



23 temporal and spatial dimensions. The ST-LSTM-SA model not only predicts the ocean sound  
24 velocity field (SVF) accurately, but also provides valuable insights for spatiotemporal  
25 prediction of other oceanic environmental variables.

26 **Key words:** sound velocity field, spatiotemporal prediction, deep learning, self-attention

27 <https://doi.org/10.1007/s00376-024-3219-6>

28 **Article Highlights:**

- 29 ● A prediction model for 3D ocean sound velocity fields was developed based on deep  
30 learning.
- 31 ● Employing transfer learning, the ST-LSTM-SA is initially trained on reanalysis data and  
32 further refined on in-situ analysis data.
- 33 ● ST-LSTM-SA shows promising prediction ability by effectively capturing the spatial and  
34 temporal variability of sound speed.

35

36

37

## 38 1. Introduction

39 Sound waves are the main medium of underwater information transmission and have  
40 various applications in marine engineering, ocean navigation positioning and underwater  
41 communication (Akyildiz et al., 2005; Stojanovic and Preisig, 2009). In order to study these  
42 applications, it is essential to obtain accurate marine sound environmental parameters. Among  
43 them, sound velocity in seawater is one of the parameters that determines the sound propagation  
44 characteristics (Kinsler et al., 2000; Heidemann et al., 2012). Sound velocity in seawater is a  
45 function of seawater temperature, salinity and pressure, among which the temperature change  
46 has the most significant effect on the sound velocity (Chen and Millero, 1977; Mackenzie,  
47 1981). Therefore, sound velocity varies with the ocean dynamic environment in both the time  
48 and spatial domains. Due to the vertical stratification of the ocean environment, which in turn  
49 makes the sound velocity exhibit a vertically layered structure (Kinsler et al., 2000). In  
50 addition, short-term and long-term physical processes in the ocean, such as waves, internal  
51 waves, currents and seasonal changes, can alter the marine environment. The superposition of  
52 these different periodic physical processes results in complex temporal and spatial variations  
53 in sound velocity (Storto et al., 2020).

54 In current marine research, real-time sound velocity information is predominantly derived  
55 from in-situ measurements of sound velocity profiles (SVPs), which capture the variation of  
56 sound velocity from the water surface to the seabed (Liu et al., 2023). However, the ocean SVF  
57 offers a more comprehensive description of sound velocity distribution in three-dimensional  
58 space, which provides a refined representation of the spatial variations of sound velocity. The  
59 construction of real-time SVFs is often challenging due to limited observational methods (Dai  
60 et al., 2019; Wang et al., 2020). Traditional offshore measurements provide only sparse point-  
61 by-point SVPs, which are costly and inefficient to collect frequently. With the development of  
62 multiple technology, methods that rely on raw data to predict and invert the sound velocity  
63 have been widely studied in recent decades.

64 Ocean acoustic tomography, systematically introduced by Munk and Wunsch (1979),  
65 plays a vital role in marine research, which has paved the way for the development of various  
66 SVP inversion methods, including matched acoustic peak arrivals (Skarsoulis et al., 1996) and  
67 matched field inversion methods (Tolstoy et al., 1991; Goncharov et al., 1993). Kalman  
68 filtering is an optimization algorithm for state estimation and it has been shown to be impactful  
69 in ocean forecasting problems (Candy and Sullivan, 1993; Carrière et al., 2009). Compressive  
70 sensing (CS) in acoustics has garnered significant attention as an emerging technology in the  
71 past decade (Gerstoft et al., 2018). Unlike conventional SVP inversion methods, the CS  
72 inversion method effectively estimates fine-scale SVPs through sparse representation using a  
73 limited number of SVPs (Bianco and Gerstoft, 2016; Choo and Seong, 2018). Furthermore,  
74 machine learning has emerged as an effective way to tackle challenges in marine science,  
75 providing fresh avenues for employing data-driven methodologies to make predictions about  
76 marine environment (Park and Kennedy, 1996; Jain and Ali, 2006; Chen et al., 2016; Huang et  
77 al., 2021). Specifically, a Convolutional Long Short-Term Memory (ConvLSTM) model based  
78 on deep learning has been applied into SVP prediction over a three-dimensional sea area, with  
79 an average prediction error of less than  $1.7 \text{ m s}^{-1}$  (Li and Zhai, 2022).

80 Over the past four decades, researchers have extensively investigated various methods for  
81 ocean sound velocity inversion and prediction. Due to the intricate nature of the ocean  
82 environment, accurately predicting the ocean SVFs still poses a significant challenge.  
83 Traditionally, approaches for spatiotemporal prediction of marine environmental variables rely  
84 on ocean numerical simulations, which suffer from significant computational demands, leading  
85 to inefficiency in prediction. In fact, the time series of observed ocean data already contains  
86 valuable information regarding the internal dynamics and external drivers of the ocean  
87 (Espeholt et al., 2022; Shao et al., 2021). Deep learning models, which have the capability to  
88 learn from large datasets, can effectively extract the intrinsic characteristics and physical laws  
89 inherent in the data (LeCun et al., 2015). As a highly popular and influential technique, deep  
90 learning has been successfully applied in various marine prediction researches (Shao et al.,  
91 2021; Xiao et al., 2019; Ham et al., 2019; Andersson et al., 2021).

92 In this research, we propose a new spatiotemporal prediction model (ST-LSTM-SA) for  
93 ocean SVFs from a data-driven perspective. Our model combines deep artificial neural  
94 networks, including convolutional operations, recurrent neural networks, and self-attention  
95 mechanisms, to effectively capture the spatiotemporal variability of sound velocity and enable  
96 end-to-end prediction. The model employs an encoding-forecasting network structure that  
97 directly outputs future SVFs based on historical observation sequences. During model training,  
98 we employ transfer learning by firstly training the model using reanalysis datasets, followed  
99 by fine-tuning with the in-situ analysis data to obtain the final prediction model. In terms of  
100 accuracy, our model outperforms ANN, LSTM, ConvLSTM and ST-LSTM models,  
101 demonstrating superior performance across multiple evaluation metrics, and exhibits enhanced  
102 stability in predicting both temporal and spatial dimensions.

## 103 **2. Data and Data Preprocessing**

### 104 2.1. Data

105 The reanalysis dataset is a continuously integrated dataset created by merging  
106 observational data with advanced numerical modeling and assimilation techniques (Cummings  
107 and Smedstad, 2013). The reanalysis dataset used in this study is the Simple Ocean Data  
108 Assimilation System version 2.24, SODA2.24 (Giese and Ray, 2011). This dataset covers the  
109 assimilation period of 1871-2008, with a spatial range spanning  $0.25^{\circ}\text{E}$  to  $0.25^{\circ}\text{W}$  and  $75.25^{\circ}\text{S}$   
110 to  $89.25^{\circ}\text{N}$ . It has a horizontal resolution of  $0.5^{\circ}\times 0.5^{\circ}$  and a monthly temporal resolution. The  
111 vertical resolution varies from 10 m in the surface layer to 250 m in the bottom layer, divided  
112 into a total of 40 unequally spaced vertical layers, available at  
113 <https://www2.atmos.umd.edu/~ocean/>.

