
NIERT: Accurate Numerical Interpolation through Unifying Scattered Data Representations using Transformer Encoder

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Numerical interpolation for scattered data aims to estimate values for target points
2 based on those of some observed points. Traditional approaches produce estima-
3 tions through constructing an interpolation function that combines multiple basis
4 functions. These approaches require the basis functions to be pre-defined explicitly,
5 thus greatly limiting their applications in practical scenarios. Recent advances
6 exhibit an alternative strategy that learns interpolation functions directly from
7 observed points using machine learning techniques, say deep neural networks. This
8 strategy, although promising, cannot effectively exploit the correlations between
9 observed points and target points as it treats these types of points separately. Here,
10 we present a learning-based approach to numerical interpolation using encoder
11 representations of Transformers (thus called NIERT). NIERT treats the value of
12 each target point as a masked token, which enables processing target points and
13 observed points in a unified fashion. By calculating the partial self-attention be-
14 tween target points and observed points at each layer, NIERT gains advantages
15 of exploiting the correlations among these points and, more importantly, avoiding
16 the unexpected interference of target points on observed points. NIERT also uses
17 the pre-training technique to further improve its accuracy. On three representative
18 datasets, including two synthetic datasets and a real-world dataset, NIERT outper-
19 forms the existing approaches, e.g., on the TFRD-ADlet dataset for temperature
20 field reconstruction, NIERT achieves an MAE of 1.897×10^{-3} , substantially better
21 than the transformer-based approach (MAE: 27.074×10^{-3}). These results clearly
22 demonstrate the accuracy of NIERT and its potential to apply in multiple practical
23 fields.

24 1 Introduction

25 Numerical interpolation for scattered data plays important and fundamental roles in a wide range
26 of practical scenarios, including solving partial differential equations (PDEs) [1], temperature field
27 reconstruction [2], time series interpolation [3, 4]. In meshfree PDE solvers, the interpolation
28 error often leads to deviations in subsequent calculations, which seriously affects the solution’s
29 accuracy [5]. In the task of temperature field reconstruction for micro-scale electronics, interpolation
30 methods are used to obtain the real-time working environment of electronic components from
31 limited measurements, and imprecise interpolation will significantly increase the cost of predictive
32 maintenance [2]. Thus, accurate approaches to numerical interpolation are highly desirable.

33 A large number of approaches have been proposed for interpolation of scattered data, which can
34 be divided into two categories, namely, traditional non-learning based methods and recent learning-
35 based methods. The typical traditional interpolation schemes construct the target function by a linear

36 combination of basis functions [6]. These schemes require explicitly-defined basis functions to
37 model the target function space, and various types of basis functions have been devised by algorithm
38 designers to adapt to different scenarios. Nevertheless, such methods still suffer from limitations
39 of high requirement of sufficient observed points, and the limited complexity of the target function.

40 Recent progress has exhibited an alternative strategy that uses neural networks to learn interpolation
41 functions directly from the given observed points. For example, conditional neural processes (CNPs)
42 [7] and their extensions [8–10] model the conditional distribution of regression functions given
43 the observed points. In addition, Chen et al. [2] proposed to use vanilla Transformer [11] to solve
44 interpolation task in temperature field reconstruction. All of these approaches use an “encoder-
45 decoder” architecture, in which the encoder learns the representations of observed points while the
46 decoder estimates values for target points. Intuitively, observed points and target points should be
47 processed in a unified fashion because they are from the same domain. However, this architecture
48 treats them separately and cannot effectively exploit the correlation between them.

49 Inspired by the recent advances of language/image models, especially BERT [12] and BEiT [13], we
50 designed an approach to numerical interpolation that can effectively exploit the correlations between
51 observed points and target points. Our approach is a learning-based approach using the encoder
52 representations of Transformers (thus called NIERT). The key elements of NIERT include: *i*) the
53 use of the mask mechanism, which enables processing target points and observed points in a unified
54 fashion, *ii*) a novel partial self-attention model, which calculates attentions between target points and
55 observed points at each layer, thus gaining the advantages of exploiting the correlations between these
56 two types of points and, more importantly, avoiding the unexpected interference of target points on
57 observed points simultaneously, and *iii*) the use of the pre-training technique, which further improves
58 the interpolation accuracy of NIERT.

59 The main contributions of this study are summarized as follows.

- 60 1. We propose an accurate approach to numerical interpolation for scattered data. On represen-
61 tative datasets, including both synthetic and real-world datasets, our approach outperforms
62 existing approaches. The experimental results demonstrate the potential of our approach
63 in a wide range of application fields. The source code of NIERT will be released for open
64 source use.
- 65 2. We propose a novel partial self-attention mechanism to make Transformer incorporated with
66 strong inductive bias for interpolation tasks; i.e., it can effectively exploit the correlation
67 among two types of points but simultaneously avoid the interference of one type of points
68 onto the others.
- 69 3. We propose to use the pre-training technique to enhance interpolation approaches. When
70 facing an interpolation task in a newly-appearing application field, we can benefit from the
71 experience learned from low-cost synthesized interpolation tasks.

