



ViTAE-SL: A vision transformer-based autoencoder and spatial interpolation learner for field reconstruction

Hongwei Fan, Sib0 Cheng, Audrey J de Nazelle, Rossella Arcucci

► To cite this version:

Hongwei Fan, Sib0 Cheng, Audrey J de Nazelle, Rossella Arcucci. ViTAE-SL: A vision transformer-based autoencoder and spatial interpolation learner for field reconstruction. *Computer Physics Communications*, 2025, 308, pp.109464. <10.1016/j.cpc.2024.109464>. <hal-04876671>

HAL Id: hal-04876671

<https://hal.science/hal-04876671v1>

Submitted on 9 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

1 ViTAE-SL: a vision transformer-based autoencoder and 2 spatial interpolation learner for field reconstruction

3 Hongwei Fan^{a,b}, Sibor Cheng^d, Audrey J de Nazelle^b, Rossella Arcucci^{c,a}

^a*Data Science Institute, Imperial College London, London, UK*

^b*Centre for Environmental Policy, Imperial College London, London, UK*

^c*Department of Earth Science and Engineering, Imperial College London, London, UK*

^d*CEREA, Ecole des Ponts and EDF R&D, Ile-de-France, France*

4 **Abstract**

Reliable and accurate reconstruction for large-scale and complex physical fields in real-time from limited observations has been a longstanding challenge. In recent years, sensors have been increasingly deployed in numerous physical systems. However, the locations of these sensors can shift over time, such as with mobile sensors, or when sensors are deployed and removed. These sparse and randomly located sensors further exacerbate the difficulty of reconstructing the physical field. In this paper, we present a new deep learning model called Vision Transformer-based Autoencoder (ViTAE) for reconstructing large-scale and complex fields. The proposed network structure is based on a novel core design: vision transformer encoder and Convolutional Neural Network (CNN) decoder. First, we split a two-dimensional field into patches and developed a vision transformer encoder to transfer patches into latent representations. We then reshape the linear latent representations to patches before concatenation, along with a CNN decoder, to reconstruct the field. The proposed model is tested in four different numerical experiments, using generated synthetic data, spatially distributed PM2.5 data, Computational Fluid Dynamics (CFD) simulation data and National Oceanic and Atmospheric Administration (NOAA) sea surface temperature data. The numerical results highlight the strength of ViTAE-SL compared to Kriging and state-of-the-art deep-learning models with significantly higher reconstruction

accuracy, computational efficiency, and robust scaling behavior.

5 *Keywords:* spatial interpolation, field reconstruction, vision transformer,
6 autoencoder

7 **1. Introduction**

8 Advances in sensor technologies have enabled sensor deployments in physics
9 system at a large scale. However, the number of monitoring sites within
10 study areas is often relatively limited due to the high maintenance and setup
11 costs. Sparsely scattered sensors are unable to capture whole physical fields
12 in detail [1], and time-varying sensors further increase the difficulty of re-
13 constructing the physical field; therefore robust and efficient spatial inter-
14 polation models are needed to reconstruct physical fields from limited and
15 randomly localised sensor information. Spatial interpolation, which is pre-
16 dicting values of a spatial process in unmonitored areas from local sensor
17 observations, is a major challenge in physics systems. Traditional applica-
18 tions of spatial interpolation include physics systems, environmental science,
19 geophysics, astrophysics, atmospheric science, and fluid dynamics, and have
20 been extended to other fields such as computer vision, public health, and
21 biological sciences [2, 3, 4].

22 Kriging [5, 6] is a spatial interpolation method that offers linear unbiased
23 prediction based on observations. Functioning as a Gaussian process [7] gov-
24 erned by covariance, Kriging extrapolates values in unobserved areas through
25 a weighted average of observations. The computation of these weights neces-
26 sitates the estimation of spatial covariance functions, conventionally assumed
27 to be stationary within the framework of Kriging. However, real physical
28 fields often have non-Gaussian and non-stationary spatial covariance func-
29 tions [8, 9]. For example, the spatial covariance of PM2.5 concentrations
30 usually changes a lot from one place to another, such as between cities and
31 rural areas[10]. Therefore, Kriging may not always be optimal in practical
32 real-world scenarios. In addition, kriging has limitations in handling large

33 spatial datasets due to its computational intensity. Finding the inversion of
 34 a $N \times N$ covariance matrix is needed to use Kriging. The number of ob-
 35 servations, N , determines the size of the matrix, and the computation takes
 36 $O(N^3)$ time and $O(N^2)$ memory using the standard Cholesky decomposition
 37 method. Due to these limitations, it is challenging to utilize Kriging for
 38 real-time large-field reconstruction.

39 In recent years, deep learning (DL) [11] and neural network (NN) [12]
 40 have been extensively used in a wide spectrum of applications [13, 11]. DL
 41 has witnessed an explosion of architectures that are continuously evolving in
 42 both size and complexity, improving their capacity and capabilities in many
 43 different tasks [14, 15, 16, 17, 18, 19, 20, 21]. NNs are effective for predictions
 44 with complex features such as non-linearity and non-stationary. The use of
 45 GPUs has made NNs computationally efficient for analyzing large datasets.

46 Deep neural networks [11] have emerged as hopeful methods to recon-
 47 struct physical field from sparse measurements in an efficient manner [22,
 48 23, 24, 19, 25]. Most neural networks require structure data, such as grid
 49 data. However, practical measurements are typically uneven or unstructured,
 50 such as air quality monitoring [26], hydrology measurement [27], and oceano-
 51 graphic observations [23]. Although graph neural network (GCN) [28] and
 52 multi-layer perception (MLP) [11] have the ability to process unstructured
 53 data, their scalability is limited due to the high computational cost. Further-
 54 more, these methods require fixed numbers and positions of measurements as
 55 input data and cannot accommodate time-varying sensors, rendering them
 56 impractical in real-world physical field reconstruction, where the numbers
 57 and positions of sensors typically change over time [29].

58 To address these two bottlenecks, Fukami et al. [30] employed Voronoi
 59 tessellation to convert observations into a structured grid format that is com-
 60 patible with convolutional neural network (CNN) and then utilized CNN to
 61 reconstruct the field from the structured grid. However, convolutions typi-
 62 cally operate on a small patch of field to extract features but do not encode

the relative position of different features. This makes it difficult for CNN to explore spatial dependencies for large field reconstruction [31], especially when the observations are very sparse. For example, air pollution monitoring sites are often sparsely distributed in large areas, making the reconstruction of air pollution maps on a fine scale a longstanding challenge [32].

The introduction of Vision Transformers (ViT) [33] has the potential to address this challenge. Transformers, originally proposed by Vaswani et al. (2017) [14], have emerged as a prominent technique in the field of Natural Language Processing (NLP) due to their ability to capture long-range dependencies and learn contextual relationships effectively. Besides the advanced structure of transformers, the training strategy is also crucial for the success of this NLP models: They involve removing a portion of the text and learning to predict the removed content. This strategy, known as auto-encoding, is applied in various fields, such as computer vision and flow problems, and demonstrates encouraging outcomes [34, 35]. For example, He et al. [36] showed that by employing masked autoencoders (AEs) to remove random sections of the input image and then rebuild those missing sections, it is possible to rebuild images that appear realistic, even when more than 90% of the image is masked. Rebuilding images from random visible patches is conceptually similar to field reconstruction from observations.