114 The Array for Real-time Geostrophic Oceanography (Argo) program has significantly  
115 enhanced oceanic observations, has yielded over 2.5 million ocean profiles to date (Johnson et  
116 al., 2022). Starting with these raw observations, researchers have produced numerous in-situ  
117 analysis datasets through the application of statistical analyses, optimal interpolation processes,  
118 and quality control techniques. (Zhang et al., 2022; Good et al., 2013; Gaillard et al., 2016). In  
119 this study, we utilize the Global Gridded Argo Dataset Based on Gradient-Dependent Optimal  
120 Interpolation (GDCSM\_Argo) (Zhang et al., 2022), covering so far the time range from January  
121 2004 to September 2022 with a monthly temporal resolution. The dataset encompasses the

122 entire global ocean with a horizontal resolution of  $1^\circ \times 1^\circ$  and consists of 58 unequally spaced  
 123 vertical layers. The vertical resolution ranges from 5m in the surface layer to 100m in the  
 124 bottom layer, available at <ftp://data.argo.org.cn/pub/ARGO/GDCSM/>.

125 The SODA2.2.4 dataset utilizes a simpler assimilation method and has limited early ocean  
 126 observation data. In contrast, the GDCSM\_Argo dataset is derived from Argo buoy  
 127 observations, providing an objective representation of the ocean interior. To effectively  
 128 leverage both datasets, we conducted vertical interpolation on the SODA2.2.4 dataset using  
 129 cubic spline interpolation, aligning it with the 58 layers of the GDCSM\_Argo dataset.  
 130 Afterward, the pre-training is conducted using the reanalysis dataset, delineated into training,  
 131 validation, and test sets covering the time spans of 1871-1980, 1981-1994, and 1995-2008,  
 132 respectively. Subsequent to this, the model undergoes additional training utilizing the  
 133 GDCSM\_Argo dataset. The training set encompasses the time frame from 2004 to 2018, while  
 134 the test set covers the period of 2019-2022. During this phase, it is noteworthy that there are  
 135 no alterations made to the hyperparameters of the model (Ham et al., 2019; Pan and Yang,  
 136 2010).

## 137 2.2. Data preprocessing

138 1. Data clipping. The dataset is initially cropped to extract the data within the study area  
 139 range, which spans from  $15^\circ\text{S}$  to  $15^\circ\text{N}$  and  $150^\circ\text{W}$  to  $180^\circ\text{W}$ . The study area is shown in Fig.  
 140 1 based on ETOPO1 bathymetric model (Amante and Eakins, 2009). This area is situated in  
 141 the central region of the Pacific Ocean, known for its dynamic climate change and ocean  
 142 environment.

143 2. Calculate sound velocity. Reanalysis and in-situ analysis data commonly include  
 144 variables such as seawater temperature and salinity, which allow for the calculation of sound  
 145 velocity using empirical equations. To determine sound velocity at each location, the water  
 146 depth values in the vertical direction were converted from pressure using the pressure-to-depth  
 147 conversion method proposed by Saunders (1981). Subsequently, the Del Grosso empirical  
 148 equation (Del Grosso, 1974) for sound velocity was used to calculate the sound velocity  
 149 information.

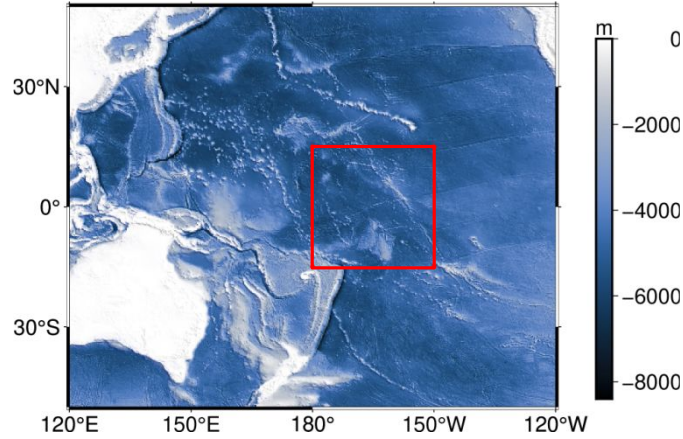
150 3. Data normalization. It involves linearly transforming input data to ensure they are  
 151 distributed within a specific range. This process helps balance the weights between different  
 152 features and enhances both the training effectiveness and generalization capability of the  
 153 model. In our study, we employed the maximum-minimum normalization operation, which  
 154 scaled all training data to the range of  $[0, 1]$ . The calculation procedure for this normalization  
 155 is as follows:

$$x^* = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad (1)$$

156 where  $x$  denotes the sample data,  $x_{\max}$  and  $x_{\min}$  are the maximum and minimum values of the  
 157 sample data,  $x^*$  represents the normalized sample data.

158 4. Slide sampling. We performed slide sampling on the normalized SVF data using a  
 159 window size of 15 and a step size of 1. Each sample consists of a sequence of 15 consecutive  
 160 monthly SVFs. Since the SVFs represent monthly averaged data, we utilize a 12-months input

161 series with an annual cycle to predict the SVFs for the subsequent 3-months, thereby achieving  
 162 seasonal forecasting.



163  
 164 **Fig. 1.** The bathymetric conditions obtained from ETOPO1 in the central Pacific Ocean, and  
 165 the study area (150°W-180°W, 15°S-15°N) is delineated by the red box.

### 166 3. Methodology

#### 167 3.1. Problem Definition

168 In terms of spatial representation, the ocean SVF comprises three dimensions. A three-  
 169 dimensional grid can be used to represent the spatial location of the sea, wherein each grid cell  
 170 contains time-dependent sound velocity information. By converting this grid into a tensor, we  
 171 can express the ocean SVF at a certain time as a three-dimensional tensor  $\mathbf{X}_t \in \mathbb{R}^{(M \times N \times D)}$ , with  
 172  $M$ ,  $N$  and  $D$  denoting longitude, latitude, and water depth, respectively.

173 Under the action of complex ocean dynamics processes, the ocean SVFs has obvious time  
 174 evolution characteristics. Therefore, the prediction problem of SVFs can be regarded as a  
 175 nonlinear time series prediction problem (Li and Zhai, 2022). Specifically, by leveraging  
 176 previously observed SVFs series within a given ocean area, we can forecast SVFs for future  
 177 time intervals with prediction models. Consequently, the temporal prediction problem for the  
 178 ocean SVFs is defined as follows: constructing a time series  $(\mathbf{X}_{t-n+1}, \mathbf{X}_{t-n+2}, \dots, \mathbf{X}_t)$  based on  $n$   
 179 consecutive past SVFs observations to predict the most probable SVFs  $(\hat{\mathbf{X}}_{t+1}, \dots, \hat{\mathbf{X}}_{t+k})$  for the  
 180 future time range  $(t+1, \dots, t+k)$  as expressed by the following equation:

$$\hat{\mathbf{X}}_{t+1}, \dots, \hat{\mathbf{X}}_{t+k} = f_{\theta}(\mathbf{X}_{t+1}, \dots, \mathbf{X}_{t+k} | \mathbf{X}_{t-n+1}, \mathbf{X}_{t-n+2}, \dots, \mathbf{X}_t) \quad (2)$$

181 where  $f$  denotes the spatiotemporal prediction model and  $\theta$  denotes the parameter that is  
 182 gradually optimized during the training process.

