the QM correction, which does not enforce temporal frame-by-frame
445 matching.

446 In contrast, the U-Net yields the strongest gains across all metrics, including a 31.7% reduc-
447 tion in RMSE and a 3.4% increase in SSIM. These improvements demonstrate that the superior
448 performance of the deep interpolation framework cannot be explained by bias correction alone,
449 but instead arises from its ability to learn spatial and temporal coherent relationships between
450 meteorological variables and underlying characteristics.

TABLE 2: Relative performance of different methods compared to baseline Barnes interpolation using RTMA as reference.

Method	RMSE reduction (%)	PSNR improvement (%)	SSIM improvement (%)
Barnes.QM	7.5	1.1	1.1
UNet	31.7	8.9	3.4

451 *c. Sensitivity to amount of training data*

452 The next step of the sensitivity study is to successively reduce the training data, one year at
453 a time. Specifically, the models are trained on data from 2018–2021 (4 years, same as before),
454 2018–2020 (3 years), 2018–2019 (2 years), and 2018 (1 year). These experiments aim to assess
455 the influence of training data size on model performance. In all cases, 2022 and 2023 are kept for
456 validation and testing, respectively, and the optimal UNet architecture with all inputs is used. The
457 motivation behind this step is to assess the applicability of our framework for regions where no
458 extensive proxy dataset like RTMA is available. Generating a similar high-resolution dataset from
459 scratch using numerical weather models is expensive, so running as few simulations as possible
460 for training data generation is desirable.

461 Tab. 3 presents the improvement of RMSE, PSNR, and SSIM over Barnes interpolation for four
462 models trained on different ranges of training data. As expected, performance decreases across
463 all metrics when fewer training data are used. We attribute this to the variability of weather
464 across different years, so datasets of only a few years are expected to contain less diverse weather,
465 resulting in less generalized models. Even then, the UNet trained on only one year (2018) achieves
466 a significantly improved statistics (reduced RSME, and increased PSNR and SSIM) compared

TABLE 3: Change of relative performance of our deep interpolation technique compared to Barnes interpolation when one, two, three, or four years of gridded RTMA data are used for training. The marked (*) configuration is used in other parts of the sensitivity study.

Train years	RMSE reduction (%)	PSNR improvement (%)	SSIM improvement (%)
2018-2021 ^(*)	31.7	8.9	3.4
2018-2020	30.5	8.4	3.2
2018-2019	29.2	7.9	3.1
2018	26.4	6.9	2.8

467 to the traditional Barnes method. This result is very promising for applications in regions with
 468 fewer available proxy data. Comparing the change in scores from one year of training data to four
 469 years also indicates a diminishing return on using more and more data. For RMSE, for example,
 470 training on two years compared to one year yields an improvement of 2.8 percent points, whereas
 471 going from three to four years only improves the scores by 1.2 percent points. Similar patterns
 472 are observed in the other metrics. We expect this trend to continue with performance saturating at
 473 some point.

474 **Intermediate result:** In conclusion, we demonstrate that a single year of training data yields
 475 a model that outperforms the traditional Barnes interpolation, and that adding more years of data
 476 has a positive but gradually diminishing effect. For the rest of the study, we continue with all four
 477 training years and the all-variable UNet to achieve the best possible performance.

478 *d. Random dropout of training stations*

479 The final step of the sensitivity analysis assesses the effect of data missing from some stations
480 during inference, i.e., when the analysis field is generated from incomplete sparse inputs. In
481 an extensive observational network, such as NYSM with 100+ stations, it is common that data is
482 missing at a few locations for certain periods due to sensor failures, maintenance, or communication
483 problems. In general, our framework is robust against that, as the nearest-neighbor interpolation
484 (cf. Fig 1b* and c*) can always generate a 2D field, which can be refined by the DL model.
485 However, having fewer sparse inputs results in different polygon patterns of the interpolated fields.
486 In the following, we aim to assess the impact of these different patterns on the final output and
487 stress that dropout during training is, in fact, crucial to make the model robust against different
488 polygon patterns, and thus, missing data during deep interpolation.

489 To assess the impact of missing stations on model performance, we drop stations during training
490 and during inference. As training happens on synthetic gridded RTMA data, missing values do
491 not typically appear at this stage. The motivation for still dropping stations is to simulate different
492 network configurations during training and to assess the robustness of the training. We begin by
493 utilizing all 124 NYSM stations during training and successively reduce the number of stations to
494 $n = 100$, $n = 75$, and $n = 50$. More specifically, during each epoch of the training, we randomly
495 select only n stations of each sample and interpolate them using nearest neighbors, resulting
496 in different polygon patterns for each sample and epoch. The resulting four models are called
497 $DL_{n_{all}}$, $DL_{n_{100}}$, $DL_{n_{75}}$, and $DL_{n_{50}}$ in the following. Note that these training stations are virtual
498 stations, where RTMA data are extracted at the NYSM locations and the DL models are trained
499 to recover the full field (training phase, cf. Fig 1). Next, missing data during inference, i.e., deep
500 interpolation, is simulated following the cross-validation strategy explained in Section 4(c). The
501 strategy provides dataset versions with $m = 124$ stations, $m = 93$ stations, or only $m = 62$ stations for
502 inference. Evaluating all four trained models with all three station subsets yields a total of twelve
503 experimental combinations, which are discussed below. The CV strategy enables us to evaluate
504 the performance of all experiments on the full 124 stations, as we always make predictions for
505 the left-out folds, making results comparable between experiments. Barnes interpolation is again
506 employed as the baseline. As performance is naturally expected to drop when fewer observations are
507 available, we compare DL-interpolated fields and Barnes-interpolated fields in the same situations.

508 For example, when only $m = 93$ stations are selected for inference for the DL models, the same
 509 stations are used for Barnes interpolation. That way, we can compare the advantage of our models
 510 over Barnes interpolation fairly. For an absolute picture, we also compare all our missing data
 511 experiments against the best-performing Barnes-interpolated fields, utilizing all 124 stations, to
 512 demonstrate the overall advantage of our framework.

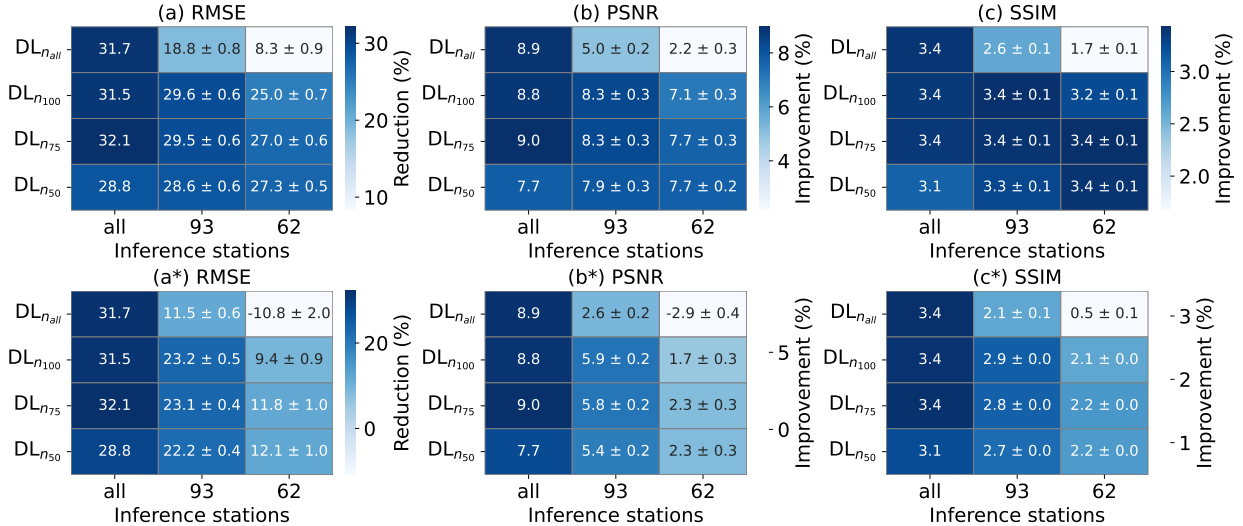


FIG. 8: Performance of deep learning models relative to Barnes interpolation across different numbers of training stations (rows) and inference stations (columns), simulating missing data in the observational network. Colors show improvement in (a) RMSE, (b) PSNR, and (c) SSIM metrics. Scores in the first row are normalized column-wise, i.e., with scores achieved by Barnes interpolation with the same number of missing stations. Second row scores (a* – c*) are all relative to the best-performing Barnes interpolated fields obtained when using observations from all stations ($m = 124$). The rows of each matrix present how the performance changes when one specific trained model (e.g., $DL_{n_{all}}$ trained on all stations) is presented with different fractions of missing data during inference. Similarly, the columns illustrate how randomly dropping different ratios of stations already during training affects the performance when a fixed number of stations (e.g., 93) is available during inference.

513 The results are presented in Fig. 8. Qualitatively, the patterns of changing performance in Fig. 8
 514 are the same across all panels. Unsurprisingly, the best extrapolation performance is achieved
 515 when all stations are available during inference (“all” columns). However, except for the $DL_{n_{all}}$,
 516 the relative performance (first row) barely changes when 93 or 62 stations are used during inference.
 517 These station numbers correspond to 25% and 50% of missing data, which are significant amounts.
 518 Even in absolute terms (second row), our framework retains a significant advantage over the all-
 519 station Barnes interpolation results in these situations, as deep interpolation with 50% of missing
 520 values still outperforms all-station Barnes interpolation by more than 10%.

521 Only the $DL_{n_{all}}$ experiments show a significant drop in performance for an increasing number of
522 missing stations. This behavior is attributed to a lack of generalization of the trained model. If all
523 station locations are used during training, i.e., no randomization happens, the model always sees
524 the same polygon pattern, only with different cell values. When fewer stations are available during
525 inference, this pattern changes, which the model cannot handle well. In contrast, the models trained
526 on random station subsets exhibited diverse polygon patterns during training and are consequently
527 robust against changes in the number of stations during inference, emphasizing that training station
528 randomization is crucial.

529 **Intermediate result:** Our framework maintains a consistent advantage over the Barnes interpo-
530 lation, even when data is missing, *if* station locations are already randomized during training to
531 make the model more robust. Based on the metrics presented in Fig. 8, the $n = 75$ model is selected
532 for the rest of this work and is treated as the final S2G-DI configuration. The results also suggest
533 that our framework can be used to fill gaps in observations across the network.

534 6. Results: Evaluation based on observational data (NYSM)

535 The S2G-DI model optimized in the sensitivity study is used here to assess the quality of analysis
536 fields generated from real observations of the New York State Mesonet (NYSM). In this section,
537 NYSM observations are provided as inputs at stage b* (Fig. 1) to generate the analysis fields. So
538 far, all evaluations of the deep learning framework have been conducted using synthetic RTMA
539 proxy data, which also served as input for analysis generation. By replacing these proxy inputs
540 with real observations, this evaluation tests whether the model produces realistic analysis fields
541 when applied to data that differ from the training distribution.

542 a. Performance in extreme wind gust conditions

543 Three example cases of extreme wind gust ($G_{10} > 17.5$ m/s) are selected from our held-out
544 testing year 2023 to provide a comprehensive assessment across different types of extreme wind
545 phenomena, from synoptic-scale frontal systems to mesoscale severe weather outbreaks. Focusing
546 on extreme gusts allows us to evaluate our approach under challenging conditions. The events are
547 listed in Tab. 4 and visualized in Fig. 9 as time series at two stations affected by the respective
548 events. The curves presented for our approach (blue) are always taken from the out-of-fold

549 predictions following our CV strategy (cf. Section c). In other words, a subset of stations *not*
550 *containing* the presented stations is used for the deep interpolation to obtain a gridded analysis
551 field. The presented time series are extracted from these analysis fields at the locations of the
552 respective stations (blue lines) and compared against the unseen true observations (black lines).
553 The Barnes interpolation curves (red) are obtained following the same strategy, and RTMA data
554 (cyan diamonds) are displayed for reference. The uncertainty intervals around the DL curves and
555 Barnes curves stem from the ensemble prediction generated during cross-validation.