72 **2 Related works**

73 **2.1 Traditional interpolation approaches for scattered data**

74 Traditional interpolation approaches for scattered data use explicit basis functions to construct
75 interpolation function, e.g., Lagrange interpolation, Newton interpolation [6], B-spline interpolation
76 [14], Shepard’s method [15], Kriging [16], and radial basis function interpolation (RBF) [17, 18].
77 Among these approaches, the classical Lagrange interpolation, Newton interpolation and B-splines
78 interpolation are usually used for univariate interpolation. Wang et al. [19] proposed a high order
79 multivariate approximation scheme for scattered data sets, in which approximation error is represented
80 with Taylor expansions at data points, and basis functions are determined through minimizing the
81 squares of approximation error.

82 **2.2 Neural network-based interpolation approaches**

83 Equipped with deep neural networks, data-driven interpolation and reconstruction methods show great
84 advantages and potential. For instance, convolutional neural networks (CNNs) have been applied
85 in the interpolation tasks of single image super-resolution [20, 21], and recurrent neural networks
86 (RNNs) and Transformers have been used for interpolation of sequences like time series data [4, 22].

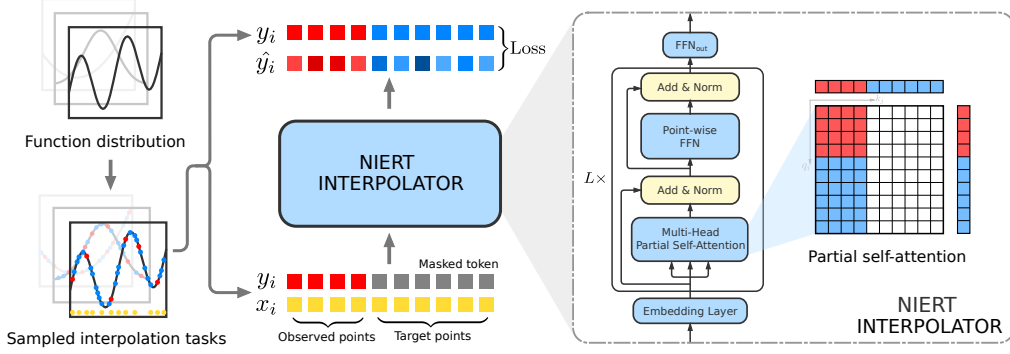


Figure 1: Overview of NIERT training process. Here, x_i represents the position of a point, and y_i represents its value. The predicted values of the point is denoted as \hat{y}_i . We prepare the training interpolation tasks by first sampling functions from a function distribution \mathcal{F} and then sampling observed points O and target points T on each function. NIERT trains an interpolator over this data. The partial self-attention mechanism facilitates exploiting the correlations between observed points and target points and avoiding the unexpected interference of target points on observed points

87 Recently, Garnelo et al. [7] proposed to model the conditional distribution of regression functions
 88 given observed points. The proposed approach, conditional neural processes (CNPs), has shown
 89 increased estimation accuracy and generalizing ability. Kim et al. [8] designed an enhanced model,
 90 attentive neural processes (ANPs), with improved accuracy. Lee et al. [10] leveraged Bayesian
 91 last layer (BLL) [23] for faster training and better prediction. In addition, the bootstrap technique
 92 was also employed for further improvement [9]. To solve the interpolation task in 2D temperature
 93 field reconstruction, Chen et al. [2] proposed an Transformer-based approach, referred to as TFR-
 94 transformer, which can also be applied to solve interpolation tasks for scattered data with higher
 95 dimensions. Note that although TFR-transformer and our NIERT are both based on transformer,
 96 they are fundamentally different: *i*) TFR-transformer adopts an encoder-decoder structure and treats
 97 observed points and target points respectively, while NIERT adopts only a Transformer encoder
 98 (equipped with partial self-attention) to encode and learn correlation of the scattered data in a unified
 99 fashion; *ii*) TFR-transformer is trained by minimizing the prediction error of target points, while
 100 NIERT’s training objective considers the prediction error of both observation points and target points.

101 2.3 Masked language/image models and the pre-training technique

102 The design of NIERT is also inspired by the recent advances in masked language/image models
 103 [12, 13, 24, 25] and pre-trained models for symbolic regression [26, 27]. The masked pre-trained
 104 models have been shown to be successful in learning representations of languages and images
 105 from large-scale data and improving the performance of downstream tasks. In addition, the mask
 106 mechanism makes the model able to reconstruct missing data from their context. Utilizing large-scale
 107 synthetic symbolic functions and sampled scattered data, Biggio et al. [26] and Valipour et al. [27]
 108 pre-trained Transformers to learn the map from scattered data to corresponding symbolic formulas.
 109 Different from these approaches, NIERT uses synthetic data to learn interpolating functions from
 110 scattered data numerically.