#### 183 3.2. Basic Deep Learning Models

184 Long Short-Term Memory (LSTM) is a unique recurrent neural network used for time  
 185 series problems. It excels at solving long-term dependencies between time series and is  
 186 applicable to the sound velocity time series prediction problem (Bengio et al., 1994; Hochreiter  
 187 and Schmidhuber, 1997). However, when it comes to spatiotemporal prediction, the fully  
 188 connected LSTM networks often struggle to capture spatial features effectively. To overcome

189 this limitation, Shi et al. (2015) introduces the ConvLSTM neural network, which combines  
 190 convolutional operations with LSTM and has shown success in precipitation nowcasting. The  
 191 ConvLSTM network incorporates convolutional operations in both input-to-state and state-to-  
 192 state transitions, allowing for the extraction of spatial features while capturing the dynamic  
 193 changes of the sequence. In a pioneering study by Li and Zhai (2022), ConvLSTM was applied  
 194 for the first time to predict SVPs. The experimental results revealed that ConvLSTM  
 195 outperformed LSTM, providing prediction results that closely aligned with the actual data.

196 ConvLSTM shares a similar internal structure with LSTM, featuring three gates within  
 197 each cell: the input gate  $i_t$ , forgetting gate  $f_t$ , and output gate  $o_t$ . The forgetting gate determines  
 198 the extent to which the previous memory state  $C_{t-1}$  is forgotten, while the input gate controls the  
 199 degree to which the current memory state  $C_t$  is updated. The output gate regulates the influence  
 200 of the current memory state  $C_t$  on the output hidden state  $H_t$ . Figure 2 provides a visual  
 201 representation of the internal structure of the ConvLSTM cell. The key equations for  
 202 ConvLSTM are as follows:

$$i_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \mathbf{e} C_{t-1} + b_i) \quad (3)$$

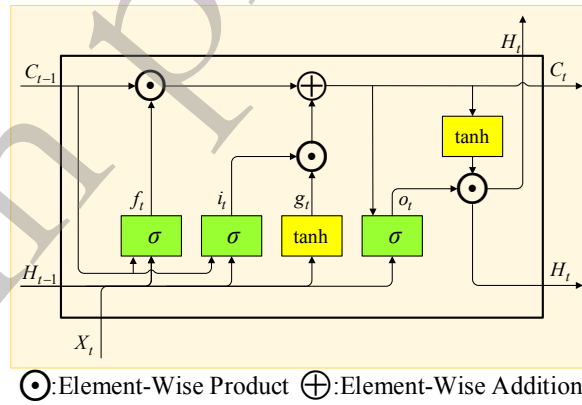
$$f_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \mathbf{e} C_{t-1} + b_f) \quad (4)$$

$$C_t = f_t \mathbf{e} C_{t-1} + i_t \mathbf{e} \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + b_c) \quad (5)$$

$$o_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \mathbf{e} C_t + b_o) \quad (6)$$

$$H_t = o_t \mathbf{e} \tanh(C_t) \quad (7)$$

203 Where  $X_t$  represents the input at the current time step,  $W$  and  $b$  are the weight and bias  
 204 coefficients that are continuously updated during model training. The \* symbol denotes the  
 205 convolution operation,  $\mathbf{e}$  represents the Hadamard (element-wise) operation and  $\sigma$  is the  
 206 sigmoid activation function.



207  
 208 **Fig. 2.** A demonstration of ConvLSTM cell structure.

209 The ConvLSTM network is a significant advancement in spatiotemporal prediction  
 210 research, considering the spatial correlation and temporal variation of data. It forms the basis  
 211 for further studies in this field. An enhanced variant of ConvLSTM, known as the PredRNN  
 212 network (Wang et al., 2017, 2022), further improves the internal structure to enhance  
 213 spatiotemporal prediction capabilities. PredRNN introduces a new spatiotemporal LSTM (ST-  
 214 LSTM) cell, which consists of two memory states: temporal memory state  $C_t'$  and  
 215 spatiotemporal memory state  $M_t'$ . In the ST-LSTM cell,  $C_t'$  is transmitted within adjacent time

216 steps on the same layer, while  $M_t^l$  is initially transmitted within the same time step, reaching  
 217 the top layer at the same moment after passing through the bottom layer at the next moment.  
 218 This transmission process is illustrated in Fig. 3. This unique method of memory state transfer  
 219 enables the memory state at the bottom level to depend on both the temporal memory state of  
 220 the previous moment at same layer and the spatiotemporal memory state from higher layer at  
 221 historical moments. Consequently, it enhances the interrelation of memory states across  
 222 different spatial levels. The specific equations of ST-LSTM cells are outlined below:

$$\mathbf{g}_t = \tanh(\mathbf{W}_{xg} * \mathbf{X}_t + \mathbf{W}_{hg} * \mathbf{H}_{t-1}^l + \mathbf{b}_g) \quad (8)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_{xi} * \mathbf{X}_t + \mathbf{W}_{hi} * \mathbf{H}_{t-1}^l + \mathbf{b}_i) \quad (9)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{xf} * \mathbf{X}_t + \mathbf{W}_{hf} * \mathbf{H}_{t-1}^l + \mathbf{b}_f) \quad (10)$$

$$\mathbf{C}_t^l = \mathbf{f}_t \mathbf{e} \mathbf{C}_{t-1}^l + \mathbf{i}_t \mathbf{e} \mathbf{g}_t \quad (11)$$

$$\mathbf{g}'_t = \tanh(\mathbf{W}'_{xg} * \mathbf{X}_t + \mathbf{W}'_{mg} * \mathbf{M}_{t-1}^l + \mathbf{b}'_g) \quad (12)$$

$$\mathbf{i}'_t = \sigma(\mathbf{W}'_{xi} * \mathbf{X}_t + \mathbf{W}'_{mi} * \mathbf{M}_{t-1}^l + \mathbf{b}'_i) \quad (13)$$

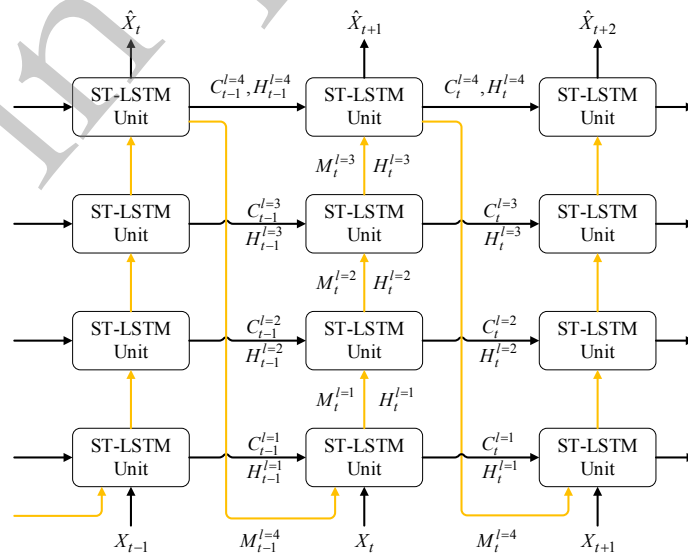
$$\mathbf{f}'_t = \sigma(\mathbf{W}'_{xf} * \mathbf{X}_t + \mathbf{W}'_{mf} * \mathbf{M}_{t-1}^l + \mathbf{b}'_f) \quad (14)$$

$$\mathbf{M}_t^l = \mathbf{f}'_t \mathbf{e} \mathbf{M}_{t-1}^l + \mathbf{i}'_t \mathbf{e} \mathbf{g}'_t \quad (15)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{xo} * \mathbf{X}_t + \mathbf{W}_{ho} * \mathbf{H}_{t-1}^l + \mathbf{W}_{mo} * \mathbf{M}_t^l + \mathbf{b}_o) \quad (16)$$

$$\mathbf{H}_t^l = \mathbf{o}_t \mathbf{e} \tanh(\mathbf{W}_{1 \times l} * [\mathbf{C}_t^l, \mathbf{M}_t^l]) \quad (17)$$