TABLE 4: Extreme wind gust events ($G_{10} > 17.5 \text{ m s}^{-1}$) selected to assess the extreme prediction capability of the DL models.

#	Period	Synoptic conditions	Peak Wind Gust	Affected Regions
1	Feb 3–4, 2023	Powerful arctic cold front, leading to strong gusts and gradient winds.	20–29 m s^{-1}	North-West and South-East of NYS
2	Mar 25–26, 2023	Strong low-pressure system, producing widespread gusty winds before and after frontal passage (CNY Central 2023a,b).	20–27 m s^{-1}	Central NYS
3	Apr 1–2, 2023	Rare early-spring tornado outbreak, leading to severe thunderstorms, damaging wind gusts, and multiple confirmed tornadoes (National Weather Service Philadelphia/Mt. Holly 2023; New York Post 2023a,b).	25–30 m s^{-1}	Central and South-Eastern NYS

556 On first sight, our DL approach matches the observed gusts well in magnitude and timing,
557 although very strong gust magnitudes are often underestimated. That is visible, for example, in
558 panel (a) between 9:00 and 13:00 or panel (c) around 15:00 during the first day. It seems that
559 our approach particularly struggles with very dynamic gusts with short peak gusts in a short time
560 interval, while less dynamic, more sustained strong gust events, as depicted in panels (b) and
561 (d), are captured better. Barnes interpolation, on the other hand, generally captures trends but
562 exhibits more underestimation and less dynamic behaviour. In some cases, such as from 15:00 to
563 18:00 in panel (f), the Barnes interpolation completely misses the extreme gust. This performance
564 advantage over Barnes interpolation is attributed to the auxiliary variables, such as M_{10} , T_2 , and
565 the terrain information, which our approach can utilize. However, extreme gust peaks are very
566 rare events with less than 0.5% of yearly samples satisfying $G_{10} > 17.5 \text{ m/s}$. That makes the
567 extreme gust regression a very imbalanced regression problem, which is inherently challenging
568 to handle (He and Garcia 2009). Machine learning models typically favor the core of the data

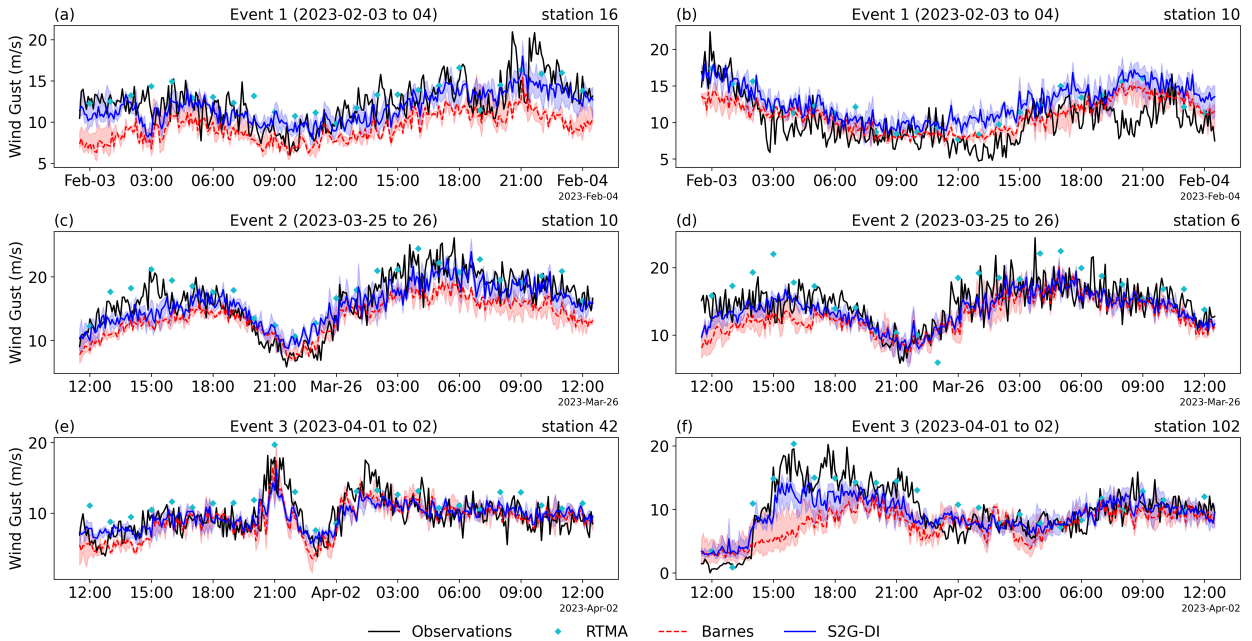


FIG. 9: Time series comparison of 10 m wind gusts during three extreme events (rows), with two representative (randomly selected) out-of-sample stations shown per event (columns). The curves for Barnes interpolation (red) and our S2G-DI (blue) represent ensemble medians, with shaded areas indicating the minima and maxima of the ensemble. RTMA (cyan diamonds) and observations (black) are shown for reference.

569 distribution, not its tails, as tails could be outliers, and because missing a few extreme values
 570 does not result in a significant penalty during optimization. While we recognize that further
 571 improvements in representing extreme gust magnitudes are needed, the current performance of the
 572 deep interpolation framework is encouraging, as it reliably captures the temporal evolution of gust
 573 events at unseen locations.

574 After discussing the interpolation performance at individual locations, Fig. 10 presents the
 575 spatial extents of the three events as captured by (a) our deep interpolation and by (b) traditional
 576 Barnes interpolation. Comparing deep interpolation and Barnes interpolation reveals that the
 577 Barnes interpolated fields suffer much more from missing stations (going from columns (i) to (iii))
 578 than the S2G-DI approach. Particularly for events two and three, it is evident that the estimated
 579 area affected by wind gust (orange and red) increases drastically for Barnes interpolation when
 580 compared to deep interpolation. Our approach, however, retains a remarkably consistent spatial
 581 pattern compared to the full network (a, i), even when observations are available at only half of the
 582 stations (a, iii). We view this spatial consistency and the similarity of the deep interpolated analysis

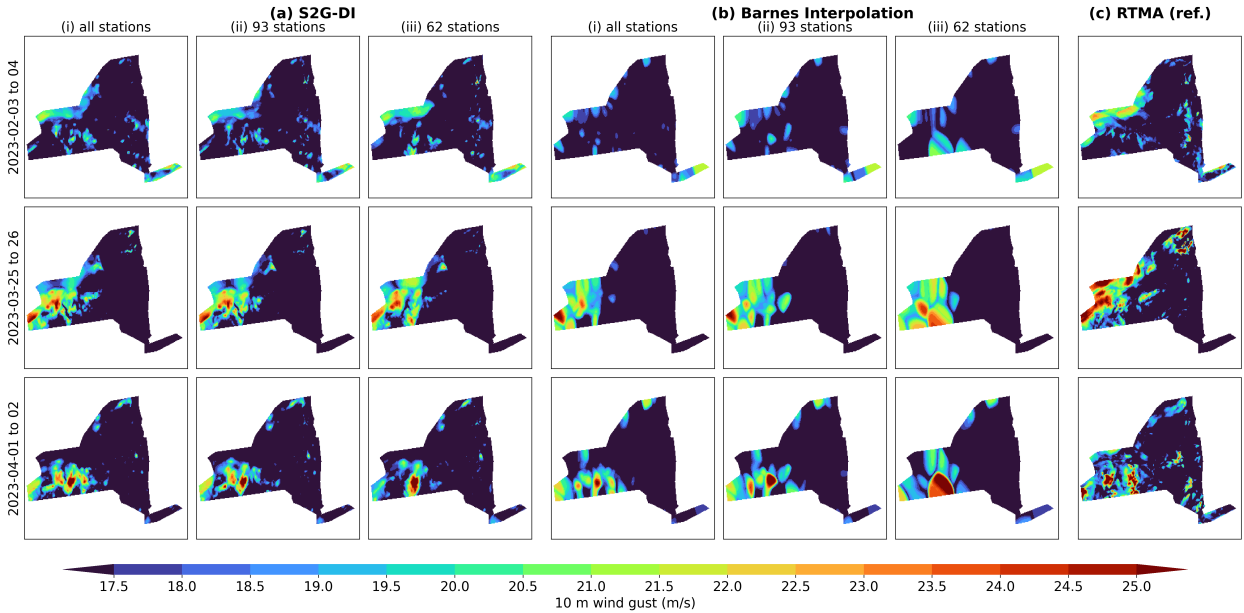


FIG. 10: Comparison of 10 m maximum wind gust analysis fields for three extreme events (rows) using (a) our deep interpolation framework and (b) traditional Barnes interpolation. The three columns per interpolation method (i – iii) present the change of the spatial pattern if observations are only available at a subset of stations for interpolation. In the absence of an absolute 2D ground truth, (c) RTMA is given for reference. The maps show the maximum wind gust per grid point for the entire event, thereby illustrating the aggregated geographical extent affected by the event. The colorbar is set to only show areas where gust exceeds the extreme threshold of 17.5 m/s.

583 fields to RTMA as a strong indication of the correctness and robustness of our method. The spatial
 584 findings are consistent with the underestimations discussed for the time series (cf. Fig. 9) and the
 585 missing-station experiments presented in the sensitivity study.

586 The spatial behavior of the analysis fields is explored in more detail in Fig. 11 using PSDs. The
 587 displayed curves are the time averages of the PSDs obtained for each snapshot of the interpolated
 588 analysis field from deep interpolation and Barnes interpolation, respectively. As such, the PSDs
 589 reflect the average detail with which an event is represented by either method. Again, RTMA
 590 is shown in black for reference. As observed during the sensitivity study, the observation-based
 591 analysis fields produced by Barnes interpolation also lose details for wavelengths of 40 km and
 592 smaller. The S2G-DI fields, on the other hand, retain more details, corroborating the advantage
 593 and value of our method.

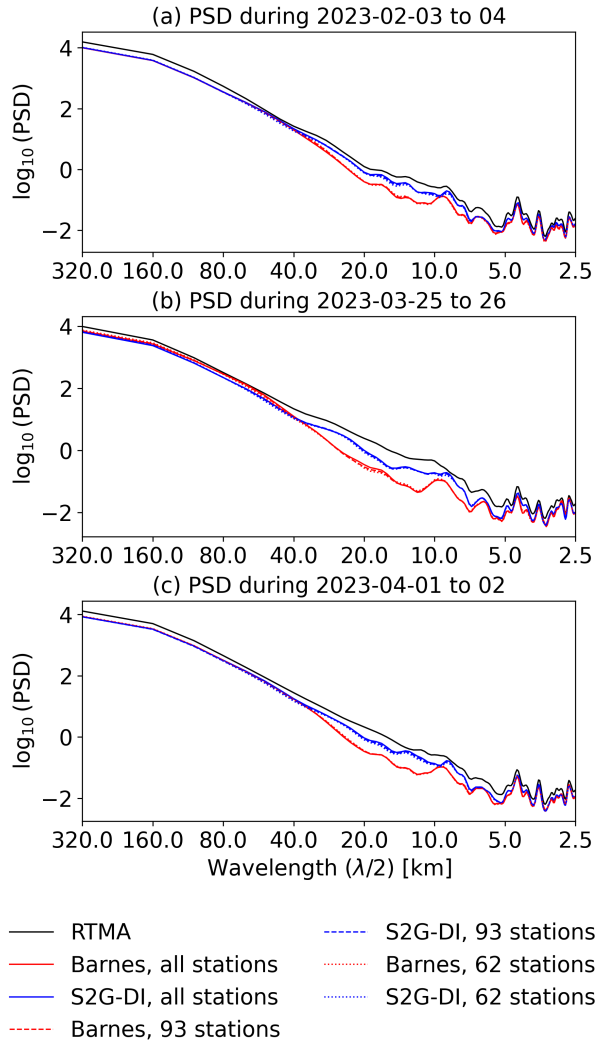


FIG. 11: Comparison of average power spectral density (PSD) for deep interpolated wind gust fields (blue) and fields from Barnes interpolation (red). Solid, dashed, and dotted lines indicate generation of analysis fields using all stations, 93 stations, and 62 stations, respectively. RTMA data is shown for reference (black), but cannot be viewed as absolute ground truth.