111 3 Method

112 3.1 Overview of NIERT

113 In the study, we focus on the interpolation task that can be formally described as follows: We are
 114 given n observed points with known values $O = \{(x_i, y_i)\}_{i=1}^n$, and m target points with values to be
 115 determined, denoted as $T = \{x_i\}_{i=n+1}^{n+m}$. Here, $x_i \in X$ denotes position of a point, $y_i = f(x_i) \in Y$
 116 denotes the value of a point, and $f: X \rightarrow Y$ denotes a function mapping positions to values. The
 117 function f is from a function distribution \mathcal{F} , which can be explicitly defined using a mathematical
 118 formula or implicitly represented using a set of scattered data in the form (x_i, y_i) . The goal of

Algorithm 1 NIERT training process

Require: NIERT model M_θ with parameters θ , epoch number N , batch size B , domain X , function distribution \mathcal{F} and error metric function $\text{Error}(\cdot, \cdot)$

for k in $\{1..N\}$ **do**
 $J \leftarrow 0$
 for b in $\{1..B\}$ **do**
 $f, \{x_i\}_i \leftarrow$ sample a function and scatter points from \mathcal{F}, X
 $\{y_i\}_i \leftarrow$ calculate f on $\{x_i\}_i$
 $O, T \leftarrow$ split and mask scatter points with values $\{(x_i, y_i)\}_i$
 $\{\hat{y}_i\}_i \leftarrow M_\theta(O, T)$
 $J \leftarrow J + \sum_i \text{Error}(\hat{y}_i, y_i)$
 end for
 Compute the gradient $\nabla_\theta J$ and update θ
end for

119 interpolation task is to accurately estimate the values $f(x)$ for each target point $x \in T$ according to
120 the observed points in O .

121 Figure 1 depicts the schematic diagram of our NIERT approach. Briefly speaking, our approach
122 employs a data-driven approach to numerical interpolation using encoder representations of Trans-
123 formers. The main element of our approach is a neural interpolator that learns to estimate values for
124 target points. The interpolator is featured by the characteristic that it treats the value of each target
125 point as a masked token, thus enabling the unifying fashion to process both target points and observed
126 points in the subsequent encoding and estimation procedures.

127 To suit the interpolation task, we design a *partial self-attention* mechanism: on one side, we calculate
128 the attention between target points and observed points at each layer, which gains NIERT the
129 advantage to effectively exploit the correlations between these two types of points. On the other side,
130 the attention is a partial one as we do not consider the effects of a target point on all other points.
131 This way, the unexpected interference of target points onto the observed points, and the interference
132 among target points, are completely avoided.

133 The training process is depicted in Algorithm 1. Specifically, we prepare the training interpolation
134 tasks by first sampling functions from a distribution \mathcal{F} and then sampling observed points O and
135 target points T on each function. When training NIERT, we set the loss function as the error between
136 the estimated values and the corresponding ground-truth. It should be pointed out that errors acquired
137 on both observed points and target points are accounted into loss function.

138 3.2 Architecture of the NIERT interpolator

139 The neural interpolator in NIERT adopts the Transformer encoder framework; however, to suit the
140 interpolation task, significant modifications and extensions were made in embedding, Transformer
141 and output layers, which are described in details below.

142 **Embedding with masked tokens:** NIERT embeds both observed points and target points into the
143 unified high-dimensional embedding space. As the position x of a data point and its value y are from
144 different domains, we use two linear modules : Linear_x embeds the positions while Linear_y embeds
145 the values.

146 It should be noted that for target points, their values are absent when embedding as they are to be
147 determined. In this case, we use a masked token as substitutes, which is embedded as a trainable
148 parameter MASK_y as performed in BERT [12]. This way, the interpolator processes both target
149 points and observed points in a unifying fashion.

150 We concatenate the embeddings of position and value of a data point as the point’s embedding,
151 denoted as h_i^0 , i.e.,

$$h_i^0 = \begin{cases} [\text{Linear}_x(x_i), \text{Linear}_y(y_i)], & \text{if } (x_i, y_i) \in O \\ [\text{Linear}_x(x_i), \text{MASK}_y], & \text{if } x_i \in T \end{cases}$$

152 **Transformer layer with partial self-attention mechanism:** NIERT feeds the embeddings of the
 153 points into a stack of L Transformer layers, producing encodings of these points as results. Each
 154 Transformer layer contains two subsequent sub-layers, namely, a multi-head self-attention module,
 155 and a point-wise fully-connected network. These sub-layers are interlaced with residual connections
 156 and layer normalization between them.

157 To avoid the unexpected interference of target points on observed points and target points themselves,
 158 NIERT replaces the original self-attention in Transformer layer with a *partial self-attention*, which
 159 calculates the feature of point i at the $l + 1$ -st layer as follows:

$$160 \quad h^{l+1} = \text{LayerNorm}(\tilde{v}^l + \text{MLP}(\tilde{v}^l)),$$

$$\tilde{v}_i^l = \text{LayerNorm}\left(v_i^l + \sum_j w_{ij} \alpha_{ij}^l v_j^l\right)$$

161 where v_i^l and α_{ij}^l represent the value embedding and ordinary attention weights at the l -th layer as
 162 calculated in Transformer [11]. In this formula, we introduce a new term w_{ij} that represents the
 163 partial self-attention pattern, i.e.,

$$w_{ij} = \begin{cases} 1, & \text{if } (x_j, y_j) \in O \\ 0, & \text{if } x_j \in T \end{cases}.$$

164 By forcing the weight w_{ij} to be 0 for a target point i and any point j , we completely avoid the
 165 unexpected interference of target points on the other points.

166 **Estimating values for target points:** For each target point i , we estimate its value \hat{y}_i through feeding
 167 its features at the final Transformer layer into a fully connected feed-forward network, i.e.,

$$\hat{y}_i = \text{MLP}_{\text{out}}(h_i^L).$$

168 We calculate the error between the estimation and the corresponding ground-truth value, and compose
 169 the errors for all points into a loss function to be minimized.