223 The memory state  $\mathbf{C}_t^l$  in the ST-LSTM unit follows the gate structures from the  
 224 ConvLSTM unit, which include the input gate  $\mathbf{i}_t$  and the forgetting gate  $\mathbf{f}_t$ . Additionally, a new  
 225 set of input gate  $\mathbf{i}'_t$  and forgetting gate  $\mathbf{f}'_t$  are introduced to control the information flow across  
 226 the memory state  $\mathbf{M}_t^l$ . The output gate  $\mathbf{o}_t$  is shared by the two memory units to facilitate memory  
 227 fusion for the storage state  $\mathbf{H}_t^l$ . The input modulation gates  $\mathbf{g}_t$  and  $\mathbf{g}'_t$  are involved in the  
 228 computation of memory states. The coefficients  $\mathbf{W}$  and  $\mathbf{b}$  represent the weight and bias terms in  
 229 the model.

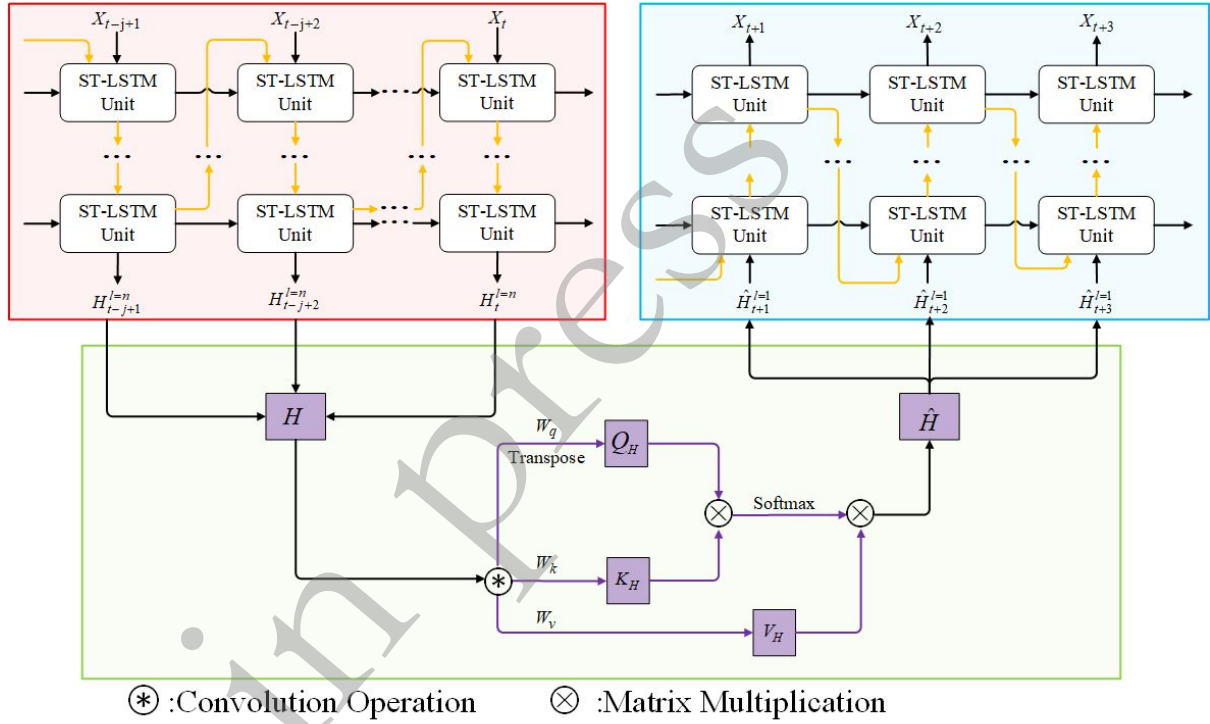




231 **Fig. 3.** The memory flow architecture of ST-LSTM, the orange arrows indicate that the  
 232 spatiotemporal memory state  $M_t^l$  is propagated in a zigzag pattern throughout the network, and  
 233 the black arrows denote the temporal memory state transition paths of  $C_t^l$ .

### 234 3.3. New ST-LSTM-SA Model

235 Numerous studies have demonstrated that ST-LSTM is well-suited for addressing  
 236 spatiotemporal prediction problems. Building upon this, we propose a novel ST-LSTM-SA  
 237 prediction model for SVFs prediction. The architecture of our model follows the encoding-  
 238 forecasting structure commonly used in earlier studies (Shi et al., 2015, 2017). In our model,  
 239 we incorporate a self-attention mechanism (Vaswani et al., 2017) between the encoding module  
 240 and the forecasting module to address temporal dependence issues in the prediction process.  
 241 This self-attention mechanism dynamically adjusts the weights of the encoding module's  
 242 outputs at different time steps, enabling us to obtain optimal inputs for the forecasting module  
 243 at each time step. Figure 4 illustrates the structure of our proposed ST-LSTM-SA model.



244

245 **Fig. 4.** Structure of the ST-LSTM-SA network. The red block represents the encoding  
 246 module. The green block presents the attention mechanism module, which demonstrates the  
 247 operation principle of the self-attention mechanism. The blue block represents the forecasting  
 248 module.

249 For the definition of the SVFs prediction problem in Section 3.1, at the current time  
 250 step  $t$ , the model is able to predict the SVFs for the next  $k$  time steps based on  $j$  historical  
 251 observations. In the encoding module, highlighted in the red boxed area in the Fig. 4, the input  
 252 SVF sequence  $(X_{t-j+1}, X_{t-j+2}, \dots, X_t)$  is encoded by  $n$  layers ST-LSTM cells to output  $j$   
 253 hidden states  $(H_{t-j+1}, H_{t-j+2}, \dots, H_t)$ . The attention mechanism module corresponds to the green  
 254 boxed area in the Fig. 4, and this part first superimposes the results of the encoding module on

255 the channels to obtain  $H$ . Then, the query  $Q_H$ , the key  $K_H$  and the value  $V_H$  are obtained by  
 256 mapping to different feature spaces through convolution operations, and  $\{W_q, W_k, W_v\}$   
 257 represent the weight parameters of the  $1 \times 1$  convolution operation. The output  $\hat{H}$  after  
 258 attention weight assignment can be obtained by Eq. (18-19):

$$\alpha = \text{softmax}(Q_H^T K_H) \quad (18)$$

$$\hat{H} = \alpha V_H \quad (19)$$

259 where  $Q_H^T$  denotes the transpose of  $Q_H$ , softmax is the nonlinear activation function, and  $\alpha$   
 260 denotes the attention weight distribution located between  $[0,1]$ . The forecasting module,  
 261 depicted within the blue boxed area in the Fig. 4, consists of  $n$  layers of ST-LSTM units with  
 262 the same structure as the encoding module. The output of the attention mechanism module  
 263 serves as the input for the forecasting module. At the last layer, the forecasting module  
 264 generates the prediction results for the next  $k$  time steps.