594 *b. General performance across all wind gust conditions*

595 While the selected extreme events provide detailed insight into the temporal and spatial behavior
 596 of the interpolation methods, it is important to assess whether these findings generalize across the
 597 full evaluation period. To this end, we aggregate the results over the entire year 2023 and evaluate
 598 the performance at held-out stations using complementary statistical metrics (Fig. 12).

599 The evaluation is performed at two temporal resolutions: native 5-minute observations
 600 (Fig. 12(a–b)) and hourly maximum gusts aggregated from the same data (Fig. 12(a*–b*)).

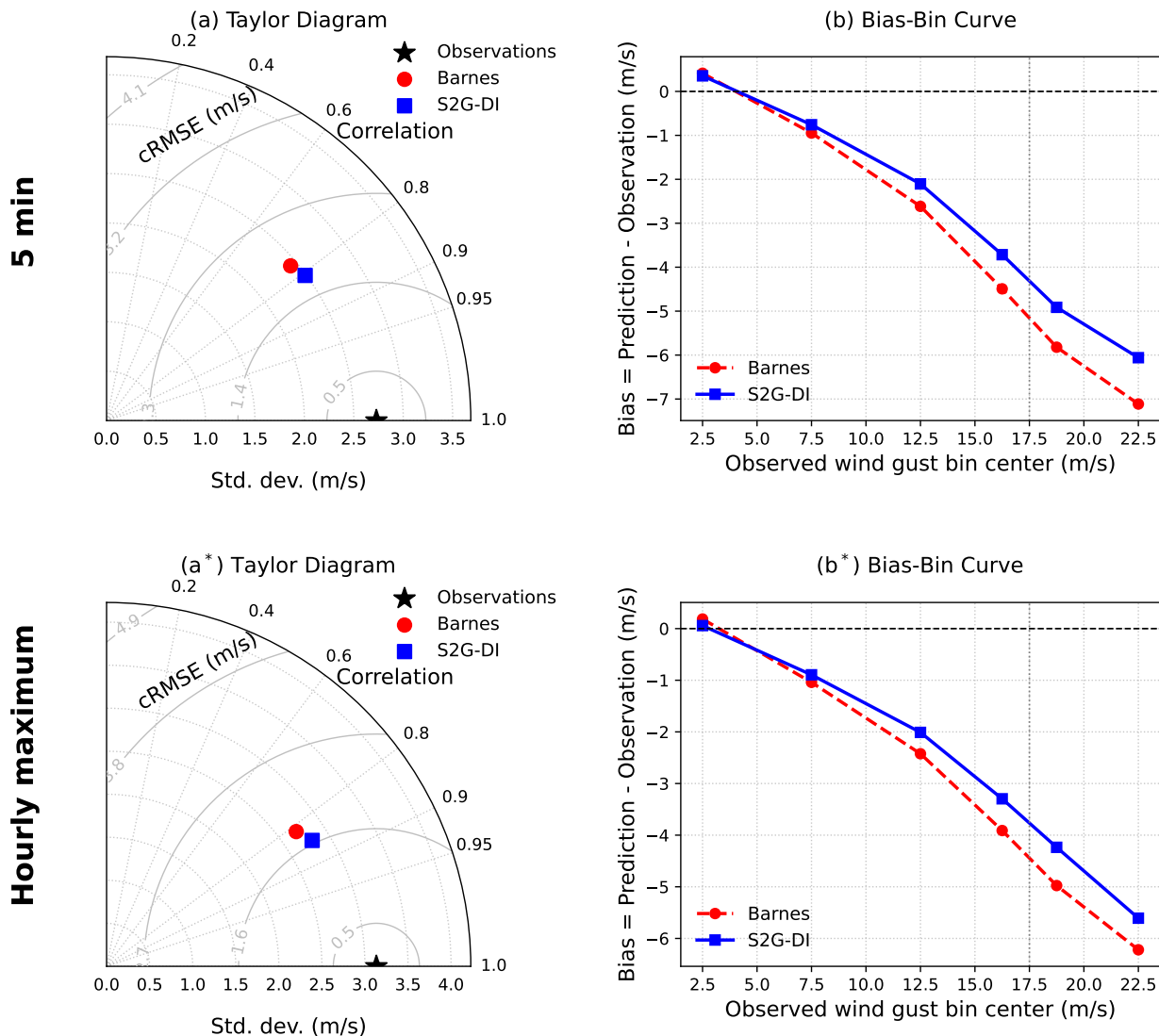


FIG. 12: Statistical evaluation of interpolation performance at held-out NYSM stations aggregated over the full year 2023 for two temporal resolutions. The top row shows results at native 5-minute resolution, while the bottom row shows results after aggregation to hourly maximum gusts. Panels (a, a*) display Taylor diagrams summarizing correlation, standard deviation, and centered RMSE, while (b, b*) show bias as a function of observed wind gust magnitude. All results are computed at held-out stations using out-of-fold predictions following the cross-validation strategy.

601 Considering the native 5-minute evaluation (Fig. 12(a-b)), the Taylor diagram summarizes the
 602 overall statistical performance, showing that S2G-DI achieves higher correlation with observations,
 603 a standard deviation closer to the observed variability, and a reduced centered RMSE compared
 604 to Barnes interpolation. The bias-bin analysis reveals systematic differences between the S2G-DI
 605 approach and the Barnes interpolation. While both approaches exhibit increasing underestimation

606 with higher wind gust magnitudes, S2G-DI consistently reduces the magnitude of this bias across
607 all bins, particularly in the extreme range beyond 17.5 m/s. These results collectively demonstrate
608 that the improvements of the deep interpolation framework are not limited to specific events but
609 persist across the full temporal domain, reinforcing its robustness and generalization capability
610 when applied to real observational data.

611 While both methods exhibit improved absolute performance at hourly resolution (Fig. 12(a*-b*))
612 due to the smoothing of short-lived gust variability, the relative performance between S2G-DI
613 and Barnes interpolation remains consistent across all metrics. In particular, the Taylor diagram
614 indicates unchanged ranking in terms of correlation and error characteristics and the bias-bin
615 structure is preserved. This consistency demonstrates that the spatial relationships learned from
616 hourly RTMA training data generalize well to sub-hourly inference, and that the performance gains
617 of the deep interpolation framework are robust to the temporal mismatch between training and
618 inference resolutions.

619 **7. Results: Evaluation against independent HRRR analysis**

620 To assess the dependence of the S2G-DI framework on the choice of reference dataset, we
621 perform an additional experiment using the HRRR dataset. In this experiment, HRRR is used
622 to as inference dataset, following the same experimental protocol as in the NYSM-based setup.
623 Specifically, sparse samples are extracted from HRRR fields to mimic station observations, and
624 both Barnes interpolation and deep interpolation are applied to reconstruct the full fields. The
625 generated fields are then compared with source HRRR dataset. Because HRRR is not used during
626 training, this setup provides an independent test of whether the relative performance of the methods
627 generalizes across datasets.

628 Figure 13 shows the annual mean wind gust fields and corresponding bias distributions of Barnes
629 interpolation and the S2G-DI model optimized relative to HRRR. Compared to the RTMA-based
630 evaluation in Fig. 5, both methods again reproduce the large-scale spatial patterns of the reference
631 field, indicating consistent behavior across datasets. However, as in the RTMA-based case, Barnes
632 interpolation produces spatially smooth fields with limited variability, while the S2G-DI fields retain
633 enhanced spatial structure that more closely resembles the variability present in the reference.

634 The bias maps further reveal systematic differences between the two approaches. Similar to
 635 the RTMA-based results, Barnes interpolation exhibits spatially coherent regions of over- and
 636 underestimation, particularly aligned with orographic features. In contrast, the S2G-DI reduces the
 637 magnitude of these structured biases and distributes residual errors more evenly across the domain.
 638 Notably, while the spatial patterns of bias differ between RTMA and HRRR—reflecting differences
 639 between the two reference products, the relative behavior of the methods remains consistent.

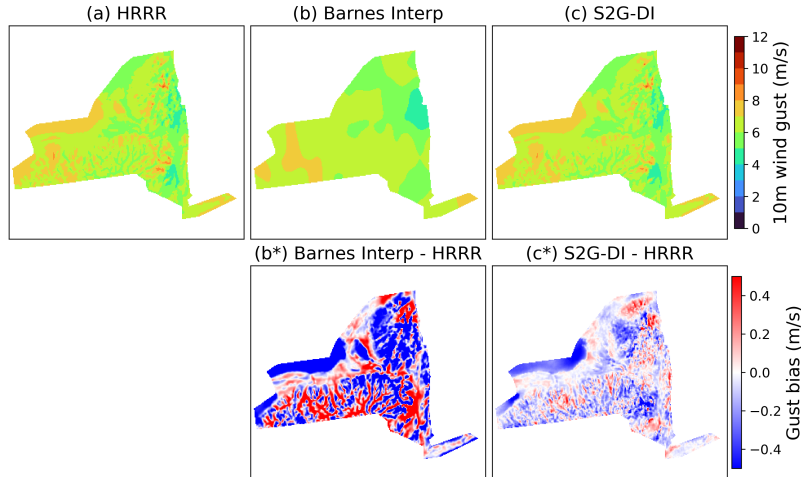


FIG. 13: Comparison of annual mean 10 m wind gust fields over New York State using (a) HRRR, (b) Barnes interpolation, and (c) S2G-DI. Panels (b*) and (c*) show the corresponding bias fields relative to HRRR. All fields are aggregated over the evaluation period 2023.

640 Table 5 quantifies these differences, showing that S2G-DI achieves substantial improvements
 641 over Barnes interpolation, including a 27.9% reduction in RMSE, a 7.3% increase in PSNR, and a
 642 2.7% increase in SSIM. These improvements are comparable in magnitude to those observed in the
 643 RTMA-based evaluation, indicating that the performance gains of the deep interpolation framework
 644 persist when the entire evaluation is conducted on an independent high-resolution dataset.

TABLE 5: Relative performance of S2G-DI compared to Barnes interpolation using HRRR as reference.