While the ViT model and the autoencoder (AE) method [36] have achieved success in image reconstruction, as far as we know, there is no prior research on field reconstruction by ViT. There are two problems that hinder the application of these methods in the field reconstruction task:

- ViT typically predicts masked image patches using visible patches. Compared to image patches with values on each pixel, the observations in this study are sparse in each patch. Setting the patch size of 1×1 pixel in ViT will lead to high computational cost, which makes the model applicable only to small resolution fields.
- ViT predicts each patch separately. Concatenating the predicted patches

93 can result in heterogeneous and inconsistent reconstructions [36].

94 This paper proposes using ViT and AE for large physical field reconstruc-
95 tion from sparse and time-varying observations.

96 **2. Contribution of the present work**

97 Inspired by the success of ViT and AE methods, we introduce a new deep
98 learning model called Vision Transformer-based autoencoder and spatial in-
99 terpolation learner (ViTAE-SL) for reconstructing large-scale and complex
100 fields. More precisely, the proposed framework consists of a ViT as the en-
101 coder and a CNN as the decoder. In order to tackle the challenges mentioned
102 earlier, we propose a method that involves incorporating sparse sensor data
103 into a Transformer. This is achieved by mapping the observations in the grid
104 field and masking the unmonitored grids. The next step involves splitting
105 the observations field into patches and feeding them into the transformer
106 encoder to obtain features. These features are then reshaped into patches
107 and concatenated together. To bridge the gap between patches and make
108 the reconstruction field more consistent, we use the CNN decoder to pre-
109 dict grid values based on concatenated patches. Following this structure,
110 our ViTAE-SL is capable of fast and accurate field construction using time-
111 varying and unstructured observations. Four test cases, including synthetic
112 data and real-world applications, are implemented in this study to compare
113 the performance of the proposed ViTAE-SL against state-of-the-art field re-
114 construction methods.

115 The subsequent sections of the paper are structured in the following
116 manner: Section 3 introduces the construction and properties of ViTAE-SL
117 model. Section 4 showcases four case studies that demonstrate the efficacy
118 of ViTAE-SL. Section 5 summarizes our primary findings and proposes po-
119 tential areas for future research.

3. Methodology

Our objective is to reconstruct a two-dimensional global field variable $Q \in \mathbb{R}^{n_x \times n_y}$ from a vector of local sensor measurements $s \in \mathbb{R}^n$, locations $x_{si} \in \mathbb{R}^2, i = 1, \dots, n$. Here, n_x and n_y respectively denote the number of grid points in the horizontal and vertical directions on a high-resolution field, and n indicates the number of local sensor measurements. ViTAE-SL is an autoencoding approach that reconstructs the original field-given observations. To handle the sparse observations and learn efficient representations from them, we adopt a ViT-based encoder that allows us to operate on the observations. After that, a lightweight CNN decoder is used to reconstruct the full grid field from the latent representations. Figure 1 illustrates the flowchart of the proposed approach. In what follows, we introduce the ViT-based encoder and CNN decoder, which are the two key components in the present approach.

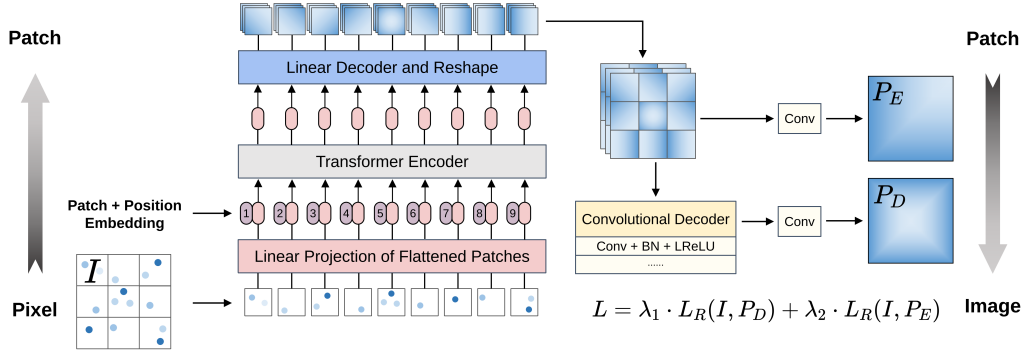


Figure 1: Model overview of ViTAE-SL. We split the grid field into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. After that, we used a CNN decoder to predict the grid values from the encoder’s latent.

134 *ViT-based encoder.* First, we transform observations into grid fields based
 135 on their location to allow use of ViT, defined as:

$$I(\{\mathbf{x}_{si}\}_{i=1}^n) = \begin{cases} x_i & \text{for any location } \mathbf{x}_{si} \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

136 n is the number of observations, where x_i denotes the value of the i -th
 137 observation. We start by splitting the grid field $I \in \mathbb{R}^{n_x \times n_y}$ into patches
 138 $\mathbf{x}_p \in \mathbb{R}^{N \times (P_h \times P_w)}$, where $P_h \times P_w$ is the resolution of each patch. N is the
 139 corresponding number of split patches and is also the input sequence length
 140 of ViT. Using a patch embedding projection \mathbf{E} , the patches are mapped to
 141 the embedding of dimensions $D = P_h \times P_w$. Then, positional embeddings
 142 \mathbf{E}_{pos} are added and fed into a series of Transformer blocks to obtain latent
 143 representations.

144 The processing of encoder can mathematically be formulated as:

$$\begin{aligned} \mathbf{z}_0 &= [\mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos}, \quad \mathbf{E} \in \mathbb{R}^{(P_h \cdot P_w) \times D}, \mathbf{E}_{pos} \in \mathbb{R}^{N \times D} \\ \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \quad \ell = 1 \dots L \\ \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell, \quad \ell = 1 \dots L \\ \mathbf{y} &= \text{LN}(\mathbf{z}_L^0) \end{aligned} \quad (2)$$

145 The Transformer encoder consists of L layers of multiheaded self-attention
 146 (MSA) and multilayer perceptron (MLP) blocks. Layernorm (LN) is applied
 147 before every block, and residual connections are applied after every block.
 148 \mathbf{y} is the output of the encoder. The length of \mathbf{y} is equal to the number of
 149 grids in a patch. We reshape the vectors into patches and concatenate these
 150 patches to form a field-shaped representation, which is used as the input to a
 151 CNN-based decoder. Additionally, a convolution layer is applied to produce
 152 $P_E \in \mathbb{R}^{n_x \times n_y}$ based on the field-shaped representation.

153 *CNN-based decoder.* The CNN-based decoder consists of sequences of convo-
 154 lutional layers, batch normalization, and activation layers. Batch normaliza-
 155 tion is applied after each convolutional layer to normalize the output before
 156 passing it through the activation layers. The parameters of the CNN-based
 157 decoder are shown in 3. The CNN-based decoder receives the field-shaped
 158 representation produced by the encoder as input, and the output of the de-
 159 coder is the predicted field $P_D \in \mathbb{R}^{n_x \times n_y}$ that corresponds to the entire field
 160 \mathbf{Q} .

$$\begin{aligned} P'_\ell &= \text{Conv2D}(P_{\ell-1}, \mathbf{W}_\ell), \quad \ell = 1 \dots L \\ P_\ell &= \text{ReLU}(P'_\ell), \quad \ell = 1 \dots L \end{aligned} \quad (3)$$

161 where \mathbf{W} denotes weights (filters) of CNN. The output of each filter
 162 operation is passed through an activation function which is chosen to be the
 163 rectified linear unit (ReLU).