### 265 3.4. Implementation Details

266 The experiment was conducted on a server with the following configuration: Windows  
 267 operating system, 5.10 GHz CPU, 16 GB RAM, and an RTX 3060Ti GPU. The experiment  
 268 utilized Python 3.8 and the PyTorch 1.11 machine learning framework with CUDA version  
 269 11.3 for efficient GPU acceleration. Four neural network models, namely ANN, LSTM,  
 270 ConvLSTM, and ST-LSTM, were selected and compared with the proposed ST-LSTM-SA  
 271 model. The purpose of this comparison was to validate the superior performance of the ST-  
 272 LSTM-SA algorithm.

273 During the training process, we employed the Adam optimizer (Kingma and Ba, 2014)  
 274 to optimize the models. Each iteration utilized a batch size of 4, and the initial learning rate  
 275 was set to 0.001. The learning rate decayed as the number of training sessions increased. To  
 276 prevent overfitting, regularization was applied to effectively regularize the models. This  
 277 regularization technique helped enhance the generalization ability of the models and avoid  
 278 excessive fitting to the training data. Table 1 presents the implementation details for each model.  
 279 The ANN model was constructed using three fully connected layers and employed the MSE  
 280 loss function. Due to its limitation in handling only one-dimensional data, the historical SVF  
 281 sequence needs to be transformed into a one-dimensional tensor for input. The LSTM model  
 282 followed the encoding-forecasting structure and consisted of four layers of LSTM units with a  
 283 loss function of MSE. Unlike the ANN model, the LSTM model preserves the time dimension,  
 284 and the shape of the input tensor is  $(B, T, M \times N \times D)$ . The ConvLSTM, ST-LSTM, and ST-  
 285 LSTM-SA models all adopted the encoding-forecasting structure. Each module in these models  
 286 consisted of two corresponding layers with a uniform hidden state and memory state of 256  
 287 channels. The convolutional kernel size was set to  $3 \times 3$ . The loss function employed for these  
 288 models was MSE loss. The initial input tensor shape for these three models is  $(B, T, M, N, D)$ .  
 289 Considering the vertical stratification of ocean sound speed, we sequentially arrange the sound  
 290 speed values of the nine consecutive bathymetry layers in a  $3 \times 3$  order, filling insufficient  
 291 spaces with 0 values. This process results in the input tensor shape of  
 292  $(B, T, M \times 3, N \times 3, \lfloor D/9 \rfloor)$ .

293

294

295

**Table 1.** Implementation details of models.

Models	Implementation details	Input shape
ANN	Layers=4, hidden_dim=1024	$(B, T \times M \times N \times D)$
LSTM	Encoder: layers=2, hidden_dim=1024 Decoder: layers=2, hidden_dim=1024	$(B, T, M \times N \times D)$
ConvLSTM	Encoder: layers=2, kernel size = (3,3), channels = [256,256] Decoder: layers=2, kernel size = (3,3), channels = [256,56]	$(B, T, M \times 3, N \times 3, \lfloor D/9 \rfloor)$
ST-LSTM	Same as ConvLSTM	$(B, T, M \times 3, N \times 3, \lfloor D/9 \rfloor)$
ST-LSTM-SA	Same as ConvLSTM	$(B, T, M \times 3, N \times 3, \lfloor D/9 \rfloor)$

## 296 3.5. Evaluation methods

297 To evaluate the performance of the different prediction models, we employed several  
 298 evaluation metrics: root mean square error (RMSE), mean absolute error (MAE), mean  
 299 absolute percentage error (MAPE), and coefficient of determination ( $R^2$ ). RMSE, MAE, and  
 300 MAPE provide insights into the magnitude of the errors between the predicted and true values,  
 301 with smaller values indicating better model performance.  $R^2$  measures the strength of  
 302 correlation between the predicted and true values, with a value ranging from 0 to 1. A value  
 303 closer to 1 indicates a stronger correlation and better model performance. These metrics are  
 304 calculated as follows, where  $\hat{y}$  and  $y$  are the true value and the predicted value,  $n$  represents  
 305 the number of values.

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}} \quad (20)$$

$$\text{MAE} = \frac{\sum_{i=1}^n |\hat{y}_i - y_i|}{n} \quad (21)$$

$$\text{MAPE} = \frac{100\%}{N} \sum_{i=1}^n \left| \frac{\hat{y}_i - y_i}{y_i} \right| \quad (22)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (\bar{y}_i - y_i)^2} \quad (23)$$

306 **4. Results and Discussion**

## 307 4.1. Overall accuracy evaluation

308 To assess the effectiveness of the proposed algorithm for SVFs prediction, we  
 309 conducted an analysis and evaluation of the experimental results obtained from the new ST-

310 LSTM-SA model and four other models. Table 2 presents an overview of the prediction results  
 311 from different models. The spatiotemporal prediction models (ConvLSTM, ST-LSTM and ST-  
 312 LSTM-SA), which incorporate convolutional operations, consistently outperform the  
 313 traditional ANN and LSTM models. The improved performance across all evaluation metrics  
 314 suggests that the effective extraction of spatial features enhances the accuracy of predicting the  
 315 three-dimensional structure of the ocean SVFs.

316 Among the three spatiotemporal prediction models examined in this study, the ST-  
 317 LSTM model slightly outperforms the ConvLSTM model. This outcome indicates that the  
 318 introduction of spatiotemporal memory units, which facilitate information exchange across  
 319 different layers, is crucial for achieving favorable performance. Furthermore, the ST-LSTM-  
 320 SA model demonstrates further improvement compared to the ST-LSTM model. This finding  
 321 indicates that the attention mechanism module effectively enhances the quality of information  
 322 transfer between the encoding module and the forecasting module. By assigning weights to the  
 323 historical SVF sequences in the encoding module, the attention mechanism module contributes  
 324 to more realistic predictions from the forecasting module.

325 **Table 2.** Overall evaluation indicators for the five models prediction results.

Models	RMSE	MAE	MAPE	R <sup>2</sup>
ANN	1.784	1.001	0.066	0.993
LSTM	1.806	1.030	0.068	0.993
ConvLSTM	1.507	0.825	0.055	0.995
ST-LSTM	1.454	0.793	0.052	0.995
ST-LSTM-SA	1.315	0.728	0.048	0.996

326 Table 3 displays the RMSE and MAPE values for each model's prediction results across  
 327 different time steps. It is evident that, except for the ANN model, the prediction performance  
 328 of all models deteriorates over time. This decline can be attributed to the accumulation of errors  
 329 in the recurrent neural networks used by the other models, whereas the multiple time steps  
 330 prediction of the ANN model is independent. Notably, the ST-LSTM-SA model consistently  
 331 outperforms the other models in terms of prediction accuracy. It achieves the lowest prediction  
 332 errors for the next three months, with reduced increases in prediction errors between adjacent  
 333 months. Figure 5 presents a statistical histogram of the prediction RMSE for sound velocity.  
 334 The ST-LSTM-SA model exhibits the smallest error statistics, with approximately 80% of the  
 335 predicted sound velocity values having the RMSE of less than  $1\text{ m s}^{-1}$ . This indicates the model's  
 336 strong spatial and temporal prediction capability, which consistently produces stable and  
 337 accurate predictions.