RMSE reduction (%)	PSNR improvement (%)	SSIM improvement (%)
27.9	7.3	2.7

645 **8. Results: Robustness across meteorological variables**

646 Finally, four additional runs are conducted for target variables other than wind gust to assess the
 647 generalization capability of the best-performing DL configuration. More specifically, the UNet
 648 with all input variables, all training years, and $n = 75$ random training stations is used to generate
 649 analysis fields of M_{10} , T_2 , and Q_2 . The performance compared to Barnes interpolation, with respect
 650 to RMSE, PSNR, and SSIM, is presented in Tab. 6. The difference in spatial patterns (not shown
 651 here) between S2G-DI and Barnes is qualitatively similar to the wind gust case discussed previously:
 652 Barnes interpolation misses terrain details and fine-scale features, which our framework captures,
 653 leading to better performance. S2G-DI achieves the highest performance for the T_2 fields, possibly
 654 because temperature strongly depends on terrain slopes and elevation, which S2G-DI captures
 655 well compared to Barnes interpolation. Specific humidity, on the other hand, exhibits a very fine
 656 granular structure, which is captured less well by S2G-DI, resulting in only 25% improvement in
 657 RMSE over Barnes. The deep interpolation performance for mean wind M_{10} is a bit better than
 658 for wind gust G_{10} . This performance difference is attributed to M_{10} fields behaving qualitatively
 659 similar to G_{10} fields while exhibiting less complex patterns and dynamics, which we consider
 660 easier to predict. All in all, we conclude that our S2G-DI is well-suited to generate analysis fields
 661 of various variables.

TABLE 6: Deep interpolation performance of S2G-DI compared to traditional Barnes interpolation for 10 m wind speed (M_{10}), 2 m temperature (T_2), and 2 m specific humidity (Q_2). G_{10} scores are repeated for reference.

Variable	Method	RMSE (lower better)	Reduction (%)	PSNR (higher better)	Improvement (%)	SSIM (higher better)	Improvement (%)
M_{10} (m/s)	Barnes	1.269		31.362		0.893	
	S2G-DI	0.710	44.0%	35.818	14.2%	0.951	6.5%
T_2 (K)	Barnes	1.271		53.764		0.999	
	S2G-DI	0.548	56.9%	60.828	13.1%	1.000	0.1%
Q_2 (kg/kg)	Barnes	4e-4		42.090		0.985	
	S2G-DI	3e-4	29.8%	45.240	7.5%	0.993	0.8%
G_{10} (m/s)	Barnes	1.309		34.471		0.934	
	S2G-DI	0.888	32.1%	37.562	9.0%	0.966	3.4%

662 **9. Conclusions and Future Scope**

663 In this study, we presented the *Sparse-to-Gridded Deep Interpolation (S2G-DI)* framework,
664 which aims to obtain high-quality gridded analysis fields of meteorological variables from sparse
665 observations coming from a network of weather stations. The primary focus of our study was on
666 10 m wind gusts. Wind gusts are rapid changes in wind speed or short-duration maxima, typically
667 defined using the three-second average wind speed maxima (World Meteorological Organization
668 (WMO) 2023). These are associated with diverse high-impact weather phenomena, including
669 tornadoes, thunderstorms, cyclones, and other high-impact weather systems (Sheridan 2018).
670 Gusts can severely damage infrastructure, energy systems, transportation, and fire management
671 operations. They may topple transmission lines, damage wind and solar farms, and bring down
672 trees or branches onto power infrastructure (Mitchell 2013). Such failures can spark wildfires,
673 which strong gusts can then spread rapidly over large areas (Mass and Ovens 2021). Wind gusts
674 also pose hazards to aviation (Gultepe et al. 2019). While cyclones can sustain damaging winds
675 for extended periods, convective systems often produce the most extreme short-duration gusts,
676 sometimes exceeding 25 m s^{-1} (El Rafei et al. 2023). Capturing these short-lived and spatially
677 localized events, therefore, remains a major challenge for both observation networks and modeling
678 systems. We also demonstrated good performance for 2 m temperature, 2 m specific humidity, and
679 10 m wind speed.

680 Our S2G-DI is a two-step framework: first, sparse observations from weather networks, such as
681 the New York State Mesonet (Brotzge et al. 2020), are nearest-neighbor-interpolated to gridded
682 but low-fidelity 2D fields. Next, these low-fidelity inputs are passed through a deep learning (DL)
683 model, which produces refined analysis fields by considering interactions between meteorological
684 variables and accounting for terrain effects. The DL model for refinement is obtained through a
685 decoupled precursor step, where it is trained on a gridded high-resolution proxy dataset covering
686 the area of the observational network. We make the core assumption that the local meteorology
687 is captured sufficiently well in the proxy dataset, such that later inference with a different dataset,
688 i.e., the refinement of the observed data, produces realistic spatial fields. This assumption was
689 confirmed throughout the study, with the help of HRRR analysis acting as an independent evaluation
690 dataset. Our deep interpolation is also shown to be robust against missing data from a subset of
691 weather stations, which is common due to maintenance or instrument malfunctions.

692 Decoupling the training of the DL model on proxy data from the deep interpolation (inference)
693 on observations means the model never uses both datasets at the same time. Therefore, the gridded
694 training data and observations do not require a temporal overlap and can have different temporal
695 resolutions. For example, the DL model can be trained on archive data with a 1 h time step
696 (common for gridded numerical data), while later S2G-DI can use real-time observations at the
697 common 5-10 min time steps. Additionally, the common spatiotemporal misalignment between
698 numerical and observational datasets is irrelevant, as only observations are used during operational
699 deep interpolation. We believe these distinct features are highly valuable for the community.

700 In a thorough sensitivity study, we found that a medium-complexity UNet performs best compared
701 to a lower-complexity deep convolutional neural network or a high-complexity transformer-based
702 UNet. The quality of the analysis fields also increased when auxiliary variables, such as temperature
703 and wind speed, were added as inputs. We attribute this to meteorological interactions between
704 inputs, resulting in a performance increase in predicting the analysis target, such as wind gust.
705 To demonstrate how powerful our deep interpolation is, we benchmarked it against the commonly
706 used traditional Barnes interpolation (Barnes 1964). In comparison with Barnes interpolation, our
707 framework not only achieved higher scores but also yielded more spatially detailed analysis fields
708 because it accounts for terrain.

709 As gridded proxy data are crucial for training the DL model, we assessed how performance
710 depends on the available amount of training data. We found that even just one year of gridded data
711 yields a model that significantly outperforms Barnes interpolation and that adding more data has a
712 positive but slowly diminishing effect. This result stresses that our framework is also applicable for
713 regions where large gridded datasets are not readily available. While generating at least one year
714 of training data using numerical weather prediction models is still a notable computational effort,
715 we view it as generally feasible for the community. As previously mentioned, the decoupling of
716 training and inference/deep interpolation enables training of the DL model on gridded data from
717 archives, provided that the distributions of the gridded data match those of the observations.

718 Similar to the performance dependency on the temporal amount of training data, it is crucial
719 to check the quality of the generated analysis fields when observations are only available at a
720 subset of stations. Missing data are common in real-world networks due to sensor malfunctions or
721 unavailability during maintenance. From a technical perspective, the S2G-DI framework inherently

722 allows deep interpolation from varying numbers of stations, as the nearest-neighbor interpolation
723 (first step) can always generate low-fidelity gridded fields as input for the DL model. As part of the
724 sensitivity study, we demonstrated that S2G-DI retains a significant performance advantage over
725 Barnes interpolation at varying degrees of missing stations. More specifically, we outperformed
726 Barnes interpolation in the same missing data situations, but also when Barnes interpolation had
727 access to *all* observations, while S2G-DI faced missing data. We found that the crucial step to
728 achieve this level of robustness is to randomly drop stations already during training. Dropping
729 stations at training time was found to be crucial, as it ensures that the DL model sees a variety of
730 network configurations and learns to robustly generate accurate results. Overall, we believe that
731 the high performance of deep interpolation even under missing data makes our framework practical
732 and attractive for operational use.

733 To assess how our optimized S2G-DI model handles challenging situations and real-world
734 observations, we analyzed three observed extreme wind gust events that occurred in New York
735 State in 2023. Extreme gusts are defined as gust magnitudes exceeding 17.5 m/s. Evaluating
736 the time series at held-out stations showed that our model outperformed the traditional Barnes
737 interpolation in many cases. We attribute this to the fact that the deep learning model can capture
738 spatial meteorological patterns and variable interactions in contrast to the recursive weighted
739 interpolation used by Barnes.

740 However, while our S2G-DI framework captures the timings and general trends of the events well,
741 it misses the magnitude of the extreme gust peaks in many cases. As less than 0.5% of observed
742 yearly gust samples at each station have gust magnitudes larger than our threshold, the extreme
743 gust estimation is a highly imbalanced regression problem. Imbalanced data pose a common
744 challenge in machine learning (He and Garcia 2009; Branco et al. 2016), as standard models
745 and training procedures tend to favor the center of data distributions rather than the tails where
746 extreme events are located. One strategy is to adjust the modeling procedure to account for the data
747 imbalance, for example, through over- or undersampling techniques, by assigning higher weights
748 to rarer samples, or by specifically adjusting the model architecture (Branco et al. 2016). But the
749 imbalance issue can also be addressed explicitly from a physics perspective. In future work, we
750 aim to incorporate weather radar data, such as from the Multi-Radar/Multi-Sensor system (Smith
751 et al. 2016). Radar data would provide crucial information about thunderstorms, which are directly

752 associated with the convective processes that generate extreme convective wind gusts (Sheridan
753 2018). In other words, strong radar echoes are expected to appear only during extreme gust events,
754 which form the tail of the gust distribution. Consequently, we expect to enhance the model's ability
755 to capture the distribution tails and, thus, its ability to capture extreme gust events, by utilizing this
756 complementary dataset.

757 Overall, our S2G-DI was shown to outperform the traditional Barnes interpolation in different
758 tests for wind gust, wind magnitude, 2 m temperature, and 2 m specific humidity. Further, the setup
759 of the S2G-DI framework is generic, so we expect it can be trained with alternative mesoscale
760 products and geographical regions. For example, preliminary results Martinez-Roig (2025) showed
761 the successful application of S2G-DI to interpolate observed 10m wind speed in Spain based
762 on the New European Wind Atlas Dörenkämper et al. (2020). We, therefore, do not see any
763 fundamental limitations that would prevent the application of our approach to other regions,
764 datasets, or observations.

765 The framework possesses several practical features, including decoupling between training and
766 interpolation, efficiency with respect to the volume of gridded training data required, and the
767 capability to robustly handle missing observations while maintaining high performance. We
768 believe that our approach to obtaining analysis fields is a valuable addition for the community, as
769 many countries and states operate networks of weather stations. In the future, we plan to expand
770 this work to capture extreme events more reliably, and we view the S2G-DI framework as a first
771 step toward purely observation-driven forecasting, which we also actively research.

772 *Acknowledgments.* The authors acknowledge support from the Weather Innovation and Smart
773 Energy and Resilience (WISER) center, funded by the U.S. National Science Foundation Indus-
774 try–University Cooperative Research Centers (IUCRC) program. We thank the New York State
775 Mesonet (NYSM) for providing the observational data used in this research. Original funding for
776 the NYSM buildup was provided by Federal Emergency Management Agency grant FEMA-4085-
777 DR-NY. The continued operation and maintenance of the NYSM is supported by National Mesonet
778 Program, University at Albany, Federal and private grants, and others. The authors are grateful to
779 Daniel Kirk-Davidoff and Evan Duffey for useful feedback on an earlier version of this work.

780 *Data availability statement.* The RTMA dataset is publicly available on Amazon Web Ser-
781 vices (AWS) and can be accessed using the AWS CLI: `s3://noaa-urma-pds/`. The HRRR
782 hourly analysis is publicly available on AWS and can be accessed using the AWS CLI: `s3:`
783 `//noaa-hrrr-bdp-pds/`. The NYSM observations are hosted at <https://nysmesonet.org>
784 and can be requested through their data request page at [https://nysmesonet.org/weather/](https://nysmesonet.org/weather/requestdata)
785 `requestdata`. The deep learning model training and inference scripts are available in the Git-
786 Lab project **S2G-DI** at <https://gitlab.com/HarishBaki/S2G-DI.git>, which can be freely
787 accessed.