164 *Reconstruction target.* Our ViTAE-SL reconstructs the field by predicting
 165 the values for each grid. The loss function computes the mean squared error
 166 (MSE) between the reconstructed P_E , P_D and original fields \mathbf{Q} in pixel space.
 167 L_R denotes the mean squared error. The loss is defined as:

$$L = \lambda_1 \cdot L_R(Q, P_D) + \lambda_2 \cdot L_R(Q, P_E) \quad (4)$$

168 Where λ_1 and λ_2 are the weights for balancing the multiple objectives.

169 4. Test cases and results

170 In this section we thoroughly demonstrate the performance of present
 171 ViTAE-SL for field reconstruction with extensive numerical experiments un-
 172 der different circumstances. Our examples include:

- 173 • Simulation data produced by Gaussian covariance models
- 174 • China air quality dataset

175 • Fluid dynamics data of unsteady wake flow

176 • NOAA global sea surface temperature

177 These datasets have different characteristics in terms of structured or un-
 178 structured field, time varying data and observation noise as shown in Table 1.
 179 The performance of the field reconstruction methods: ViTAE-SL, Kriging
 180 and Voronoi tessellation-assisted convolutional neural network (VCNN) have
 181 been evaluated on both structured and unstructured data. For example,
 182 the unsteady wake flow is based on two-dimensional irregular meshes. How-
 183 ever, all observations used in this study are unstructured. The first test
 184 case consists of simulation data generated by Gaussian covariance model and
 185 the observable points are randomly selected and time-varying. Spatially dis-
 186 tributed air quality observations come from the Chinese regulatory monitor
 187 stations [37] which change over time. Unsteady wake flow has fixed but very
 188 sparse observations. As for the National Oceanic and Atmospheric Admin-
 189 istration (NOAA) sea surface temperature data, we select different numbers
 190 of observations as the train and test dataset, which are time-varying. These
 191 datasets are chosen to evaluate the performance of ViTAE-SL for handling
 192 sparse, unstructured and time-varying data.

Table 1: Different test cases performed in this paper

| Test case | unstructured field | unstructured observation | time-varying sensors | observation noise |
|--------------------|-----------------------|-----------------------------|-------------------------|----------------------|
| simulation data | X | ✓ | ✓ | X |
| Air quality | ✓ | ✓ | ✓ | ✓ |
| Unsteady wake flow | ✓ | ✓ | X | X |
| NOAA | ✓ | ✓ | ✓ | X |

193 *4.1. Test case 1: Simulation data*

194 In this section, we evaluate the performance of ViTAE-SL against Kriging
195 and VCNN in stationary simulation data.

196 In order to replicate a physically uniform field, the simulation field data
197 are created using the *gstool* [38], with a Gaussian covariance kernel with a
198 correlation scale length of $L \in (40, 80)$. This means that the correlation
199 between two points in the field only depends on their spatial distance, which
200 is ideal for the Kriging method. These simulations are commonly used to
201 compare different field reconstruction methods [39].

202 After creating the grid field, grid points are randomly selected as obser-
203 vations. Unlike VCNN and ViTAE-SL, the Kriging method requires prior
204 knowledge of the covariance kernel. Numerical experiments on Kriging are
205 conducted using two kernel functions: Gaussian and Exponential, both with
206 the exact correlation length used for data generation. The latter is used to
207 simulate situations where the kernel function is misjudged, as in real-world
208 scenarios, the exact covariance kernel is often difficult to determine [40].

209 Figure 2 illustrates the computational time of Kriging (using the Gaus-
210 sian kernel) for reconstructing various field sizes based on different numbers
211 of observations. It’s observed that as the field size and the number of ob-
212 servations increase, the computational time for Kriging grows exponentially.
213 This poses computational challenges for large-scale field reconstruction with
214 Kriging. For instance, when the field size exceeds 256 and the number of
215 observations is more than 0.1%, Kriging requires thousands of seconds to fit
216 and predict.

217 In order to compare the performance of ViTAE-SL against VCNN and
218 Kriging, we conducted experiments on a field size of 512 x 512. The number
219 of observations used for training was set to 0.5%, 1%, 2%, and 5% of the
220 total number of grid points in the field. For each observation ratio, we
221 generated 10,000 field snapshots and randomly selected observations from
222 each snapshot, resulting in time-varying observation placement. We used an

223 80/10/10 split for training, validation, and testing datasets.

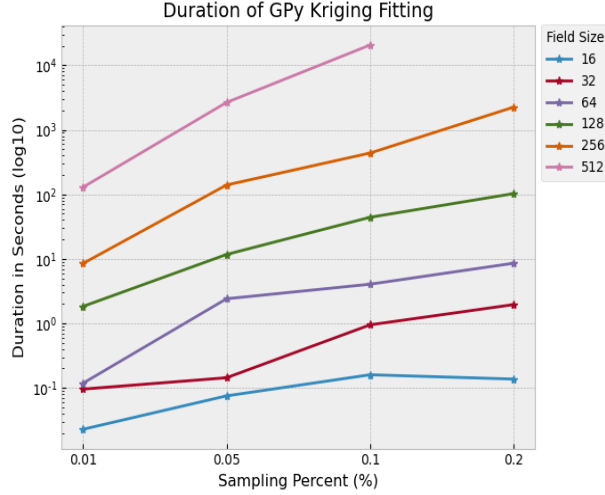


Figure 2: Execution time of Kriging for different size of the field.

224 We use the original ViT setup [33] for the ViT-based encoder design and
 225 employ “Lite”, “Base”, and “Large” models as summarized in Table 2. In the
 226 following sections, we’ll use abbreviated notation to denote the model size and
 227 the input patch size. For instance, ViTAE-large/16 represents combination
 228 of the “Large” version of ViT in Table 2 and CNN in Table 3 with 16×16
 229 input patch size.

| Model | Layers | Hidden Size D | Heads | Channel | Patch Size |
|-----------|--------|---------------|-------|---------|------------|
| ViT-Lite | 8 | 32 | 8 | 16 | 16 |
| ViT-Base | 8 | 64 | 8 | 32 | 16 |
| ViT-Large | 8 | 128 | 8 | 64 | 16 |

Table 2: Details of Vision Transformer encoder variants.

230 For the CNN-based decoder, we use a sample convolutional neural net-
 231 work; the configurations of each decoder for “Lite”, “Base”, and “Large”
 232 models are shown in Table 3. The CNN-based decoder is trained to estimate

233 the whole simulation field based on intermediate features from ViT-based
 234 encoder.

| Model | Layers | Channel D | filter size |
|-----------|--------|-----------|-------------|
| CNN-Lite | 5 | 16 | 3 |
| CNN-Base | 5 | 32 | 3 |
| CNN-Large | 5 | 64 | 3 |

Table 3: Details of CNN decoder variants.

234
 235 As shown in Figure 3, the field reconstruction from the ViTAE-SL closely
 236 resembles the ground truth (GT) without prior knowledge.

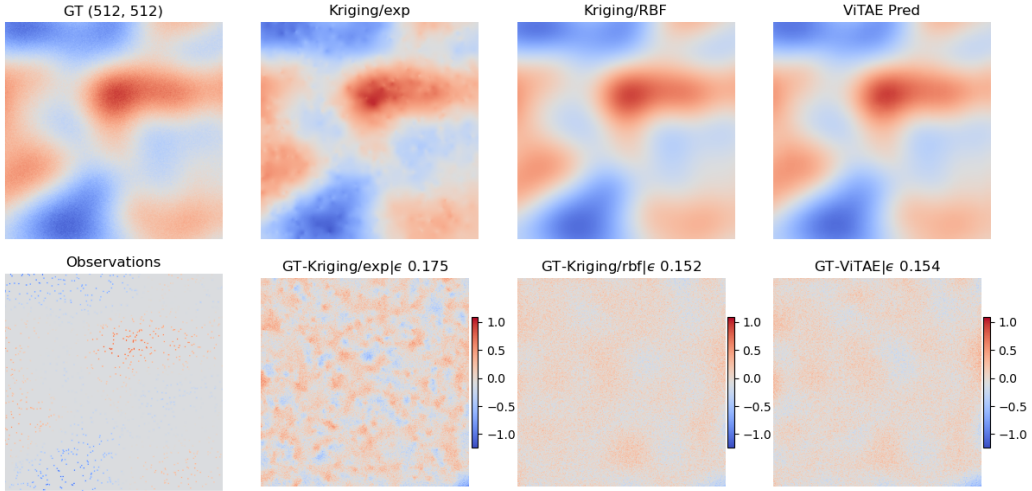


Figure 3: 512×512 Gaussian field reconstruction results of ViTAE-SL and Kriging , 0.5% sampling rate compare to the GT.