338

339

340

341

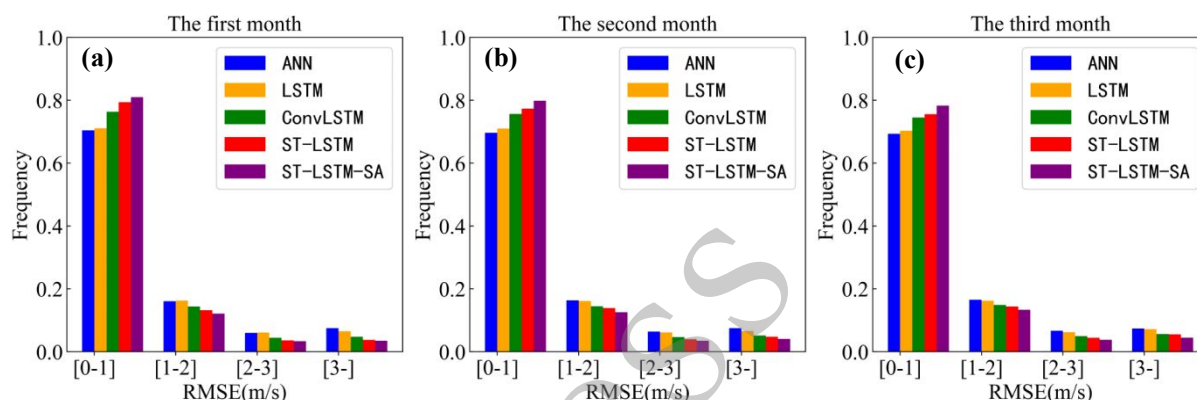
342

343

344 **Table 3.** Statistical of sound velocity prediction error for all test samples in future 3 months.

Mdels	RMSE			MAPE		
	1st month	2nd month	3rd month	1st month	2nd month	3rd month
ANN	1.788	1.789	1.775	0.066	0.066	0.066
LSTM	1.609	1.662	1.758	0.066	0.067	0.068
ConvLSTM	1.432	1.514	1.571	0.053	0.055	0.056
ST-LSTM	1.279	1.478	1.588	0.048	0.053	0.056
ST-LSTM-SA	1.211	1.329	1.399	0.045	0.048	0.051

345



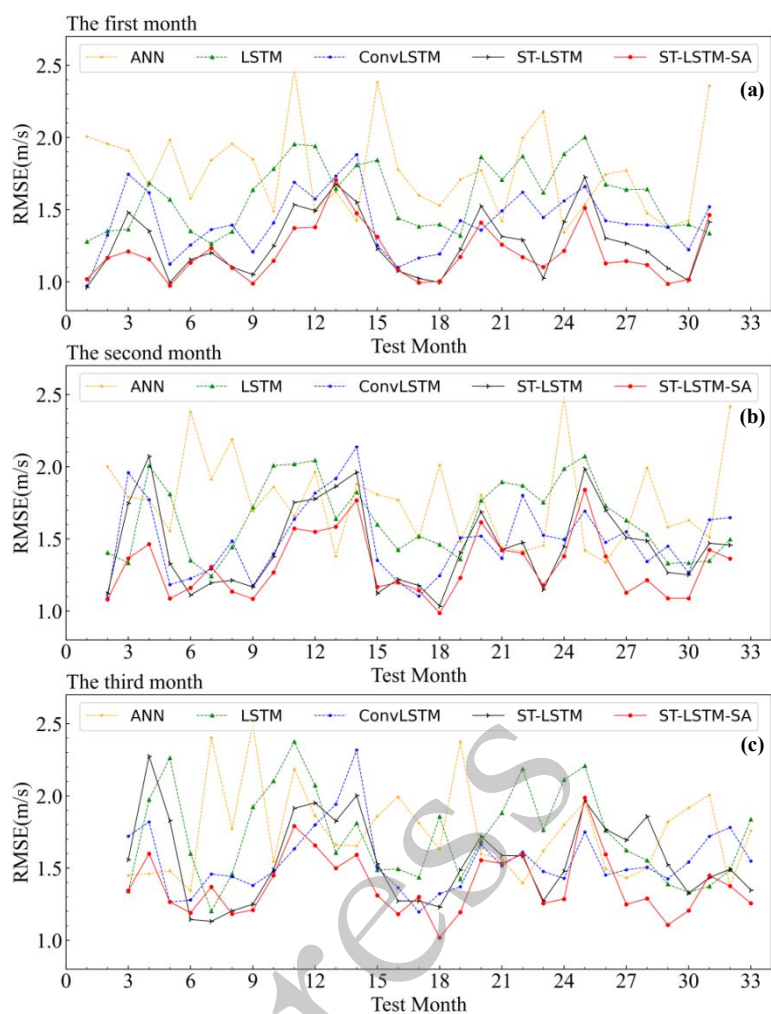
346

347 **Fig. 5.** Frequency distribution of prediction RMSE for all test samples in future 3 months: (a)  
348 the first month; (b) the second month; (c) the third month.

349 In Fig. 6, we present the RMSE of the prediction results at different time steps. For the  
350 prediction of the SVFs over the next 31 months, the error curves of the ANN and LSTM models  
351 exhibit more pronounced fluctuations. Conversely, the RMSE curves of the three  
352 spatiotemporal prediction models demonstrate consistent periodic patterns. Moreover, we  
353 observe smoother transitions between adjacent months, leading to a notable reduction in  
354 prediction errors. Notably, due to their similar network structures, the error curves of the ST-  
355 LSTM and ST-LSTM-SA models closely align with each other. However, in most cases, the  
356 ST-LSTM-SA model exhibits further improvement in prediction accuracy compared to ST-  
357 LSTM. This finding indicates that the network structure designed in this paper is more suitable  
358 for predicting the sound velocity field.

359

360



361

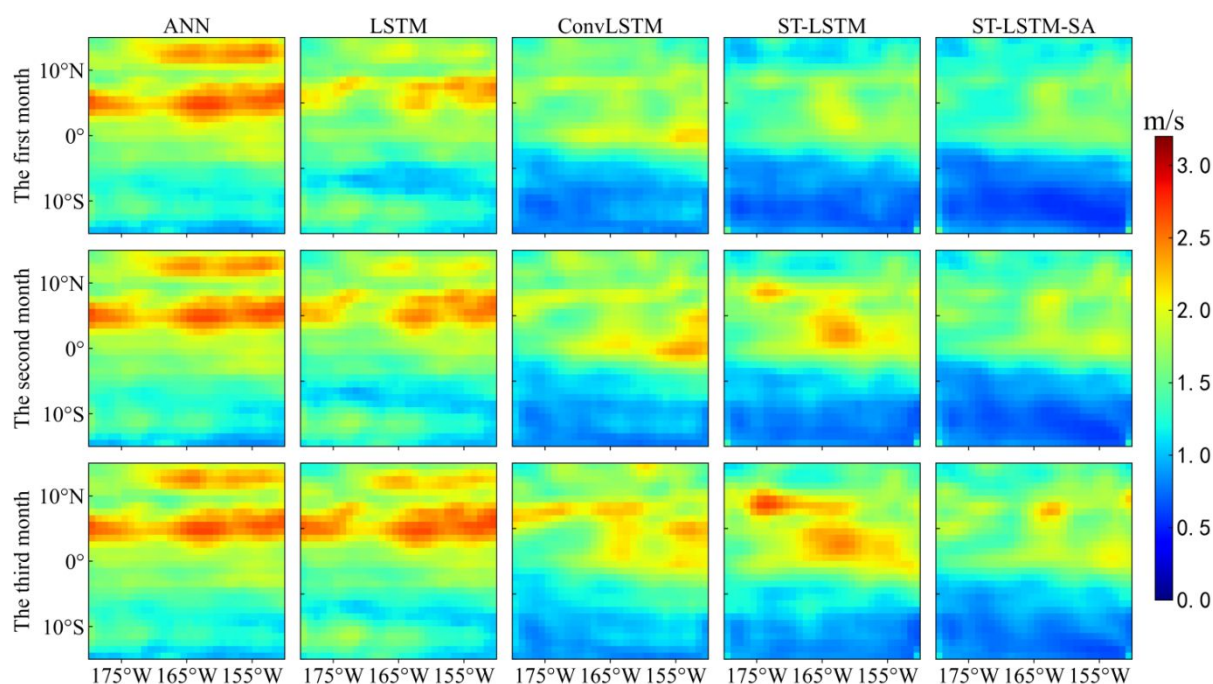
362 **Fig. 6.** The RMSE of different models versus the test dataset across different months: (a) the  
 363 first month; (b) the second month; (c) the third month.