788 APPENDIX A

789 **Additional Background on Deep Learning Methods**

790 In recent years, data-driven approaches, encompassing both classical machine learning (ML) and
791 deep learning (DL), have gained increasing popularity as a promising alternative for generating
792 high-resolution gridded meteorological fields from a combination of observations and model sim-
793 ulations (Franco et al. 2020; Boomgard-Zagrodnik and Brown 2022; Chen et al. 2024; Jeong et al.
794 2025; Dujardin and Lehning 2022). Depending on the input–output structure, these approaches can
795 be categorized into *point-to-point* learning, *image-to-point* learning, and *image-to-image* learning
796 frameworks.

797 In a *point-to-point* learning framework, each input is expressed as a feature vector of station-level
798 attributes such as spatial coordinates, distances between stations, or other predictors. This setup
799 allows interpolation at any location and has been applied to several meteorological variables. For
800 temperature and precipitation, studies have used a variety of methods, including artificial neural net-

801 works in the Virtual Weather Station framework (Franco et al. 2020), random forests with mesonet
802 data (Boomgard-Zagrodnik and Brown 2022), and more recent deep neural network–based kriging
803 approaches (Chen et al. 2024). Hybrid designs have also appeared, for example, combining sparse
804 and dense station networks before training LSTM models to reconstruct daily maximum/minimum
805 temperatures and precipitation (Jeong et al. 2025). Similar ideas have been explored for wind
806 gusts: early work used ANNs with station variables (Sallis et al. 2011), followed by the use of
807 reanalysis predictors (Coburn and Pryor 2022), and later decision trees, ensemble methods, and
808 neural network–Gaussian process combinations (Kartal et al. 2023; Jahan et al. 2024; Zanetta et al.
809 2025). While these approaches work well for filling gaps at specific stations, they are less effective
810 at representing spatial correlations in areas with complex terrain.

811 In an *image-to-point* learning framework, each input sample is represented by a gridded image
812 containing inherent coordinate information, while predictions are made for a single grid point. For
813 example, Dujardin and Lehning (2022) used a convolutional neural network (CNN) to predict u
814 and v wind components at target locations, given COSMO1 model outputs within a 19×19 pixel
815 window and high-resolution topography from a digital elevation model (DEM), demonstrating the
816 added value of topography inputs for capturing spatial correlations. Their approach, however,
817 focused on downscaling model-simulated fields rather than mapping station observations to grids.
818 While effective for representing local spatial correlations, the limited spatial context restricts the
819 framework’s ability to capture mesoscale phenomena such as cyclones, fronts, and thunderstorms.

820 In contrast, the *image-to-image* learning framework treats both inputs and targets as structured
821 images, enabling explicit modeling of spatial correlations and interactions between neighboring
822 locations, thus producing spatially consistent interpolations. Fukami et al. (2021) demonstrated
823 this approach by downsampling a two-dimensional (2D) target field to sensor locations, projecting
824 it onto a structured grid via Voronoi tessellation, and training a CNN to map the low-quality
825 field to the original. Applications included wake flow reconstruction, sea surface temperature
826 (SST) estimation, and turbulent channel flow modeling. Subsequent studies extended this concept:
827 Zhao et al. (2024) replaced the CNN with a Fourier neural operator, Wang et al. (2024) used
828 convolutional LSTMs for joint SST reconstruction and forecasting, and Sunderhaft et al. (2024)
829 applied it to Antarctic surface temperatures from reanalysis. Similarly, Güemes et al. (2022)
830 replaced Voronoi tessellation with cubic binning and employed a GAN for SST reconstruction.

831 These works highlight the suitability of DL-based *image-to-image* learning for mapping sparse
832 observations to gridded fields and its potential for real-time use, as shown by Achermann et al.
833 (2024) in volumetric wind speed prediction over complex terrain ($1.5 \times 1.5 \text{ km}^2$) from sparse
834 uncrewed aerial vehicle observations.

835 Building on the *image-to-image* framework proposed by Fukami et al. (2021), which projects
836 sparse observations onto a structured grid before refining them with a deep learning model, we
837 adapt and extend this approach for the real-time reconstruction of high-resolution gridded wind
838 gust fields. While prior applications have focused on variables such as sea surface temperature
839 (Güemes et al. 2022; Wang et al. 2024) and wind speed over small complex-terrain domains
840 (Achermann et al. 2024), they have not addressed the unique challenges of generating high-quality
841 gridded wind gust analyses across a large, topographically complex region such as New York State.
842 Our S2G-DI framework bridges this gap through a decoupled training–inference design, enabling
843 real-time gridded wind gust field generation from sparse mesonet observations at a 5 min temporal
844 resolution, while remaining adaptable to other meteorological variables, observational networks,
845 and geographical regions.

846 APPENDIX B

847 Deep Learning Models

848 *a. DCNN*

849 The DCNN architecture, illustrated in Fig. B1, is relatively simple and takes the same input and
850 target tensors as the UNet model. It comprises seven sequential layers, each consisting of a 2D
851 convolutional layer with a 7×7 filter and 48 feature channels, followed by a GeLU activation. A
852 final 2D convolutional layer with a 3×3 filter and a single output channel, without any activation
853 function, is applied to produce the final output. The architecture has been adopted from the studies
854 of Fukami et al. (2021) and Sunderhaft et al. (2024), with little modifications made. Since our
855 objective has been to predict at locations other than the stations, the values at the station locations
856 in the predictions are replaced with those of the input tensor. This is to make sure that the values at
857 the station locations will always be the same as what we provide as inputs. The entire architecture
858 contains approximately 0.69 million trainable parameters.

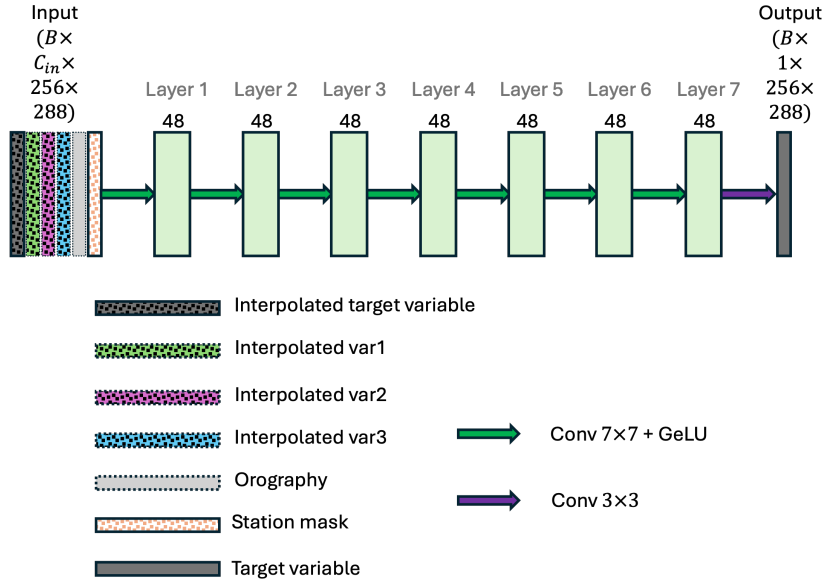


FIG. B1: An illustration of the Deep convolutional Neural Network (DCNN) model architecture.

859 *b. UNet*

860 The U-Net architecture is shown in Fig. B2, has been adopted from Wang et al. (2022) with little
 861 modifications. It consists of five hierarchical layers and follows an encoder–decoder structure.
 862 The encoder progressively extracts features by reducing the spatial resolution and increasing the
 863 feature dimensionality through a series of ConvBlocks and DownBlocks. In contrast, the decoder
 864 reconstructs the full-resolution output by gradually increasing the spatial dimensions and reducing
 865 the feature depth via a series of UpBlocks and ConvBlocks.

866 The core building block of the network is the ConvBlock, which serves as the backbone of the
 867 architecture. Each ConvBlock consists of a 2D convolutional layer with a 3×3 kernel, followed
 868 by a GeLU activation, another 3×3 convolution, and a second GeLU activation. To mitigate
 869 overfitting, a dropout layer with a rate of 0.2 and a drop-path layer with a rate of 0.2 are included
 870 as regularization components. A residual skip connection is used, which is added after a final
 871 convolutional layer of 1×1 filter. This residual skip connection facilitates improved gradient flow
 872 during training and helps alleviate vanishing gradient problems, following the approach introduced
 873 in the original Residual Network (ResNet) architecture (He et al. 2015).

874 The DownBlock reduces the spatial resolution by half while preserving the feature dimensionality.
 875 It consists of a 2D convolutional layer with a 4×4 kernel and a stride of 2, using the same number of

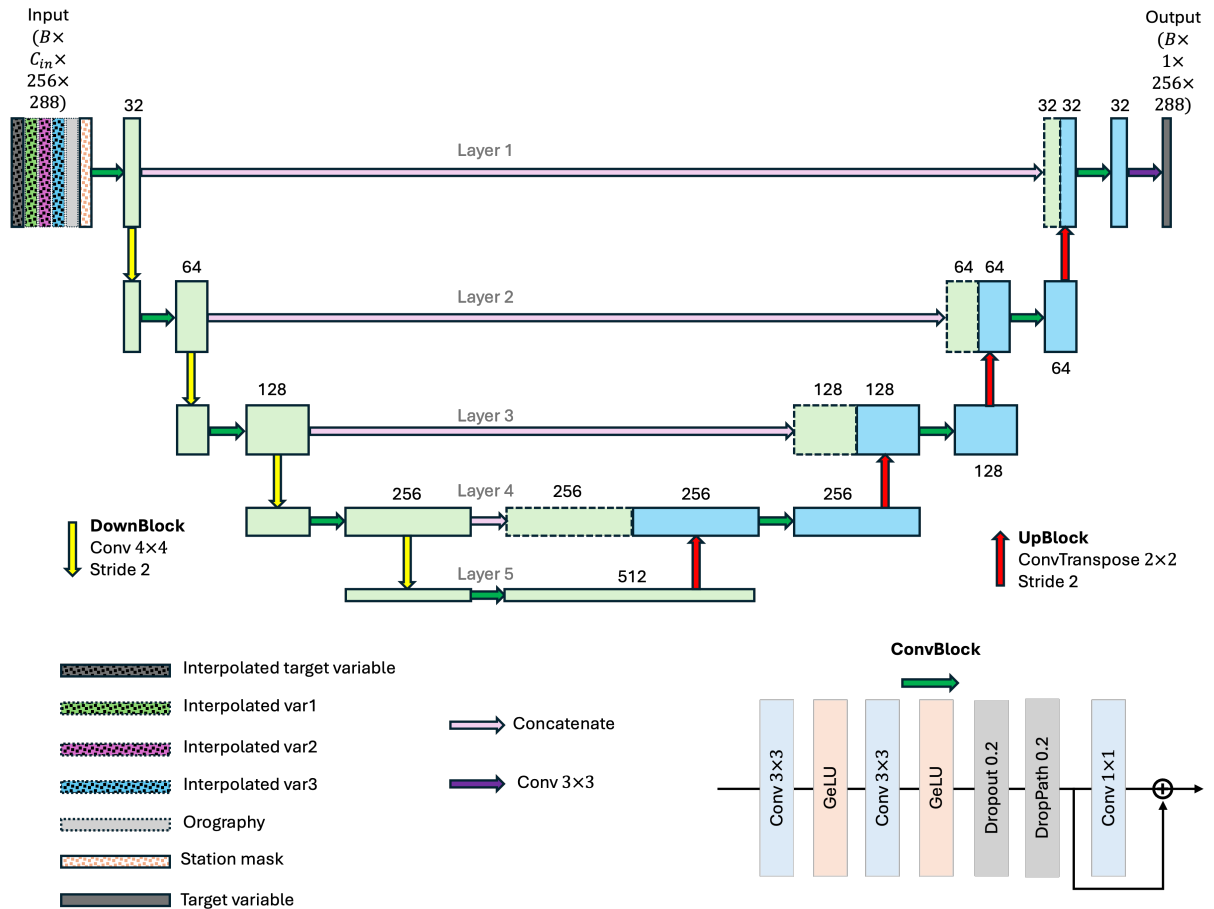


FIG. B2: An illustration of the UNet model architecture employed in this study.