237 Figure 3 also reports the relative error defined as:

$$\epsilon = \frac{\|Q_{\text{ref}} - Q_{\text{reconstruct}}\|_2}{\|Q_{\text{ref}}\|_2}, \quad (5)$$

238 where $\|\cdot\|$ denotes the L_2 norm, and Q_{ref} and $Q_{\text{reconstruct}}$ are the reference
 239 and reconstructed fields, respectively. The relative error metric is commonly

used in field reconstruction and prediction tasks [41, 42]. In our study, we compared Kriging’s reconstruction results with Gaussian and exponential covariance kernels, referred to as Kriging/RBF and Kriging/exp, respectively. Figure 3 clearly shows that the Kriging/RBF model greatly outperforms the Kriging/exp model, suggesting that Kriging is vulnerable when the covariance kernel is not accurately specified. However, ViTAE-SL created reconstruction results that were nearly as precise as those of Kriging/RBF without any prior information. Besides, ViTAE-SL is much more efficient than Kriging, as shown in Table 5. For example, ViTAE-SL-lite/16 runs 10^6 faster than Kriging when the field size is 512 and the ratio of observation grids is 5%.

| Model | ϵ | | | |
|------------------|---------------|---------------|---------------|---------------|
| Kriging/RBF | 0.2243 | 0.2221 | 0.2218 | 0.2215 |
| Kriging/Exp | 0.2553 | 0.2552 | 0.2550 | 0.2379 |
| VCNN | 0.2324 | 0.2259 | 0.2216 | 0.2160 |
| ViTAE-lite/16 | 0.2431 | 0.2346 | 0.2290 | 0.2242 |
| ViTAE-base/16 | 0.2280 | 0.2369 | 0.2250 | 0.2234 |
| ViTAE-large/16 | 0.2255 | 0.2228 | 0.2213 | 0.2202 |
| Sampling Percent | 0.5% | 1% | 2% | 5% |

Table 4: Gaussian field reconstruction result of ViTAE-SL, VCNN and Kriging.

| Model | Execution time (s) | | | |
|------------------|--------------------|--------|--------|--------|
| Kriging/RBF | 21 | 59 | 191 | 1491 |
| Kriging/Exp | 31 | 76 | 253 | 1586 |
| VCNN | 0.035 | 0.035 | 0.035 | 0.035 |
| ViTAE-lite/16 | 0.0105 | 0.0104 | 0.0105 | 0.0106 |
| ViTAE-base/16 | 0.0128 | 0.0127 | 0.0128 | 0.0128 |
| ViTAE-large/16 | 0.0150 | 0.0154 | 0.0151 | 0.0153 |
| Sampling Percent | 0.5% | 1% | 2% | 5% |

Table 5: Execution time in seconds of the Gaussian field reconstruction for ViTAE-SL, VCNN and Kriging.

We also compare ViTAE-SL with the VCNN [30]. Table 4 shows the online computational time and the ϵ values for ViTAE-SL, VCNN and Krig-

ing for different numbers of observations in 512×512 field reconstruction. These results show that ViTAE-SL outperforms VCNN when observations are sparse (less than 2% in this case), demonstrating the strength of ViTAE-SL in handling sparse observations.

4.2. Test case 2: Chinese air quality

To evaluate the performance of ViTAE-SL in real-world physical field reconstruction problems, we first applied our approach to the PM2.5 concentration dataset in China. The ChinaHighPMC [37] dataset consists of daily air pollutants at ground level PM2.5 for China with 1km resolution. These data are extracted and combined from different resources, including ground-based measurements, satellite remote sensing products, atmospheric reanalysis, and model simulations. In particular, we utilise the 5km, 10km and 20km resolution grid fields for the Chinese area, so the corresponding shapes of the entire field are 714×1229 , 357×615 and 184×312 . Unstructured observations from monitor stations [43] are projected in the grid field. Since sensors are affected by ambient factors, noise is included in these observation values. When there are several monitor stations in one grid cell, the value is calculated as the average observation values. The daily air quality field and corresponding observations span from 2013 to 2020 and the total number of snapshots is 2992. We randomly selected 80% of the snapshots as training data sets, 10% are used as validation dataset and 10 % for testing dataset. The average number of observations is 1173, resulting in 0.13%, 0.52% and 2.08% against the total number of grid points in 5km, 10km and 20km resolution grid fields. Table 6 shows the model parameters of the ViT-based encoder and Table 7 shows the corresponding parameters of the CNN-based decoder.

ViTAE-SL is compared with the Kriging method and the VCNN. Kriging has been previously applied in a similar circumstance [44]. The performance of ViTAE-SL spatial data recovery for air pollution in China is shown in Figure 4. The reconstructed field by ViTAE-SL shows great agreement with the reference data in the test dataset. The advantage of ViTAE-SL is in partic-

| Model | Layers | Hidden Size D | Heads | Channel | Patch Size |
|--------------|--------|---------------|-------|---------|------------|
| ViT-Lite/32 | 8 | 32 | 8 | 16 | 32 |
| ViT-Base/32 | 8 | 64 | 8 | 32 | 32 |
| ViT-Large/32 | 8 | 128 | 8 | 64 | 32 |

Table 6: Vision Transformer encoder Parameters used in the study case of Chinese air quality.

| Model | Layers | Channel D | filter size |
|-----------|--------|-----------|-------------|
| CNN-Lite | 5 | 16 | 3 |
| CNN-Base | 5 | 32 | 3 |
| CNN-Large | 5 | 64 | 3 |

Table 7: CNN decoder parameters used in the study case of Chinese air quality data.

283 ular significant when comparing the contours of reconstruction. Table 8 and
 284 Table 9 show the online computational time and metrics of ViTAE-SL, VCNN
 285 and Kriging in the China air quality field reconstruction task, which demon-
 286 strate that ViTAE-SL spends much less time than kriging with a much better
 287 performance. Moreover, in terms of reconstruction accuracy, ViTAE-SL out-
 288 performs VCNN when the observations are sparse (i.e., 0.13% and 0.52%).
 289 In fact, by construction, ViTAE-SL is powerful in capturing chaotic local
 290 patterns in real-world field reconstruction task. To further evaluate the re-
 291 construction performance, two image similarity metrics Peak Signal-to-Noise
 292 Ratio (PSNR) and Structural Similarity Index Measure (SSIM) [45] are em-
 293 ployed to compare the reconstructed field against the original one. These
 294 metrics are widely used for similarity measures robust to rotation and trans-
 295 lation [46]. As shown in Figure 4, ViTAE-SL also has a higher PSNR and
 296 SSIM scores compared to Kriging and VCNN. These results speak to the
 297 significant advantage of ViTAE-SL in a real-world field reconstruction prob-
 298 lem.