364 4.2. Spatial predictive accuracy assessment

365 4.2.1. Horizontal direction analysis

366 To analyze the predictive accuracy of the models at various spatial locations, we delved  
 367 into their prediction results along two dimensions: the horizontal direction and the water depth  
 368 direction. Initially, we computed the RMSE of all water layers at each latitude and longitude  
 369 grid point for every model. Figure 7 presents a visual representation of the obtained analysis  
 370 outcomes.

371 The spatial distribution of errors reveals that the ANN and LSTM models exhibit  
 372 significant prediction biases across most locations. Interestingly, the spatial distribution of  
 373 prediction errors remains relatively consistent for the upcoming three months. However, the  
 374 spatiotemporal prediction models demonstrate noticeable improvements. This improvement  
 375 can be attributed to the fact that the ANN and LSTM models convert the three-dimensional  
 376 structure of the SVFs into a one-dimensional representation during training, resulting in a  
 377 considerable reduction in spatial correlation within the input data. Consequently, these models  
 378 tend to focus on individual locations rather than the entire SVF during the prediction process.



379

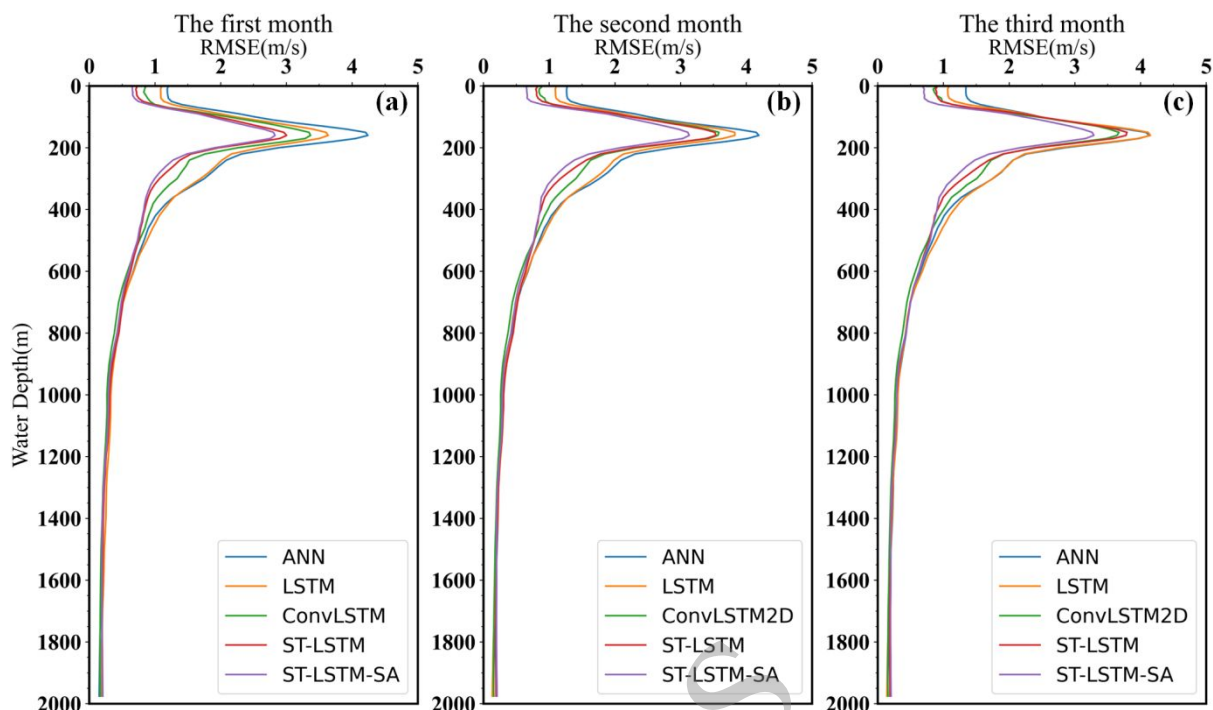
380 **Fig. 7.** Statistical predictions in horizontal direction for all test samples in 3 months  
 381 prediction. Each grid point represents the RMSE of all water layers at one location.

382 Examining the prediction results of the spatiotemporal prediction models for the  
 383 subsequent three months, we observe that the prediction error gradually expands over time. All  
 384 spatiotemporal models exhibit better performance during the first month of the forecast.  
 385 However, both the ConvLSTM and ST-LSTM models hardly maintain stable prediction  
 386 capabilities in the following second and third months. In contrast, the prediction results of the  
 387 ST-LSTM-SA model consistently maintain a balanced spatial distribution, with minimal  
 388 increases in error between adjacent months. This demonstrates the effectiveness of the ST-  
 389 LSTM-SA model in capturing both spatial and temporal variations in sound velocity. Moreover,  
 390 each module within the model plays a distinct role, making it well-suited for predicting the  
 391 marine sound velocity field.

#### 392 4.2.2. Depth direction analysis

393 In order to gain further insights into the models' prediction capabilities at different water  
 394 depths, we computed the RMSE of the prediction results for each water layer, as illustrated in  
 395 Fig. 8. It is evident that the prediction errors of all models exhibit a pattern of initially  
 396 increasing and then decreasing with increasing depth. In other words, the models perform  
 397 consistently and maintain stable prediction abilities in the surface layer and deep isothermal  
 398 layer, with the spatiotemporal prediction models achieving a prediction accuracy within  $1\text{ m s}^{-1}$ .  
 399 However, in the thermocline layer, the models' prediction accuracy fluctuates significantly.  
 400 Both the ANN and LSTM models reach maximum RMSE exceeding  $4\text{ m s}^{-1}$ , while the  
 401 spatiotemporal prediction model surpasses  $3\text{ m s}^{-1}$ . This indicates that the ocean environment  
 402 undergoes more pronounced changes in the thermocline layer. The prediction models struggle  
 403 to effectively capture the underlying patterns and mechanisms of these changes, as the marine

404 variables such as seawater temperature and salinity are influenced by light and complex  
 405 physical processes, leading to considerable uncertainties.