170 3.3 Enhancing NIERT using pre-training technique

171 The interpolation functions from different applications usually differ greatly in their forms; however,
 172 the interpolation tasks might still share some common characteristics, say the correlation between
 173 observed points and target points. These common characteristics enable enhancing NIERT using
 174 the pre-training technique. Here, we pre-train NIERT using a synthetic dataset (see 4.1 for further
 175 details) and fine-tune it on other datasets in application fields.

176 4 Experiments and results

177 We evaluated NIERT and compared it with ten representative scattered data interpolation approaches
 178 on both synthetic and real-world datasets. We also examined the effects of the key elements of NIERT,
 179 including the partial self-attention, and the pre-training technique.

180 4.1 Experiment setting

181 The datasets, metrics and approaches for comparison are briefly described below. Further details of
 182 experiment settings are provided in Supplementary text.

183 **Datasets:** Three representative datasets in various application fields, including two synthetic datasets
 184 NeSymReS and TFRD-ADlet, and real-world dataset PhysioNet, are used for evaluation.

185 NeSymReS is a synthetic dataset for mathematical function interpolation, which is built using data
 186 generator proposed by Biggio et al. [26] and Lample and Charton [28]. We construct a function set
 187 with various dimensionality of data points, including 1D, 2D, 3D, and 4D. Scattered points in each
 188 instance are randomly sampled and divided into observed points and target points.

189 TFRD-ADlet [2] is a synthetic dataset for 2D temperature field reconstruction where each instance
 190 represents a simulated 2D temperature field containing several heat source components and a specific

191 Dirichlet conditioned boundary. The goal of each task instance is to reconstruct the whole temperature
 192 field according to a limited number of observed points with measured temperature.

193 PhysioNet Challenge 2012 dataset [29] is a real world dataset collected from intensive care unit
 194 (ICU) records for time-series data interpolation. Each point in an instance represents a measurement
 195 at a specific time, where each measurement contains up to 37 physiological indices. Following the
 196 study [22], we randomly divided the points into observed points and target points by setting the ratio
 197 of observed points at five levels, i.e., 50%, 60%, 70%, 80%, and 90%. Note that this dataset is a
 198 representative of hard interpolation tasks due to the sparsity and irregularity of the records.

199 **Pre-training dataset:** In this study, NeSymReS dataset was used for pre-training NIERT to further
 200 improve its interpolation accuracy on TFRD-ADlet and PhysioNet dataset. For TFRD-ADlet dataset
 201 we directly use 2D TFRD-ADlet dataset for pre-training. As the PhysioNet dataset has a dimen-
 202 sionality of 37, we construct the pre-training instances by stacking random 37 functions from 1D
 203 TFRD-ADlet dataset and then sampling interpolation task instances.

204 **Metrics:** When evaluating NIERT and other interpolation approaches, the prediction error of target
 205 points are calculated as interpolation accuracy. For the NeSymReS and PhysioNet dataset, we adopted
 206 mean squared error (MSE) as the error metric. For TFRD-ADlet dataset, we use three error metrics:
 207 mean absolute error (MAE), MAE in component area (CMAE) and MAE at boundary (BMAE)
 208 following Chen et al. [2]. Accordingly, we use L_2 -form loss function for NeSymReS and PhysioNet
 209 dataset and L_1 -form loss function for TFRD-ADlet dataset for training.

210 **Approaches for comparison:** For NeSymReS dataset, we compared NIERT with six representative
 211 interpolation approaches, including RBF[30], MIR [19], CNP[7], ANP[8], BANP[9] and TFR-
 212 transformer [2]. For TFRD-ADlet we compared NIERT with CNP, ANP, BANP and TFR-transformer.
 213 For PhysioNet we compared NIERT with four approaches designed for time-series data interpolation,
 214 including RNN-VAE [31], L-ODE-RNN [32], L-ODE-ODE[33], and mTAND-Full[22].

215 4.2 Interpolation accuracy on synthetic and real-world datasets

216 For each instance of the test datasets, we applied the trained NIERT to estimate values for the target
 217 points. We calculate the errors between the estimation and the ground-truth as interpolation accuracy.

Interpolation approach	MSE ($\times 10^{-5}$) on NeSymReS test set			
	1D	2D	3D	4D
RBF	215.439	347.060	443.094	327.775
MIR	67.281	274.601	448.933	342.997
CNP	67.176	248.668	392.348	314.311
ANP	34.558	140.005	206.699	164.751
BANP	14.913	84.187	143.518	140.288
TFR-transformer	15.556	58.569	99.986	90.579
NIERT	8.964	45.319	77.664	72.025

218

Interpolation approach	Evaluation criteria ($\times 10^{-3}$)		
	MAE	CMAE	BMAE
CNP	96.674	109.419	56.939
ANP	54.684	62.511	26.524
BANP	28.671	29.450	19.984
TFR-transformer	27.074	29.772	18.835
NIERT	3.473	3.947	2.467
NIERT w/ pretraining	1.897	1.971	1.246

Table 1: Interpolation accuracy of NIERT and the existing approaches on NeSymReS test dataset

Table 2: Interpolation accuracy of NIERT and the existing approaches over the TFRD-ADlet dataset