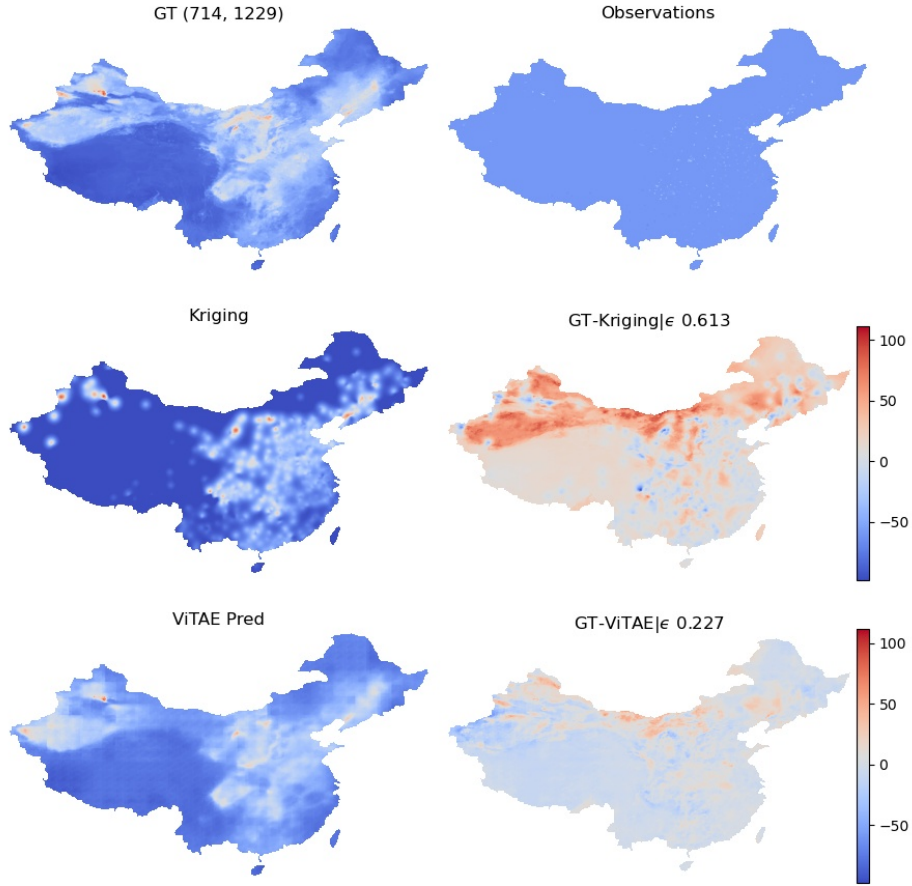


Figure 4: China air quality reconstruction results of ViTAE-SL and Kriging and comparison to the GT.

4.3. Test case 3: unsteady wake flow

To compare ViTAE-SL with the state-of-the-art VCNN in the case of unstructured grids and extremely sparse observations, we consider the two-dimensional unsteady laminar cylinder wake data. This test case was used in [30] for demonstrating the capability of VCNN. The training data set is prepared with a direct numerical simulation, which numerically solves the incompressible Navier–Stokes equations. In this study, we aim to reconstruct the vorticity field and we consider the same data that were used in [30]. The

| Model | ϵ | | | Execution time (s) |
|-----------------|---------------|---------------|---------------|--------------------|
| Kriging | 0.8599 | 0.8195 | 0.5520 | 117.5 |
| VCNN | 0.3978 | 0.3527 | 0.3261 | 0.0063 |
| ViTAE-lite | 0.3571 | 0.3389 | 0.3245 | 0.0053 |
| ViTAE-base | 0.3570 | 0.3434 | 0.3129 | 0.0058 |
| ViTAE-large | 0.3566 | 0.3371 | 0.3205 | 0.0061 |
| Monitor Percent | 0.13% | 0.52% | 2.08% | 0.13% |

Table 8: Performance in terms of Loss and execution time of ViTAE-SL, VCNN and Kriging for different observation density of China air quality reconstruction.

| Model | SSIM | | | | PSNR | |
|-----------------|---------------|---------------|---------------|----------------|----------------|----------------|
| Kriging | 0.7788 | 0.8195 | 0.8517 | 24.0098 | 25.0369 | 25.8168 |
| VCNN | 0.9385 | 0.9340 | 0.9321 | 29.7478 | 30.3167 | 30.4041 |
| ViTAE-lite | 0.9421 | 0.9330 | 0.9271 | 30.8184 | 30.8221 | 30.4489 |
| ViTAE-base | 0.9451 | 0.9366 | 0.9377 | 30.9079 | 30.7691 | 30.8588 |
| ViTAE-large | 0.9471 | 0.9399 | 0.9304 | 30.9466 | 30.9628 | 30.6220 |
| Monitor Percent | 0.13% | 0.52% | 2.08% | 0.13% | 0.52% | 2.08% |

Table 9: Performances of ViTAE-SL, VCNN and Kriging for different observation density of China air quality reconstruction.

307 shape of the entire field is 192×112 . The data span approximately 4 vortex
 308 shedding periods and the total number of snapshots is 5000. We randomly
 309 selected 80% of the snapshots as training data sets, 10% are used as vali-
 310 dation dataset and 10% for testing dataset. The number of sensors is set
 311 to 8 with fixed input sensor locations for both training and testing. Here
 312 we use the same model parameters as in Section 4.2. Figure 5 illustrates an
 313 example of the reconstruction fields in the testing dataset using VCNN and
 314 ViTAE-SL with the position of sensors. As shown in Figure 5, the vorticity
 315 fields reconstructed by ViTAE-SL are close to the GT. It can be seen from
 316 the reconstructed vorticity field that VCNN shows considerable reconstruc-
 317 tion error due to the small number of sensors; however, ViTAE-SL provides
 318 significantly more accurate reconstruction over the whole field in this exam-
 319 ple. Table 10 illustrates the averaged value of three different metrics, namely
 320 ϵ , PSNR and SSIM, over the testing dataset. Consistent with our analysis

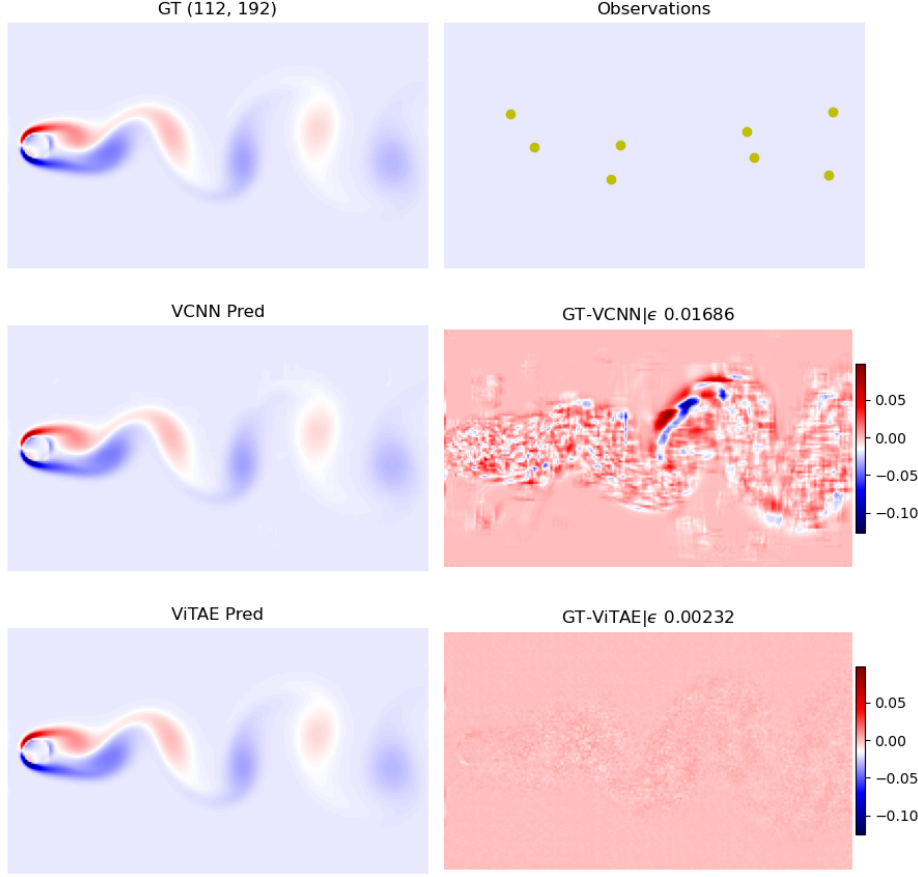


Figure 5: Unsteady wake flow reconstruction results of ViTAE-SL and VCNN compared to the GT.