876 input and output channels. This approach retains more contextual features compared to commonly
 877 used pooling techniques (He et al. 2015). Conversely, the UpBlock increases the spatial resolution
 878 by a factor of two while reducing the feature dimension by half. It is implemented using a 2D
 879 transposed convolutional layer with a 2×2 kernel and a stride of 2, where the number of output
 880 channels is half the number of input channels. While the DownBlock operates solely on the
 881 spatial dimensions and maintains feature dimensionality, the UpBlock simultaneously modifies
 882 both spatial and feature dimensions.

883 In the first layer, the input is passed through a ConvBlock with 32 filters, producing an output
 884 with the same spatial dimensions as the input but with 32 channels. The output of this layer is
 885 retained for a future skip connection. Simultaneously, it is passed through a DownBlock, which
 886 reduces the spatial resolution by half. This sequence of applying a ConvBlock followed by a

887 DownBlock is repeated up to layer 5, with the number of filters doubled at each successive layer.
888 By layer 5, the feature map reaches a channel depth of 512, while the spatial resolution is reduced
889 to 16×18 . Together, these operations constitute the encoder portion of the U-Net architecture.

890 The output of the encoder portion is passed through an UpBlock, which increases the spatial
891 resolution by a factor of 2 and reduces the feature dimension by half. The upsampled output is then
892 concatenated along the feature dimension with the corresponding ConvBlock output from layer 4
893 of the encoder (via a skip connection), resulting in a feature map with doubled channel depth at that
894 layer. The skip connection helps in extracting high-level localized contextual information along
895 with mitigating vanishing gradients in deeper models. This concatenated output is subsequently
896 passed through a ConvBlock, which reduces the feature dimensionality by half (e.g., to 256
897 channels). This sequence of applying an UpBlock, concatenation, and ConvBlock is repeated up
898 to layer 1, with the number of features halving at each successive layer. By the time the decoder
899 reaches layer 1, the feature map has a channel depth of 32, and the spatial resolution is restored
900 to the original input dimensions of 256×288 . Together, these series of operations constitute the
901 decoder portion of the U-Net.

902 Finally, a 2D convolutional layer with a 3×3 kernel, a single output feature, and no activation
903 function is applied to produce the final output, which consists of one channel corresponding to the
904 target variable. Same as the DCNN, the values at the station locations are replaced from the input
905 tensor. This UNet architecture consists of around 9.5 million parameters in total.

906 *c. SwinT2UNet*

907 Among all the models, the SwinT2UNet (shown in Fig. B3) is the most complex architecture,
908 has been adopted from Wang et al. (2022) with considerable modifications. It follows the same
909 encoder–decoder structure as the UNet; however, the backbone is replaced with a SwinBlock,
910 along with several additional architectural modifications.

911 The SwinBlock is, in fact, the SwinTransformerV2Block from , which sequentially applies
912 a non-shifted window and a shifted-window mechanism. Each block consists of a window-
913 based multi-headed self-attention (W-MSA) module, followed by layer normalization, a residual
914 connection, another layer normalization, a multilayer perceptron (MLP), and a second residual
915 connection. This is followed by a shifted-window multi-headed self-attention (SW-MSA) module,

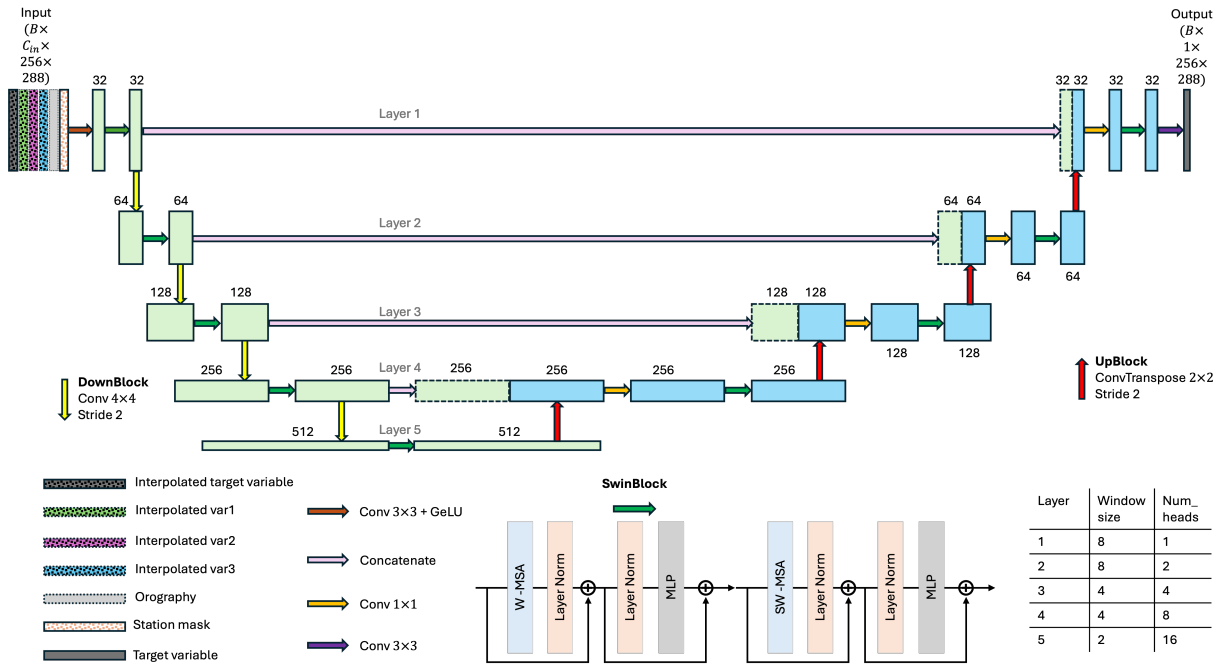


FIG. B3: An illustration of the UNet with Sifted window multiheaded self-attention transformer block (Swintransformer) as backbone (SwinT2UNet) model architecture.

916 again followed by layer normalization, a residual connection, another layer normalization, an MLP,
 917 and a final residual connection.

918 Unlike in the UNet architecture, the DownBlock reduces the spatial resolution by half while
 919 doubling the feature dimensionality. It consists of a 2D convolutional layer with a 4×4 kernel and
 920 a stride of 2, where the number of Conv 1×1 output channels is twice that of the input channels. On the other
 921 hand, the UpBlock is identical to that of the UNet, it increases the spatial resolution by a factor of
 922 two while reducing the feature dimension by half.

923 In the first layer, the input is passed through an input projection block, which consists of a 2D
 924 convolutional layer with a 3×3 kernel and 32 feature channels, followed by a GeLU activation.
 925 This block transforms the input into the required number of feature representations for subsequent
 926 operations. The output from the input projection block is then passed through a SwinBlock with
 927 a window size of 8 and one attention head, which preserves the spatial and feature dimensions
 928 while extracting key hierarchical representations. The output from this layer is retained for a future
 929 skip connection. Simultaneously, it is passed through a DownBlock, which reduces the spatial
 930 resolution by half and doubles the number of feature channels. Similar to the UNet encoder, this

931 sequence consisting of a SwinBlock followed by a DownBlock is repeated up to layer 5, with
932 the number of filters doubling at each successive layer. The specific window sizes and number of
933 attention heads used in each SwinBlock are summarized in the table shown in Fig. B3. By layer
934 5, the spatial resolution is reduced to 16×18 , and the feature dimensionality increases to 512.

935 Similar to the UNet, the output of the encoder portion is passed through an UpBlock and
936 followed by a concatenation of SwinBlock output from encoder layer 4. The resulting output
937 having a feature dimension of 512 goes through a 2D convolutional layer with 1×1 filter, to
938 reduced the features by half, which is needed since the following SwinBlock takes and gives same
939 number of (256) features. This sequence of applying UpBlock, concatenation, feature reduction
940 through 1×1 convolution, and SwinBlock is repeated up to layer 1, with the number of features
941 halving at each successive layer. By the time the decoder reaches layer 1, the feature map has
942 a channel depth of 32, and the spatial resolution is restored to the original input dimensions of
943 256×288 . Together, these series of operations constitute the decoder portion of the SwinT2UNet.

944 Similar to the UNet, the output of the encoder portion is passed through an UpBlock, followed
945 by a concatenation with the corresponding SwinBlock output from encoder layer 4. The resulting
946 feature map, with a dimensionality of 512, is then passed through a 2D convolutional layer with a
947 1×1 kernel to reduce the number of features by half. This step is necessary because the subsequent
948 SwinBlock expects and produces 256 features. This sequence of operations, i.e., UpBlock,
949 concatenation, feature reduction via 1×1 convolution, and SwinBlock, is repeated up to layer 1,
950 with the number of features halving at each successive layer. By the time the decoder reaches layer
951 1, the feature map has a channel depth of 32, and the spatial resolution is restored to the original
952 input dimensions of 256×288 . Together, these operations constitute the decoder portion of the
953 SwinT2UNet architecture.

954 Finally, a 2D convolutional layer with a 3×3 kernel, a single output channel, and no activation
955 function is applied to the output of the decoder to produce the final output, which consists of one
956 channel corresponding to the target variable. Here as well, the values at the station locations are
957 replaced from the input tensor. This SwinT2UNet architecture consists of around 14.24 million
958 parameters in total.

Barnes interpolation

Barnes Objective Analysis is a widely used two-pass interpolation method that applies Gaussian distance weighting to observations surrounding each grid point (Barnes 1964; Koch et al. 1983). In the first pass, the interpolated value at a grid point (i, j) is obtained by a weighted average of the observations $f(x_m, y_m)$, where each observation is weighted based on its distance r_m from the grid point. This is expressed as:

$$g_0(i, j) = \frac{\sum_{m=1}^M w_m f(x_m, y_m)}{\sum_{m=1}^M w_m}, \quad (\text{C1})$$

where the weight is given by $w_m = \exp\left(-\frac{r_m^2}{\kappa}\right)$. Here, κ is a smoothing parameter with units of length squared, defined as $\kappa = \kappa^* \left(\frac{2\Delta n}{\pi}\right)^2$, where κ^* is a dimensionless smoothing coefficient and Δn is a representative station spacing.

The second pass refines the field by correcting for the local bias between the observed values and the first-pass estimates. This correction is expressed as:

$$g_1(i, j) = g_0(i, j) + \frac{\sum_{m=1}^M w'_m [f(x_m, y_m) - g_0(x_m, y_m)]}{\sum_{m=1}^M w'_m}, \quad (\text{C2})$$

where the corrected weights are defined as $w'_m = \exp\left(-\frac{r_m^2}{\gamma\kappa}\right)$, and γ is a convergence (or damping) parameter, typically set between 0.2 and 1.0, to sharpen the influence of nearby residuals in the correction pass.