Interpolation approach	Ratio of observed points				
	50%	60%	70%	80%	90%
RNN-VAE	13.418 \pm 0.008	12.594 \pm 0.004	11.887 \pm 0.005	11.133 \pm 0.007	11.470 \pm 0.006
L-ODE-RNN	8.132 \pm 0.020	8.140 \pm 0.018	8.171 \pm 0.030	8.143 \pm 0.025	8.402 \pm 0.022
L-ODE-ODE	6.721 \pm 0.109	6.816 \pm 0.045	6.798 \pm 0.143	6.850 \pm 0.066	7.142 \pm 0.066
mTAND-Full	4.139 \pm 0.029	4.018 \pm 0.048	4.157 \pm 0.053	4.410 \pm 0.149	4.798 \pm 0.036
NIERT	2.868 \pm 0.021	2.811 \pm 0.032	2.656 \pm 0.041	2.598 \pm 0.078	2.709 \pm 0.157
NIERT w/ pretraining	2.831\pm0.021	2.771\pm0.019	2.641\pm0.052	2.539\pm0.085	2.596\pm0.159

Table 3: The relationship between interpolation accuracy (measured using MSE, $\times 10^{-3}$) and the ratio of observed points. Here, we use the PhysioNet dataset as representatives

219 **Accuracy on the NeSymReS dataset:** As shown in Table 1, on the 1D NeSymReS testset, RBF
 220 shows the largest interpolation error (MSE: 215.439). MIR, another approach using explicit basis
 221 functions, also shows a high interpolation error of 67.281. In contrast, BANP and TFR-transformer,
 222 which use neural networks to learn interpolation, show relatively lower errors (MSE: 14.913, 15.556).
 223 Compared with these approaches, our NIERT approach achieves the best interpolation accuracy

224 (MSE: 8.964). Table 1 also demonstrates the advantage of NIERT over the existing approach on the
 225 2D, 3D, and 4D instances.

226 Note that the number of observed points varies
 227 greatly in test instances, making the average
 228 interpolation error calculated over all test in-
 229 stances insufficient to measure interpolation per-
 230 formance. Therefore, we further divide test in-
 231 stances into subsets according to the number of
 232 observed points. As shown in Figure 2, as the
 233 number of observed points increases, the inter-
 234 polation error decreases as expected. In addition,
 235 the relative advantages of these approaches vary
 236 with the number of observed points, e.g., CNP is
 237 better than RBF and MIR initially but finally be-
 238 comes worse as the number of observed points
 239 increases. Among all approaches, NIERT stably
 240 shows the best performance over all test subsets,
 241 regardless of the number of observed points.

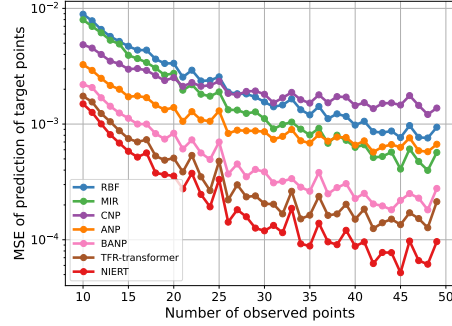


Figure 2: The relationship between the interpolation accuracy and the number of observed points. Here we use the 2D instances in the NeSymReS test dataset as representatives

242 **Accuracy on the TFRD-ADlet dataset:** As shown in Table 2, CNP, although employing the neural
 243 network technique, still performs poorly with MAE as high as of 96.674. In contrast, NIERT achieves
 244 the lowest interpolation error (MAE: 3.473), which is over one order of magnitude lower than CNP,
 245 ANP, BANP and TFR-transformer. Moreover, when enhanced with the pre-training technique, NIERT
 246 can further decrease its interpolation MAE to be 1.897. Besides MAE, other metrics, say CMAE and
 247 BMAE, also show the superior of NIERT over the existing approaches (Table 2).

248 **Accuracy on the PhysioNet dataset:** Table 3 suggests that on the PhysioNet dataset, NIERT also
 249 outperforms the existing approaches, e.g., when controlling the ratio of observed points to be 50%,
 250 NIERT achieves an average MSE ($\times 10^{-3}$) of 2.868, significantly lower than RNN-VAE (13.418),
 251 L-ODE-RNN (8.132), L-ODE-ODE (6.721) and mTAND-Full (4.139). Again, NIERT with the
 252 pre-training technique shows better performance. The advantages of NIERT hold across various
 253 settings of the ratio of the observed points.

254 Taken together, these results clearly demonstrate the power of NIERT for numerical interpolation in
 255 multiple application fields, including interpolating the scattered data generated using mathematical
 256 functions, reconstructing temperature fields, and interpolating time-series data.

257 4.3 Case studies of interpolation results

258 To further understand the advantages of NIERT, we carried out case studies through visualizing
 259 the observed points, the reconstructed interpolation functions and the interpolation errors in this
 260 subsection. More visualized cases are put in the Supplementary material.