321 of Figure 5, ViTAE-SL also outperforms VCNN in all three metrics. Both
 322 ViTAE-SL and VCNN run less than 0.01s in this reconstruction task. We
 323 here did not compare our results against Kriging, due to its poor performance
 324 with such sparse observations. This study shows that interpolation schemes,
 325 such as Kriging, and CNN-based methods, such as VCNN, are unable to
 326 reconstruct the fine-grained features of fields accurately if the observations
 327 are extremely sparse (e.g., 0.06% in this study). However, the proposed
 328 ViTAE-SL model is able to deliver a fast and accurate reconstruction of the

329 vorticity field of dimension using only 8 sensors.

| Model | ϵ | PSNR | SSIM | Execution time (s) |
|---------------|---------------|--------------|---------------|--------------------|
| VCNN | 0.0230 | 64.55 | 0.9988 | 0.0016 |
| ViTAT-lite/8 | 0.0112 | 70.63 | 0.9988 | 0.0043 |
| ViTAE-base/8 | 0.0050 | 77.56 | 0.9999 | 0.0050 |
| ViTAE-large/8 | 0.0023 | 84.37 | 0.9988 | 0.0052 |

Table 10: Unsteady wake flow reconstruction performances of ViTAE-SL and VCNN.

330 4.4. Test Case 4: NOAA global sea surface temperature

331 To show the capacity of ViTAE-SL in global physical field reconstruction,
 332 we used NOAA sea surface temperature data collected from satellite and
 333 ship-based observations. Data are made up of weekly observations of sea
 334 surface temperature with a spatial resolution of 360×180 . To make a
 335 fair comparison against VCNN [30], we use the same data set in [30]. 1040
 336 snapshots span from 1981 to 2001 are used to train the models and the test
 337 snapshots are taken from 2001 to 2018. In this test case, sensors are supposed
 338 to be randomly placed over the field. The number of sensors for training is
 339 set to $n_{\text{sensor,train}} = 10, 20, 30, 50, 100$ and for the test data, we also consider
 340 unseen cases with 70 and 200 sensors such that $n_{\text{sensor,test}} = 10, 20, 30, 50,$
 341 $70, 100, 200$.

342 An example of global sea surface temperature reconstruction using 10 sen-
 343 sors in the testing dataset is displayed in Figure 6. A significant advantage of
 344 the proposed ViTAE-SL model can be clearly observed with a considerably
 345 lower reconstruction error. The values of ϵ as defined in (5) are also indi-
 346 cated in Figure 6. As shown in Figure 6, the VCNN model fails to accurately
 347 reconstruct the details of the temperature fields when the number of sensors
 348 is severely limited (10 trained sensors in this case). In fact, the field recon-
 349 structed by VCNN also exhibits some non-realistic discontinuities due to the
 350 small number of sensors. On the other hand, the proposed ViTAE-SL is

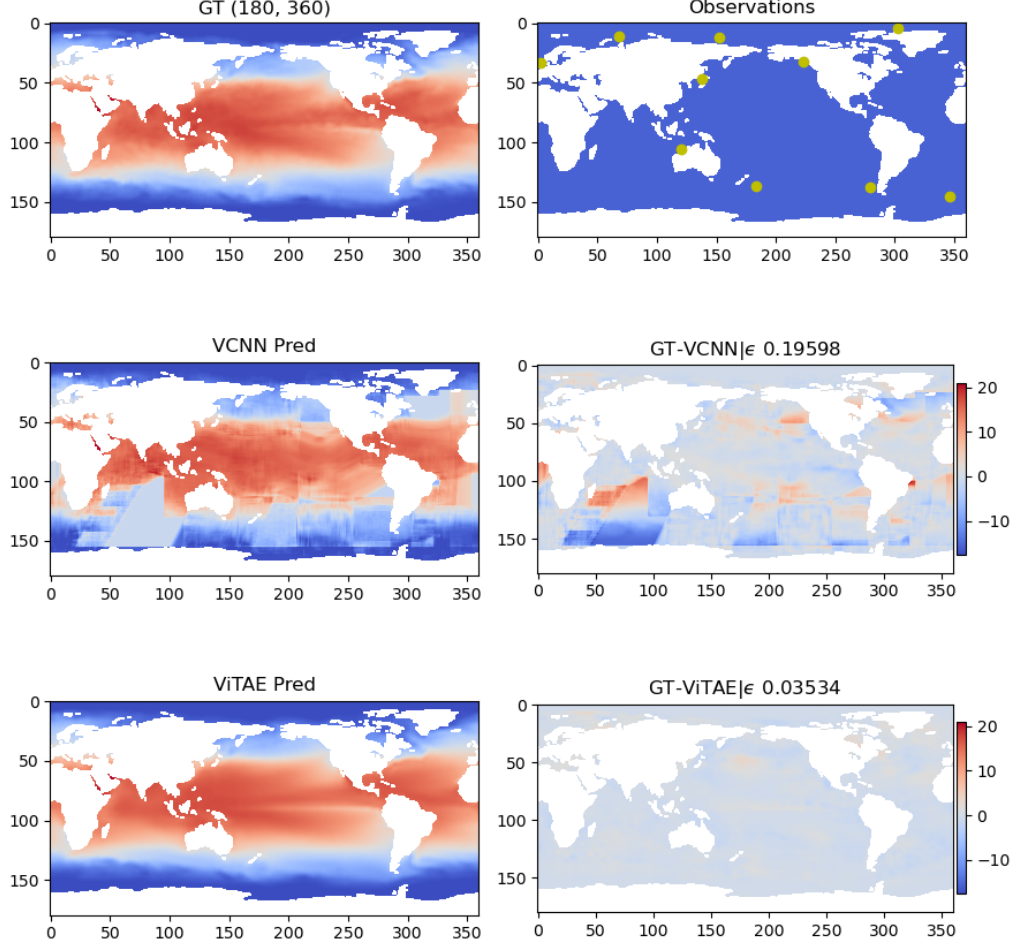


Figure 6: NOAA field reconstruction performance of ViTAE-SL and VCNN for trained 10 input observations, compared to the GT.

351 capable of delivering smooth and accurate reconstruction of the temperature
 352 field.