We adopted the Barnes interpolation implementation provided by MetPy (available at: https://unidata.github.io/MetPy/latest/api/generated/metpy.interpolate.interpolate_to_points.html), which accepts the convergence parameter (γ) and the dimensionless smoothing parameter (κ^*) as input arguments. The default values of γ and κ^* are 0.25 and 5.052, respectively. However, the values of these two parameters have been shown to influence the

979 interpolation accuracy considerably (Sun and Crook 2001). Thus, it is of paramount importance
 980 to first identify the optimal values of these parameters. For the purpose, we sampled (γ, κ^*) in
 981 the ranges of $[0.05, 0.4]$ and $[3, 7]$ for γ and κ^* , respectively, using Sobol sampling design. Using
 982 the sampled parameter values in the Barnes interpolation, we interpolated G_{10} for the period of
 983 2023, and computed RMSE, PSNR, and SSIM scores, in comparison to the RTMA data. From the
 984 scores, we found a new optimal value for the parameters as $\gamma = 0.264$ and $\kappa^* = 6.566$, as shown in
 Fig. C1.

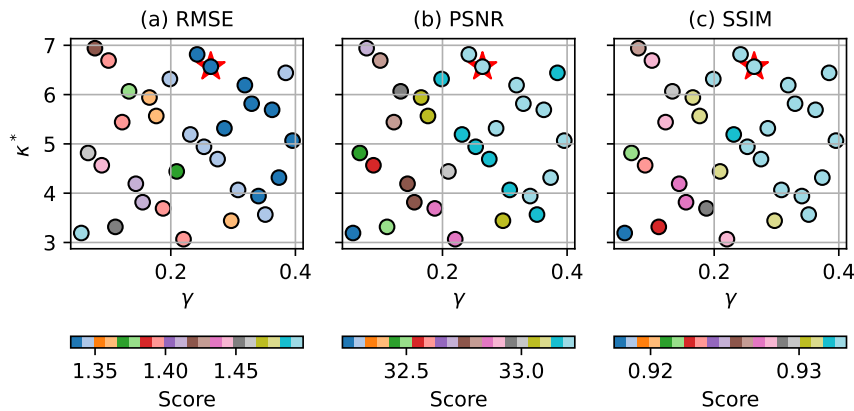


FIG. C1: Performance evaluation of Barnes interpolation accuracy across different values of the smoothing parameter κ^* and convergence parameter γ , sampled using a Sobol sequence. The subplots show the variation of (a) RMSE, (b) PSNR, and (c) SSIM scores between interpolated G_{10} and RTMA data for the year 2023. The red star in each panel denotes the location of the optimal parameter combination ($\gamma = 0.264176, \kappa^* = 6.566287$).

985

986 References

987 Achermann, F., T. Stastny, B. Danciu, A. Kolobov, J. J. Chung, R. Siegwart, and N. Lawrance,
 988 2024: WindSeer: real-time volumetric wind prediction over complex terrain aboard a small
 989 uncrewed aerial vehicle. *Nat. Commun.*, **15** (1), 3507.

990 Amponsah, W., P. A. Ayril, B. Boudevillain, C. Bouvier, I. Braud, P. Brunet, and Borga, 2018:
 991 Integrated high-resolution dataset of high-intensity european and mediterranean flash floods.
 992 *Earth System Science Data*, **10** (4), 1783–1794.

993 An, S., T.-J. Oh, E. Sohn, and D. Kim, 2025: Deep learning for precipitation nowcasting: A survey
 994 from the perspective of time series forecasting. *Expert Syst. Appl.*, **268** (126301), 126 301.

- 995 Anagun, Y., S. Isik, and E. Seke, 2019: Srlibrary: Comparing different loss functions for super-
996 resolution over various convolutional architectures. *Journal of Visual Communication and Image*
997 *Representation*, **61**, 178–187.
- 998 Barker, D., and Coauthors, 2012: The weather research and forecasting model’s community
999 variational/ensemble data assimilation system: Wrfda. *Bulletin of the American Meteorological*
1000 *Society*, **93 (6)**, 831–843.
- 1001 Barnes, S. L., 1964: A technique for maximizing details in numerical weather map analysis.
1002 *Journal of Applied Meteorology*, **3 (4)**, 396–409.
- 1003 Blankenau, P. A., A. Kilic, and R. Allen, 2020: An evaluation of gridded weather data sets for the
1004 purpose of estimating reference evapotranspiration in the united states. *Agric. Water Manag.*,
1005 **242 (106376)**, 106 376.
- 1006 Boomgard-Zagrodnik, J. P., and D. J. Brown, 2022: Machine learning imputation of missing
1007 mesonet temperature observations. *Comput. Electron. Agric.*, **192 (106580)**, 106 580.
- 1008 Branco, P., L. Torgo, and R. P. Ribeiro, 2016: A Survey of Predictive Modeling on Imbalanced
1009 Domains. *ACM Comput. Surv.*, **49 (2)**, 31:1–31:50, <https://doi.org/10.1145/2907070>.
- 1010 Brotzge, J. A., and Coauthors, 2020: A technical overview of the new york state mesonet standard
1011 network. *J. Atmos. Ocean. Technol.*, **37 (10)**, 1827–1845.
- 1012 Cai, Y., H. He, and Z. He, 2024: Loss functions analysis of performance improvements in single-
1013 image super-resolution. *IGARSS 2024-2024 IEEE International Geoscience and Remote Sensing*
1014 *Symposium*, IEEE, 3138–3141.
- 1015 Campbell Scientific, Inc., 2023: Mesonets. Accessed: 7 August 2025, <https://www.campbellsci.com.br/mesonets>.
- 1017 Charbonnier, P., L. Blanc-Feraud, G. Aubert, and M. Barlaud, 2002: Two deterministic half-
1018 quadratic regularization algorithms for computed imaging. *Proceedings of 1st International*
1019 *Conference on Image Processing*, IEEE Comput. Soc. Press, Vol. 2, 168–172 vol.2.
- 1020 Chen, K., E. Liu, M. Deng, X. Tan, J. Wang, Y. Shi, and Z. Wang, 2024: DKNN: deep kriging
1021 neural network for interpretable geospatial interpolation. *Geogr. Inf. Syst.*, **38 (8)**, 1486–1530.

- 1022 Chen, M., W. Shi, P. Xie, V. B. S. Silva, V. E. Kousky, R. Wayne Higgins, and J. E. Janowiak,
1023 2008: Assessing objective techniques for gauge-based analyses of global daily precipitation. *J.*
1024 *Geophys. Res.*, **113 (D4)**.
- 1025 Chen, Y., J. Hall, W. D. Van, N. Andela, S. Hantson, L. Giglio, and Randerson J, 2023: Multi-
1026 decadal trends and variability in burned area from the 5th version of the global fire emissions
1027 database (GFED5). *Earth System Science Data Discussions*, **2023**, 1–52.
- 1028 Christiansen, M. B., W. Koch, J. Horstmann, C. B. Hasager, and M. Nielsen, 2006: Wind resource
1029 assessment from c-band sar. *Remote Sensing of Environment*, **105 (1)**, 68–81.
- 1030 CNY Central, 2023a: Gusty wind causes minor tree damage around area, no
1031 widespread power outages. Accessed: 2025-07-28, [https://cnycentral.com/news/local/
1032 gusty-wind-causes-minor-tree-damage-around-area-no-widespread-power-outages](https://cnycentral.com/news/local/gusty-wind-causes-minor-tree-damage-around-area-no-widespread-power-outages).
- 1033 CNY Central, 2023b: Strong wind gusts continue for cny tonight with rain changing to
1034 snow showers into sunday. Accessed: 2025-07-28, [https://cnycentral.com/news/instagram/
1035 strong-wind-gusts-continue-for-cny-tonight-with-rain-changing-to-snow-showers-into-sunday](https://cnycentral.com/news/instagram/strong-wind-gusts-continue-for-cny-tonight-with-rain-changing-to-snow-showers-into-sunday). ■
- 1036 Coburn, J., and S. C. Pryor, 2022: Do machine learning approaches offer skill improvement for
1037 short-term forecasting of wind gust occurrence and magnitude? *Weather Forecast.*, **37 (5)**,
1038 525–543.
- 1039 Colle, B. A., and D. R. Novak, 2010: The new york bight jet: Climatology and dynamical evolution.
1040 *Mon. Weather Rev.*, **138 (6)**, 2385–2404.
- 1041 Cressman, G. P., 1959: An operational objective analysis system. *Mon. Weather Rev.*, **87 (10)**,
1042 367–374.
- 1043 De Pondeca, M. S., and Coauthors, 2011: The real-time mesoscale analysis at noaa’s national
1044 centers for environmental prediction: current status and development. *Weather and Forecasting*,
1045 **26 (5)**, 593–612.
- 1046 Dörenkämper, M., and Coauthors, 2020: The Making of the New European Wind Atlas –
1047 Part 2: Production and evaluation. *Geoscientific Model Development*, **13 (10)**, 5079–5102,
1048 <https://doi.org/10.5194/gmd-13-5079-2020>.

- 1049 Dowell, D. C., and Coauthors, 2022: The high-resolution rapid refresh (hrrr): An hourly updating
1050 convection-allowing forecast model. part i: Motivation and system description. *Weather and*
1051 *Forecasting*, **37 (8)**, 1371–1395.
- 1052 Dujardin, J., and M. Lehning, 2022: Wind-Topo: Downscaling near-surface wind fields to high-
1053 resolution topography in highly complex terrain with deep learning. *Q. J. R. Meteorol. Soc.*,
1054 **148 (744)**, 1368–1388.
- 1055 El Rafei, M., S. Sherwood, J. Evans, and A. Dowdy, 2023: Analysis and characterisation of extreme
1056 wind gust hazards in new south wales, australia. *Nat. Hazards (Dordr.)*, **117 (1)**, 875–895.
- 1057 Franco, B. M., L. Hernández-Callejo, and L. M. Navas-Gracia, 2020: Virtual weather stations for
1058 meteorological data estimations. *Neural Comput. Appl.*, **32 (16)**, 12 801–12 812.
- 1059 Fukami, K., R. Maulik, N. Ramachandra, K. Fukagata, and K. Taira, 2021: Global field recon-
1060 struction from sparse sensors with voronoi tessellation-assisted deep learning. *Nat. Mach. Intell.*,
1061 **3 (11)**, 945–951.
- 1062 Gandin, L. S., 1963: Objective analysis of meteorological fields. *Israel program for scientific*
1063 *translations*, **242**.
- 1064 Gilleland, E., D. Ahijevych, B. G. Brown, B. Casati, and E. E. Ebert, 2009: Intercomparison of
1065 spatial forecast verification methods. *Weather and forecasting*, **24 (5)**, 1416–1430.
- 1066 Glorot, X., and Y. Bengio, 2010: Understanding the difficulty of training deep feedforward neural
1067 networks. *Proceedings of the thirteenth international conference on artificial intelligence and*
1068 *statistics*, JMLR Workshop and Conference Proceedings, 249–256.
- 1069 Gultepe, I., and Coauthors, 2019: A review of high impact weather for aviation meteorology. *Pure*
1070 *and applied geophysics*, **176 (5)**, 1869–1921.
- 1071 Güemes, A., C. Sanmiguel Vila, and S. Discetti, 2022: Super-resolution generative adversarial
1072 networks of randomly-seeded fields. *Nat. Mach. Intell.*, **4 (12)**, 1165–1173.
- 1073 Hardy, R. L., 1971: Multiquadric equations of topography and other irregular surfaces. *Journal of*
1074 *geophysical research*, **76 (8)**, 1905–1915.