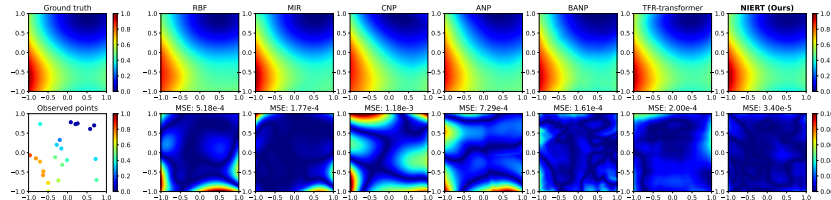


Figure 3: An example of 2D interpolation task extracted from NeSymReS test set. The up-left figure shows the ground-truth function while the bottom-left figure shows the 22 observed points. The interpolation functions reported by NIERT and the existing approaches are listed on the top panel with their differences with the ground-truth are list below

261 Figure 3 show a 2D instance in the NeSymReS test set, respectively. As illustrated by these two
 262 figures, RBF performs poorly in the application scenario with sparse observed data. In addition, RBF
 263 and MIR, especially ANP, cannot accurately predict values for the target points that fall out of the
 264 range restricted by observed points. The CNP approach can only learn the rough trend stated by the

265 observed points, thus leading to significant errors. In contrast, BANP, TFR-transformer and NIERT
 266 can accurately estimate values for target points within considerably large range, and compared with
 267 BANP and TFR-transformer, NIERT can produce more accurate results.

268 Figure4 shows an instance of temperature field reconstruction extracted from TFRD-ADlet. From
 269 this figure, we can observe that when using pre-training technique, NIERT further improves its
 270 interpolation accuracy in the whole area.

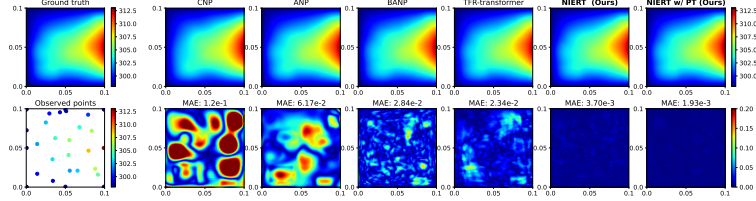


Figure 4: An example of temperature field reconstruction task extracted from TFRD-ADlet test set. The up-left figure shows the ground-truth temperature field while the bottom-left figure shows the 32 observed points. The reconstructed results reported by NIERT and the existing approaches are listed on the top panel with their differences with the ground-truth temperature field are list below

271 4.4 Contribution analysis of observed points for interpolation

272 An idealized interpolation approach is expected to effectively exploit all observed points with
 273 appropriate consideration of relative positions among observed points and target points as well. To
 274 examine this issue, we visualized the attention weight of each observed point to all target points.
 275 These attention weights provide an intuitive description of the contribution by observed points.

276 As shown in Figure 5, when using TFR-transformer, the contributions by observed points are
 277 considerably imbalanced: on one side, some observed points might affect their neighboring target
 278 points in a large region; on the other side, the other observed points have little contributions to
 279 interpolation. In contrast, when using NIERT, contributions by an observed point are much more
 280 local and thus targeted. More importantly, all observed points have contributions to interpolation.

281 These results demonstrate that NIERT can exploit the correlation between observed points and target
 282 points more effectively.

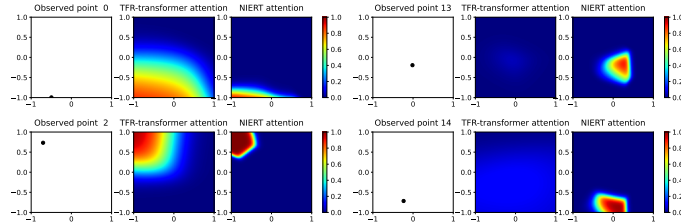


Figure 5: Contributions by observed points for interpolation. The 2D instance is same to that used in Figure 3. We randomly select 4 observed points and extract their attention weights from the final attention layer of NIERT and TFR-transformer. These attention weights provide an intuitive description of the contribution by observed points. The contributions by other 18 observed points are shown in Supplementary material

283 4.5 Ablation study

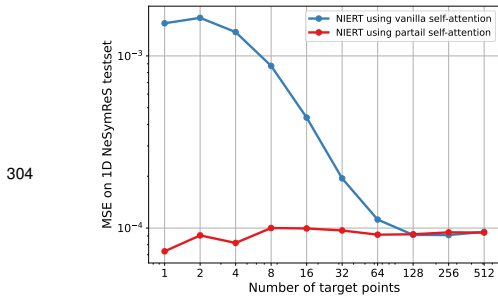
284 **The effects of partial self-attention:** For a specific interpolation task, the interpolation function is
 285 determined by the observed points only. Thus, an idealized encoding of observed points should not
 286 be affected by target points. To investigate the affects of target points, we evaluated NIERT on the
 287 test sets with various number of target points. Here, we compared two variants of NIERT, one with
 288 partial self-attention, and the other with vanilla self-attention. Both of these two variants were trained
 289 using the same training sets (the number of target points varies within [206, 246]).

290 As illustrated by Figure 6, the variant with vanilla self-attention shows poor performance for the tasks
 291 with few target points, say less than 64 target points. In contrast, the variant with partial self-attention
 292 always performs stably without significant changes of accuracy.

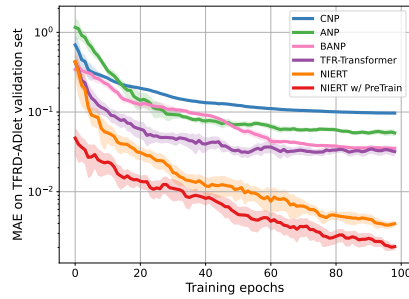
293 The results clearly demonstrate that the partial self-attention mechanism allows NIERT to be free
 294 from the unexpected affects by the target points.