353 In this test case, we also aim to evaluate the impact of the number of sen-
 354 sors on the performance of VCNN and ViTAE-SL in terms of reconstruction
 355 accuracy. We plot the evolution of the ϵ error of both training and testing
 356 data against the number of sensors in Figure 7. As shown by the dashed

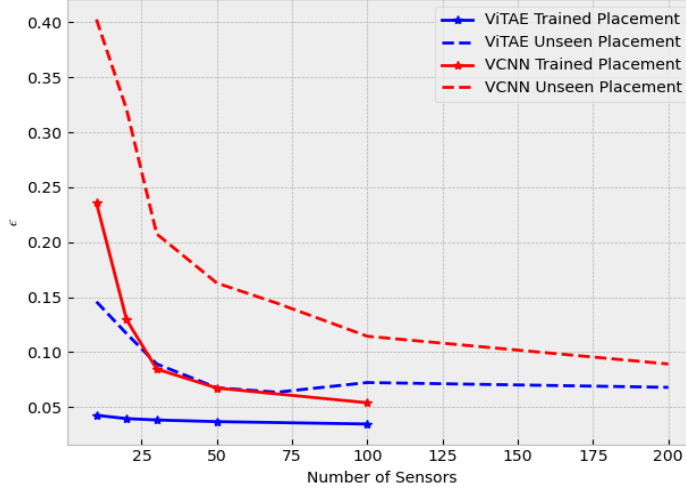


Figure 7: NOAA field reconstruction performances of ViTAE-SL and VCNN for various input observations.

blue line, when the number of sensors is over 30 for the test data, ViTAE-SL achieves a reasonably good reconstruction with the ϵ error being less than 0.1. However, as shown by the dashed red line, VCNN requires more than 200 sensors to obtain the same accuracy. In fact, the training of ViTAE-SL (solid blue line) also converges much faster compared to VCNN (solid red line).

5. Conclusion and Future Works

Spatial interpolation from sensor observations has been a longstanding challenge in physics systems. In order to address this problem, the paper proposes a novel ViT-based autoencoder which is an efficient spatial interpolation learner. The proposed novel method ViTAE-SL relies an attention mechanism in ViT to specify the sensor location and the global spatial dependence of the sensor measurements. The use of ViT splits the input sensor data to patches and translates patches into latent variables by a transformer block, which relying entirely on an attention mechanism to draw global de-

pendencies. Then the use of CNN decoder converts the latent variables to patches before concatenation. This process then enables the applications of CNNs to derive grid field reconstruction. The Vision Transformer encoder allows more parallelization and this requires significantly less computational time, which is crucial for large physical field reconstruction tasks. The CNN decoder also contributes to the physical field reconstruction consistency. The proposed method combines the advantages of both ViT and CNN, which make it outperform Kriging and the state-of-the-art VCNN.

The numerical results presented in the four test cases in this paper clearly suggests that the novel ViTAE-SL significantly outperforms the interpolation-based Kriging method and the CNN-based VCNN method. In particular, compared to Kriging, ViTAE requires no prior knowledge of the spatial distribution and reduces considerably the computational cost. The advantage of ViTAE-SL compared to VCNN is more significant when the data are sparse and unstructured, showing a great capacity of capturing spacial dependency. In summary, these experiments demonstrate that ViTAE-SL is capable of dealing with nonstationarity, nonlinear relationships, and non-Gaussian field reconstruction from scattered observations.

This work has provided a new perspective on DL in spatial prediction and could be applied to a wide range of studies of complex physics systems. In the future, the proposed method could be extended to a multivariate system to explore relationships and correlations between heterogeneous data. For example, this proposed method could be applied to integrate various types of measurements, such as low-cost air quality sensors, emission invention, traffic flow, and meteorological measurements, to enhance air quality field predictions. In this sense, it goes beyond spatial interpolation methods. This perspective allows researchers to explore the wealth of sensors measurements using data-driven technique, and will support scientific endeavor across a wide range of studies in physics system, spatio-temporal statistics and geostatistics.

402 **Code and data availability**

403 The code has been implemented using python and available at <https://github.com/fanhongweifd/ViTAE-for-field-reconstruction>. All ex-
404 periments are finished with the same RTX 3080 GPU.
405

406 **Acknowledgement**

407 This work is supported by the EPSRC grant EP/T003189/1 Health as-
408 sessment across biological length scales for personal pollution exposure and
409 its mitigation (INHALE). Sibó Cheng acknowledges the support of the French
410 Agence Nationale de la Recherche (ANR) under reference ANR-22-CPJ2-
411 0143-01.

412 **Acronyms**

| | | |
|-----|-----------------|--|
| 413 | DL | deep learning |
| 414 | NN | neural network |
| 415 | AE | autoencoder |
| 416 | NLP | Natural Language Processing |
| 417 | CNN | convolutional neural network |
| 418 | VCNN | Voronoi tessellation-assisted convolutional neural network |
| 419 | SSIM | Structural Similarity Index Measure |
| 420 | PSNR | Peak Signal-to-Noise Ratio |
| 421 | ViT | Vision Transformers |
| 422 | ViTAE-SL | Vision Transformer-based autoencoder and spatial |
| 423 | | interpolation learner |
| 424 | MLP | multi-layer perception |
| 425 | GCN | graph neural network |

426 **GT** ground truth
 427 **MSE** mean squared error

428 **References**

- 429 [1] Y. Cheng, X. Li, Z. Li, S. Jiang, Y. Li, J. Jia, X. Jiang, Aircloud: A cloud-based air-
 430 quality monitoring system for everyone, in: Proceedings of the 12th ACM Conference
 431 on Embedded Network Sensor Systems, 2014, pp. 251–265.
- 432 [2] N. Cressie, Statistics for spatial data, John Wiley & Sons, 2015.
- 433 [3] M. P. Austin, Spatial prediction of species distribution: an interface between ecolog-
 434 ical theory and statistical modelling, *Ecological modelling* 157 (2-3) (2002) 101–118.
- 435 [4] L. A. Waller, C. A. Gotway, Applied spatial statistics for public health data, John
 436 Wiley & Sons, 2004.
- 437 [5] M. A. Oliver, R. Webster, Kriging: a method of interpolation for geographical infor-
 438 mation systems, *Int. J. Geogr. Inf. Sci.* 4 (1990) 313–332.
- 439 [6] N. Cressie, The origins of kriging, *Mathematical Geology* 22 (1990) 239–252.
- 440 [7] C. E. Rasmussen, Gaussian processes in machine learning, in: Summer school on
 441 machine learning, Springer, 2003, pp. 63–71.
- 442 [8] C. J. Paciorek, M. J. Schervish, Spatial modelling using a new class of nonstation-
 443 ary covariance functions, *Environmetrics: The official journal of the International*
 444 *Environmetrics Society* 17 (5) (2006) 483–506.
- 445 [9] H. Zareifard, M. J. Khaledi, Non-gaussian modeling of spatial data using scale mixing
 446 of a unified skew gaussian process, *Journal of Multivariate Analysis* 114 (2013) 16–28.
- 447 [10] F. Karagulian, C. A. Belis, C. F. C. Dora, A. M. Prüss-Ustün, S. Bonjour, H. Adair-
 448 Rohani, M. Amann, Contributions to cities’ ambient particulate matter (pm): A
 449 systematic review of local source contributions at global level, *Atmospheric environ-*
 450 *ment* 120 (2015) 475–483.
- 451 [11] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, *nature* 521 (7553) (2015) 436–444.
- 452 [12] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convo-
 453 lutional neural networks, *Communications of the ACM* 60 (6) (2017) 84–90.