- 1075 Haylock, M. R., N. Hofstra, A. M. G. Klein Tank, E. J. Klok, P. D. Jones, and M. New, 2008:
1076 A european daily high-resolution gridded data set of surface temperature and precipitation for
1077 1950–2006. *J. Geophys. Res.*, **113** (D20).
- 1078 He, H., and E. A. Garcia, 2009: Learning from Imbalanced Data. *IEEE Transactions on Knowledge
1079 and Data Engineering*, **21** (9), 1263–1284, <https://doi.org/10.1109/TKDE.2008.239>.
- 1080 He, K., X. Zhang, S. Ren, and J. Sun, 2015: Deep residual learning for image recognition. *arXiv
1081 [cs.CV]*.
- 1082 Hersbach, H., and Coauthors, 2020: The era5 global reanalysis. *Quarterly journal of the royal
1083 meteorological society*, **146** (730), 1999–2049.
- 1084 Houston, A. L., N. A. Lock, J. Lahowetz, B. L. Barjenbruch, G. Limpert, and C. Oppermann,
1085 2015: Thunderstorm observation by radar (thor): An algorithm to develop a climatology of
1086 thunderstorms. *Journal of Atmospheric and Oceanic Technology*, **32** (5), 961–981.
- 1087 Hu, Y., G. Jia, H. Gao, Y. Li, M. Hou, J. Li, and C. Miao, 2023: Spatial characterization of global
1088 heat waves using satellite-based land surface temperature. *International Journal of Applied Earth
1089 Observation and Geoinformation*, **125**, 103 604.
- 1090 Jahan, I., D. Cerrai, and M. Astitha, 2024: Storm gust prediction with the integration of ma-
1091 chine learning algorithms and WRF model variables for the northeast united states. *Artificial
1092 Intelligence for the Earth Systems*, **3** (3).
- 1093 Jeong, Y., D. Kim, and K. Byun, 2025: A novel deep learning-based approach for reconstruc-
1094 tion of historical long-term high-quality gridded meteorological dataset. *J. Hydrol. (Amst.)*,
1095 **654** (132850), 132 850.
- 1096 Kahl, J. D. W., 2020: Forecasting peak wind gusts using meteorologically stratified gust factors
1097 and MOS guidance. *Weather Forecast.*, **35** (3), 1129–1143.
- 1098 Kalnay, E., 2003: *Atmospheric modeling, data assimilation and predictability*. Cambridge univer-
1099 sity press.

- 1100 Kartal, S., S. Basu, and S. J. Watson, 2023: A decision-tree-based measure–correlate–predict
1101 approach for peak wind gust estimation from a global reanalysis dataset. *Wind Energy Sci.*,
1102 **8 (10)**, 1533–1551.
- 1103 Knapp, K. R., and J. P. Kossin, 2007: New global tropical cyclone data set from isccp b1
1104 geostationary satellite observations. *Journal of Applied Remote Sensing*, **1 (1)**, 013 505.
- 1105 Koch, S. E., M. desJardins, and P. J. Kocin, 1983: An interactive Barnes objective map analysis
1106 scheme for use with satellite and conventional data. *J. Clim. Appl. Meteorol.*, **22 (9)**, 1487–1503.
- 1107 Kwon, D.-K., and A. Kareem, 2009: Gust-front factor: New framework for wind load effects on
1108 structures. *Journal of structural engineering*, **135 (6)**, 717–732.
- 1109 Le Toumelin, L., I. Gouttevin, N. Helbig, C. Galiez, M. Roux, and F. Karbou, 2023: Emulating the
1110 adaptation of wind fields to complex terrain with deep learning. *Artif. Intell. Earth Syst.*, **2 (1)**,
1111 1–39.
- 1112 Liston, G. E., and K. Elder, 2006: A meteorological distribution system for high-resolution
1113 terrestrial modeling (MicroMet). *J. Hydrometeorol.*, **7 (2)**, 217–234.
- 1114 Mahmood, R., and Coauthors, 2017: Mesonets: Mesoscale weather and climate observations for
1115 the united states. *Bull. Am. Meteorol. Soc.*, **98 (7)**, 1349–1361.
- 1116 Mankin, K. R., S. Mehan, T. R. Green, and D. M. Barnard, 2025: Review of gridded climate
1117 products and their use in hydrological analyses reveals overlaps, gaps, and the need for a more
1118 objective approach to selecting model forcing datasets. *Hydrol. Earth Syst. Sci.*, **29 (1)**, 85–108.
- 1119 Martinez-Roig, M., 2025: Comparative Evaluation of Deep Learning Architectures for Spatial
1120 Infilling of Meteorological Data. Geosphere Austria.
- 1121 Mass, C. F., and D. Ovens, 2021: The synoptic and mesoscale evolution accompanying the 2018
1122 camp fire of northern california. *Bull. Am. Meteorol. Soc.*, **102 (1)**, E168–E192.
- 1123 Matheron, G., 1967: Kriging or polynomial interpolation procedures. *CIMM Transactions*, **70 (1)**,
1124 240–244.
- 1125 Mitchell, J. W., 2013: Power line failures and catastrophic wildfires under extreme weather
1126 conditions. *Eng. Fail. Anal.*, **35**, 726–735.

- 1127 Müller, S. H., D. Cáceres, S. Eisner, M. Flörke, C. Herbert, C. Niemann, and Döll, 2021: The
1128 global water resources and use model WaterGAP v2. 2d: Model description and evaluation.
1129 *Geoscientific Model Development*, **14** (2), 1037–1079.
- 1130 National Weather Service Philadelphia/Mt. Holly, 2023: Event review: April 1, 2023 tornado
1131 outbreak. Accessed: 2025-07-28, <https://www.weather.gov/phi/eventreview2023401>.
- 1132 New York Post, 2023a: New york city facing threat of 60 mph
1133 winds, tornadoes. Accessed: 2025-07-28, [https://nypost.com/2023/04/01/
1134 new-york-city-facing-threat-of-60-mph-winds-tornadoes/](https://nypost.com/2023/04/01/new-york-city-facing-threat-of-60-mph-winds-tornadoes/).
- 1135 New York Post, 2023b: Northeast hit with tornadoes, damaging wind gusts. Accessed: 2025-07-28,
1136 <https://nypost.com/2023/04/02/northeast-hit-with-tornadoes-damaging-wind-gusts/>.
- 1137 Niziol, T., 1987: Operational forecasting of lake effect snowfall in western and central new york.
1138 *Weather and Forecasting*, **2**, 310–321.
- 1139 Niziol, T. A., W. R. Snyder, and J. S. Waldstreicher, 1995: Winter weather forecasting throughout
1140 the eastern united states. part IV: Lake effect snow. *Weather Forecast.*, **10** (1), 61–77.
- 1141 Peng, J., S. Dadson, F. Hirpa, E. Dyer, T. Lees, D. G. Miralles, and Funk, 2020: A pan-african
1142 high-resolution drought index dataset. *Earth System Science Data*, **12** (1), 753–769.
- 1143 Peraza, J., P. R. Rossini, and A. Patrignani, 2025: Mapping mesoscale soil moisture using a
1144 model-data fusion approach. *J. Hydrol. (Amst.)*, **654** (132768), 132768.
- 1145 Rasmussen, R., and Coauthors, 2023: Conus404: The near-usgs 4-km long-term regional hy-
1146 droclimate reanalysis over the conus. *Bulletin of the American Meteorological Society*, **104** (8),
1147 E1382–E1408.
- 1148 Ridal, M., and Coauthors, 2024: Cerra, the copernicus european regional reanalysis system.
1149 *Quarterly Journal of the Royal Meteorological Society*, **150** (763), 3385–3411.
- 1150 Sallis, P. J., W. Claster, and S. Hernández, 2011: A machine-learning algorithm for wind gust
1151 prediction. *Comput. Geosci.*, **37** (9), 1337–1344.
- 1152 Shepard, D., 1968: A two-dimensional interpolation function for irregularly-spaced data. *Proceed-
1153 ings of the 1968 23rd ACM national conference*, 517–524.

- 1154 Sheridan, P., 2018: Current gust forecasting techniques, developments and challenges. *Adv. Sci.*
1155 *Res.*, **15**, 159–172.
- 1156 Skamarock, W. C., 2004: Evaluating mesoscale nwp models using kinetic energy spectra. *Monthly*
1157 *weather review*, **132 (12)**, 3019–3032.
- 1158 Smith, T. M., and Coauthors, 2016: Multi-radar multi-sensor (mrms) severe weather and aviation
1159 products: Initial operating capabilities. *Bulletin of the American Meteorological Society*, **97 (9)**,
1160 1617–1630.
- 1161 Subedi, S., A. Kechchour, M. Kantar, V. Sharma, and B. C. Runck, 2025: Can gridded real-time
1162 weather data match direct ground observations for irrigation decision-support? *Agrosyst. Geosci.*
1163 *Environ.*, **8 (2)**.
- 1164 Sun, J., and N. A. Crook, 2001: Real-time low-level wind and temperature analysis using single
1165 WSR-88D data. *Weather Forecast.*, **16 (1)**, 117–132.
- 1166 Sunderhaft, R., L. Frank, and J. Davis, 2024: Deep learning improvements for sparse spatial field
1167 reconstruction. *arXiv [cs.CV]*.
- 1168 Themeßl, M. J., A. Gobiet, and A. Leuprecht, 2011: Empirical-statistical downscaling and er-
1169 ror correction of daily precipitation from regional climate models. *International Journal of*
1170 *Climatology*, **31 (10)**, 1530–1544.
- 1171 Venter, Z. S., A. Hassani, E. Stange, P. Schneider, and N. Castell, 2024: Reassessing the role of
1172 urban green space in air pollution control. *Proc. Natl. Acad. Sci. U. S. A.*, **121 (6)**, e2306200 121.
- 1173 Wang, H., H. Zhou, and S. Cheng, 2024: Dynamical system prediction from sparse observations
1174 using deep neural networks with voronoi tessellation and physics constraint. *Comput. Methods*
1175 *Appl. Mech. Eng.*, **432 (117339)**, 117 339.
- 1176 Wang, Z., X. Cun, J. Bao, W. Zhou, J. Liu, and H. Li, 2022: Uformer: A general U-shaped
1177 transformer for image restoration. *2022 IEEE/CVF Conference on Computer Vision and Pattern*
1178 *Recognition (CVPR)*, IEEE.

- 1179 Wasula, A. C., L. Bosart, and K. D. LaPenta, 2002: The influence of terrain on the severe weather
1180 distribution across interior eastern new york and western new england. *Weather and Forecasting*,
1181 **17**, 1277–1289.
- 1182 World Meteorological Organization (WMO), 2023: *Guide to Instruments and Methods of Obser-*
1183 *vation*. World Meteorological Organization, Geneva, Switzerland, URL [https://library.wmo.int/](https://library.wmo.int/idurl/4/41650)
1184 [idurl/4/41650](https://library.wmo.int/idurl/4/41650), wMO-No. 8.
- 1185 Xiao, H., Y. Wang, Y. Zheng, Y. Zheng, X. Zhuang, H. Wang, and M. Gao, 2023: Convective-
1186 gust nowcasting based on radar reflectivity and a deep learning algorithm. *Geosci. Model Dev.*,
1187 **16 (12)**, 3611–3628.
- 1188 Zanetta, F., D. Nerini, M. Buzzi, and H. Moss, 2025: Efficient modeling of sub-kilometer surface
1189 wind with gaussian processes and neural networks. *Artif. Intell. Earth Syst.*
- 1190 Zhao, X., X. Chen, Z. Gong, W. Zhou, W. Yao, and Y. Zhang, 2024: RecFNO: A resolution-
1191 invariant flow and heat field reconstruction method from sparse observations via fourier neural
1192 operator. *Int. J. Therm. Sci.*, **195 (108619)**, 108 619.