295 **The effects of pre-training technique:** To investigate the effects of the pre-training technique, we
 296 show in Figure 7 the training process of two versions of NIERT, one without pre-training technique,
 297 and the other enhanced with pre-training. As depicted by the figure, even at the first epoch, the
 298 pre-trained NIERT shows a sufficiently high interpolation accuracy, which is comparable with the
 299 fully-trained BANP and TFR-transformer. Moreover, the performance of the pre-trained NIERT
 300 improves in roughly the same convergence speed to the original NIERT. At the final epoch, the
 301 pre-trained NIERT decreases the interpolation error to be nearly half of that of the original NIERT.

302 These results clearly suggest that the experience learned by NIERT from the interpolation task in one
 303 application field has potential to be transferred to the interpolation tasks in other application fields.



304
 305 Figure 6: The robustness of NIERT to the number of target points. Here, two variants of NIERT are trained on 1D NeSymReS dataset



306 Figure 7: The convergence of NIERT, NIERT with pre-training, and the existing approaches. Here, models are trained on TFRD-ADlet dataset

306 5 Discussion and conclusion

307 We present in the study an accurate approach to numerical interpolation for scattered data. The
 308 specific features of our NIERT approach are highlighted by the full exploitation of the correlation
 309 between observed points and target points through unifying scattered data representation. At the same
 310 time, the use of partial self-attention mechanism can effectively avoid the interference of target points
 311 onto the observed points. The enhancement with pre-training technique is another special feature
 312 of NIERT. The advantages of NIERT in interpolation accuracy have been clearly demonstrated by
 313 experimental results on both synthetic and real-world datasets.

314 The current version of NIERT has a computational complexity of $O(n(m+n))$, thus cannot handle
 315 the interpolation tasks with extremely large amounts of observed points due to the limitations of GPU
 316 memory size. Compared with the lightweight traditional methods, our NIERT approach has a much
 317 larger model with expensive computation to learn complex function distribution, which limits its
 318 application in cost sensitive scenarios. How to reduce the memory requirement and computational
 319 cost is one of the future works. Additionally, it is interesting to combine NIERT and the traditional
 320 approaches based on basis functions to yield an approach with both high accuracy and interpretability.

321 We expect NIERT, with extensions and modifications, to greatly facilitate numerical interpolations in
 322 a wide range of engineering and science fields.

323 References

324 [1] Richard Franke and Gregory M Nielson. Scattered data interpolation and applications: A
 325 tutorial and survey. *Geometric Modeling*, pages 131–160, 1991.

326 [2] Xiaoqian Chen, Zhiqiang Gong, Xiaoyu Zhao, Weien Zhou, and Wen Yao. A Machine Learning
 327 Modelling Benchmark for Temperature Field Reconstruction of Heat-Source Systems. *arXiv*

- 328 *preprint arXiv:2108.08298*, 2021.
- 329 [3] Mathieu Lepot, Jean-Baptiste Aubin, and François HLR Clemens. Interpolation in time series:
330 An introductive overview of existing methods, their performance criteria and uncertainty
331 assessment. *Water*, 9(10):796, 2017.
- 332 [4] Satya Narayan Shukla and Benjamin M Marlin. Interpolation-prediction networks for irregularly
333 sampled time series. *arXiv preprint arXiv:1909.07782*, 2019.
- 334 [5] GR3543275 Liu. An overview on meshfree methods: for computational solid mechanics.
335 *International Journal of Computational Methods*, 13(05):1630001, 2016.
- 336 [6] Michael T Heath. *Scientific Computing: An Introductory Survey, Revised Second Edition*.
337 Society for Industrial and Applied Mathematics, revised 2nd edition edition, 2018.
- 338 [7] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray
339 Shanahan, Yee Whye Teh, Danilo Rezende, and S. M. Ali Eslami. Conditional Neural Processes.
340 In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on*
341 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1704–1713.
342 PMLR, 10–15 Jul 2018.
- 343 [8] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum,
344 Oriol Vinyals, and Yee Whye Teh. Attentive Neural Processes. In *International Conference on*
345 *Learning Representations*, 2019.
- 346 [9] Juho Lee, Yoonho Lee, Jungtaek Kim, Eunho Yang, Sung Ju Hwang, and Yee Whye Teh.
347 Bootstrapping neural processes. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan,
348 and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages
349 6606–6615. Curran Associates, Inc., 2020.
- 350 [10] Byung-Jun Lee, Seunghoon Hong, and Kee-Eung Kim. Residual neural processes. In *Pro-*
351 *ceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4545–4552,
352 2020.
- 353 [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
354 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Informa-*
355 *tion Processing Systems*, pages 5998–6008, 2017.
- 356 [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
357 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
358 2018.
- 359 [13] Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT Pre-Training of Image Transformers. *arXiv*
360 *preprint arXiv:2106.08254*, 2021.
- 361 [14] Charles A Hall and W Weston Meyer. Optimal error bounds for cubic spline interpolation.
362 *Journal of Approximation Theory*, 16(2):105–122, 1976.
- 363 [15] William J Gordon and James A Wixom. Shepard’s method of “metric interpolation” to bivariate
364 and multivariate interpolation. *Mathematics of computation*, 32(141):253–264, 1978.
- 365 [16] Hans Wackernagel. *Multivariate geostatistics: an introduction with applications*. Springer
366 Science & Business Media, 2003.
- 367 [17] Michael JD Powell. Radial basis functions for multivariable interpolation: a review. *Algorithms*
368 *for approximation*, 1987.
- 369 [18] Bengt Fornberg and Julia Zuev. The Runge phenomenon and spatially variable shape parameters
370 in RBF interpolation. *Computers & Mathematics with Applications*, 54(3):379–398, 2007.
- 371 [19] Qiqi Wang, Parviz Moin, and Gianluca Iaccarino. A high order multivariate approximation
372 scheme for scattered data sets. *Journal of Computational Physics*, 229(18):6343–6361, 2010.