- 454 [13] G. Hadash, E. Kermany, B. Carmeli, O. Lavi, G. Kour, A. Jacovi, Estimate and
455 replace: A novel approach to integrating deep neural networks with existing applica-
456 tions, arXiv preprint arXiv:1804.09028 (2018).
- 457 [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser,
458 I. Polosukhin, Attention is all you need, Advances in neural information processing
459 systems 30 (2017).
- 460 [15] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in:
461 Proceedings of the IEEE conference on computer vision and pattern recognition, 2016,
462 pp. 770–778.
- 463 [16] C. Q. Casas, R. Arcucci, P. Wu, C. Pain, Y.-K. Guo, A reduced order deep data
464 assimilation model, *Physica D: Nonlinear Phenomena* 412 (2020) 132615.
- 465 [17] S. Cheng, C. Quilodr  n-Casas, S. Ouala, A. Farchi, C. Liu, P. Tandeo, R. Fablet,
466 D. Lucor, B. Iooss, J. Brajard, et al., Machine learning with data assimilation and
467 uncertainty quantification for dynamical systems: a review, *IEEE/CAA Journal of*
468 *Automatica Sinica* 10 (6) (2023) 1361–1387.
- 469 [18] R. Arcucci, J. Zhu, S. Hu, Y.-K. Guo, Deep data assimilation: integrating deep
470 learning with data assimilation, *Applied Sciences* 11 (3) (2021) 1114.
- 471 [19] S. Cheng, C. Liu, Y. Guo, R. Arcucci, Efficient deep data assimilation with sparse
472 observations and time-varying sensors, *Journal of Computational Physics* 496 (2024)
473 112581.
- 474 [20] T. H. Dur, R. Arcucci, L. Mottet, M. M. Solana, C. Pain, Y.-K. Guo, Weak constraint
475 gaussian processes for optimal sensor placement, *Journal of Computational Science*
476 42 (2020) 101110.
- 477 [21] C. Buizza, C. Q. Casas, P. Nadler, J. Mack, S. Marrone, Z. Titus, C. Le Cornec,
478 E. Heylen, T. Dur, L. B. Ruiz, et al., Data learning: Integrating data assimilation
479 and machine learning, *Journal of Computational Science* 58 (2022) 101525.
- 480 [22] N. B. Erichson, L. Mathelin, Z. Yao, S. L. Brunton, M. W. Mahoney, J. N. Kutz, Shal-
481 low neural networks for fluid flow reconstruction with limited sensors, *Proceedings of*
482 *the Royal Society A* 476 (2238) (2020) 20200097.

- 483 [23] T. Bolton, L. Zanna, Applications of deep learning to ocean data inference and subgrid
484 parameterization, *Journal of Advances in Modeling Earth Systems* 11 (1) (2019) 376–
485 399.
- 486 [24] J. Yu, J. S. Hesthaven, Flowfield reconstruction method using artificial neural net-
487 work, *Aiaa Journal* 57 (2) (2019) 482–498.
- 488 [25] J. Wu, D. Xiao, M. Luo, Deep-learning assisted reduced order model for high-
489 dimensional flow prediction from sparse data, *Physics of Fluids* 35 (10) (2023).
- 490 [26] F. Concas, J. Mineraud, E. Lagerspetz, S. Varjonen, X. Liu, K. Puolamäki, P. Nurmi,
491 S. Tarkoma, Low-cost outdoor air quality monitoring and sensor calibration: A survey
492 and critical analysis, *ACM Transactions on Sensor Networks (TOSN)* 17 (2) (2021)
493 1–44.
- 494 [27] S. Cheng, J.-P. Argaud, B. Iooss, D. Lucor, A. Ponçot, Error covariance tuning in vari-
495 ational data assimilation: application to an operating hydrological model, *Stochastic
496 Environmental Research and Risk Assessment* 35 (5) (2021) 1019–1038.
- 497 [28] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, S. Y. Philip, A comprehensive survey on
498 graph neural networks, *IEEE transactions on neural networks and learning systems*
499 32 (1) (2020) 4–24.
- 500 [29] I. C. L. Environmental Research Group, London air quality network, [https://www.
501 londonair.org.uk/LondonAir/Default.aspx](https://www.londonair.org.uk/LondonAir/Default.aspx) (2022 (accessed Nov 8, 2022)).
- 502 [30] K. Fukami, R. Maulik, N. Ramachandra, K. Fukagata, K. Taira, Global field recon-
503 struction from sparse sensors with voronoi tessellation-assisted deep learning, *Nature
504 Machine Intelligence* 3 (11) (2021) 945–951.
- 505 [31] D. Linsley, J. Kim, V. Veerabadran, C. Windolf, T. Serre, Learning long-range spatial
506 dependencies with horizontal gated recurrent units, *Advances in neural information
507 processing systems* 31 (2018).
- 508 [32] S. Jain, A. A. Presto, N. Zimmerman, Spatial modeling of daily pm_{2.5}, no₂, and co
509 concentrations measured by a low-cost sensor network: comparison of linear, machine
510 learning, and hybrid land use models, *Environmental Science & Technology* 55 (13)
511 (2021) 8631–8641.

- [33] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, arXiv preprint arXiv:2010.11929 (2020).
- [34] R. Fu, D. Xiao, I. M. Navon, F. Fang, L. Yang, C. Wang, S. Cheng, A non-linear non-intrusive reduced order model of fluid flow by auto-encoder and self-attention deep learning methods, International Journal for Numerical Methods in Engineering 124 (13) (2023) 3087–3111.
- [35] X. Pan, D. Xiao, Domain decomposition for physics-data combined neural network based parametric reduced order modelling, Journal of Computational Physics 519 (2024) 113452.
- [36] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, R. Girshick, Masked autoencoders are scalable vision learners, arXiv preprint arXiv:2111.06377 (2021).
- [37] Z. Wei, Jing; Li, Chinahighpm2.5, <https://zenodo.org/record/5919482#.Y2ty-NJByV4> (2022 (accessed Nov 8, 2022)).
- [38] S. Müller, Geostat framework, <https://geostat-framework.org/> (2022 (accessed Nov 8, 2022)).
- [39] W. Chen, Y. Li, B. J. Reich, Y. Sun, Deepkriging: Spatially dependent deep neural networks for spatial prediction, arXiv preprint arXiv:2007.11972 (2020).
- [40] D. Ginsbourger, D. Dupuy, A. Badea, L. Carraro, O. Roustant, A note on the choice and the estimation of kriging models for the analysis of deterministic computer experiments, Applied Stochastic Models in Business and Industry 25 (2) (2009) 115–131.
- [41] S. Cheng, I. C. Prentice, Y. Huang, Y. Jin, Y.-K. Guo, R. Arcucci, Data-driven surrogate model with latent data assimilation: Application to wildfire forecasting, Journal of Computational Physics (2022) 111302.
- [42] S. Cheng, J. Chen, C. Anastasiou, P. Angeli, O. K. Matar, Y.-K. Guo, C. C. Pain, R. Arcucci, Generalised latent assimilation in heterogeneous reduced spaces with machine learning surrogate models, Journal of Scientific Computing 94 (1) (2023) 1–37.
- [43] C. N. E. M. Centre, Chinese air quality, <https://air.cnemc.cn/> (2022 (accessed Nov 8, 2022)).

- 543 [44] Y. Zhan, Y. Luo, X. Deng, K. Zhang, M. Zhang, M. L. Grieneisen, B. Di, Satellite-
544 based estimates of daily no₂ exposure in china using hybrid random forest and spa-
545 tiotemporal kriging model, *Environmental science & technology* 52 (7) (2018) 4180–
546 4189.
- 547 [45] A. Hore, D. Ziou, Image quality metrics: Psnr vs. ssim, in: 2010 20th international
548 conference on pattern recognition, IEEE, 2010, pp. 2366–2369.
- 549 [46] D. Lee, H. Park, I. K. Park, K. M. Lee, Joint blind motion deblurring and depth
550 estimation of light field, in: *Proceedings of the European Conference on Computer*
551 *Vision (ECCV)*, 2018, pp. 288–303.