- 373 [20] Ying Tai, Jian Yang, and Xiaoming Liu. Image super-resolution via deep recursive residual
374 network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
375 pages 3147–3155, 2017.
- 376 [21] Juncheng Li, Zehua Pei, and Tiejong Zeng. From Beginner to Master: A Survey for Deep
377 Learning-based Single-Image Super-Resolution. *arXiv preprint arXiv:2109.14335*, 2021.
- 378 [22] Satya Narayan Shukla and Benjamin Marlin. Multi-Time Attention Networks for Irregularly
379 Sampled Time Series. In *International Conference on Learning Representations*, 2021.
- 380 [23] Noah Weber, Janez Starc, Arpit Mittal, Roi Blanco, and Lluís Màrquez. Optimizing over a
381 bayesian last layer. In *NeurIPS workshop on Bayesian Deep Learning*, 2018.
- 382 [24] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
383 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
384 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
385 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
386 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
387 Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In H. Larochelle,
388 M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information
389 Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- 390 [25] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
391 autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.
- 392 [26] Luca Biggio, Tommaso Bendinelli, Alexander Neitz, Aurelien Lucchi, and Giambattista Parascandolo.
393 Neural Symbolic Regression that Scales. In *International Conference on Machine
394 Learning*, pages 936–945. PMLR, 2021.
- 395 [27] Mojtaba Valipour, Bowen You, Maysum Panju, and Ali Ghodsi. SymbolicGPT: A Generative
396 Transformer Model for Symbolic Regression. *arXiv preprint arXiv:2106.14131*, 2021.
- 397 [28] Guillaume Lample and François Charton. Deep learning for symbolic mathematics. *arXiv
398 preprint arXiv:1912.01412*, 2019.
- 399 [29] Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. Predicting in-
400 hospital mortality of icu patients: The physionet/computing in cardiology challenge 2012. In
401 *2012 Computing in Cardiology*, pages 245–248. IEEE, 2012.
- 402 [30] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David
403 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.
404 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew
405 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W.
406 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.
407 Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul
408 van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific
409 Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- 410 [31] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation
411 of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*,
412 2014.
- 413 [32] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural Ordinary
414 Differential Equations. *Advances in Neural Information Processing Systems*, 2018.
- 415 [33] Yulia Rubanova, Ricky T. Q. Chen, and David K Duvenaud. Latent Ordinary Differential
416 Equations for Irregularly-Sampled Time Series. In H. Wallach, H. Larochelle, A. Beygelzimer,
417 F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing
418 Systems*, volume 32. Curran Associates, Inc., 2019.

419 **Checklist**

- 420 1. For all authors...
- 421 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
422 contributions and scope? [Yes]
- 423 (b) Did you describe the limitations of your work? [Yes] Limitations of our work are
424 briefly described in Section 5.
- 425 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 426 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
427 them? [Yes]
- 428 2. If you are including theoretical results...
- 429 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 430 (b) Did you include complete proofs of all theoretical results? [N/A]
- 431 3. If you ran experiments...
- 432 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
433 perimental results (either in the supplemental material or as a URL)? [Yes] In the
434 Supplementary.
- 435 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
436 were chosen)? [Yes] All the training details are in the Supplementary.
- 437 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
438 iments multiple times)? [Yes] We reported error bars in experiments on PhysioNet
439 dataset. See Table 3.
- 440 (d) Did you include the total amount of compute and the type of resources used (e.g., type
441 of GPUs, internal cluster, or cloud provider)? [Yes] Described in the Supplementary.
- 442 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 443 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 444 (b) Did you mention the license of the assets? [Yes] Mentioned in the Supplementary.
- 445 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 446 (d) Did you discuss whether and how consent was obtained from people whose data you're
447 using/curating? [No] The datasets used in our work are all open-sourced by their
448 creators.
- 449 (e) Did you discuss whether the data you are using/curating contains personally identifiable
450 information or offensive content? [No] The data we are using doesn't contain any
451 personally identifiable information or offensive content.
- 452 5. If you used crowdsourcing or conducted research with human subjects...
- 453 (a) Did you include the full text of instructions given to participants and screenshots, if
454 applicable? [N/A]
- 455 (b) Did you describe any potential participant risks, with links to Institutional Review
456 Board (IRB) approvals, if applicable? [N/A]
- 457 (c) Did you include the estimated hourly wage paid to participants and the total amount
458 spent on participant compensation? [N/A]