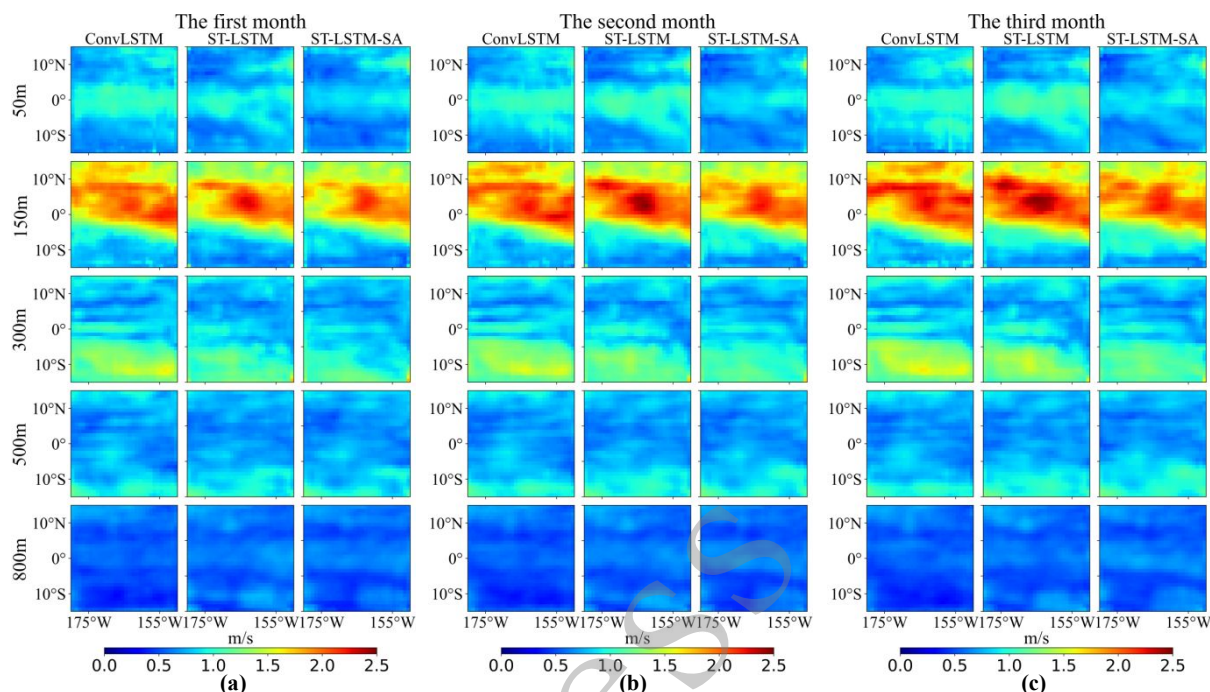
406  
 407 **Fig. 8.** Statistical predictions on depth direction for all test samples in 3 months prediction:  
 408 (a) the first month; (b) the second month; (c) the third month.

409 Although the prediction accuracy of each model in the thermocline layer falls short of  
 410 expectations, a comparison among the models reveals noteworthy findings. The errors of the  
 411 spatiotemporal prediction models, unlike those of the ANN and LSTM models, exhibit  
 412 convergence across different water depths. Notably, the ST-LSTM-SA model demonstrates  
 413 significant improvement in prediction accuracy for both the surface layer and thermocline layer.  
 414 This suggests that capturing the spatial characteristics of sound velocity is crucial in addressing  
 415 the prediction challenges associated with the SVFs. Furthermore, incorporating the attention  
 416 mechanism enhances not only the prediction accuracy but also the stability of the model across  
 417 future prediction time.

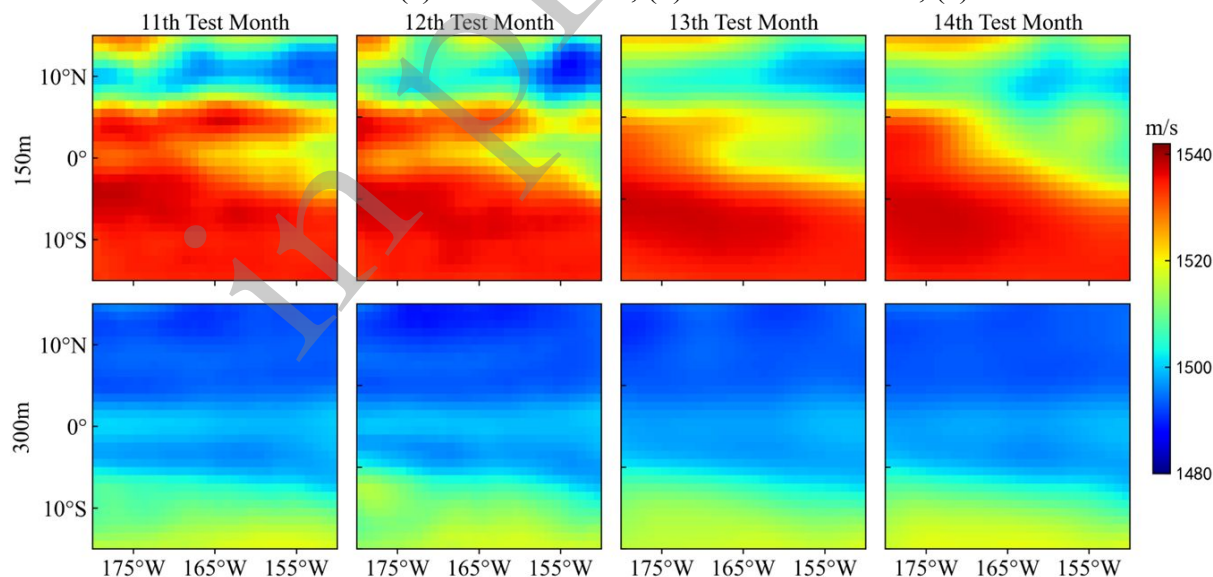
418 Figure 9 illustrates the horizontal slices of the sound velocity RMSE for the ConvLSTM,  
 419 ST-LSTM, and ST-LSTM-SA models at various depths: 50 m, 150 m, 300 m, 500 m, and 800  
 420 m. In the surface layer (50 m), the ST-LSTM-SA model demonstrates enhanced prediction  
 421 accuracy in the central region compared to the other two models. Moving deeper into the water  
 422 layers at 150 m and 300 m, notable spatial variations in sound velocity prediction errors are  
 423 observed. To showcase these differences, we select the 11th, 12th, 13th and 14th test months  
 424 from Fig. 6, where error fluctuations are more pronounced, and present their respective sound  
 425 velocity horizontal slices at 150 m and 300 m in Fig. 10. From the results, it becomes apparent  
 426 that sound velocity exhibits significant variability across most of the area north of 5°S at 150  
 427 m, while a smoother transition is observed in the southern area. At 300 m, a separation line in  
 428 sound velocity is still present, but the variation between adjacent months is relatively smoother.  
 429 Although we acknowledge the potential for doubting the accuracy of the dataset (Zhou et al.,  
 430 2023), it is important to note that the sharp fluctuations in sound speed primarily impact the



431 model's predictive ability for purely spatiotemporal prediction problems. These fluctuations  
 432 are closely linked to complex changes in the oceanic environment. As we delve deeper into the  
 433 water layers at 500 m and 800 m, seawater temperature gradually stabilizes, resulting in regular  
 434 changes in sound velocity. Consequently, the prediction abilities of the different models  
 435 become nearly indistinguishable.



436  
 437 **Fig. 9.** Horizontal slices of prediction RMSE for all test samples at 50m, 150m, 300m, 500m,  
 438 and 800m in future 3 months: (a) the first month; (b) the second month; (c) the third month.



439  
 440 **Fig. 10.** Horizontal slices of the sound velocity value at 150m and 300m for the 11th, 12th,  
 441 13th and 14th test months.

## 442 **5. Conclusions**

443 Traditionally, numerical ocean simulations are predominantly employed for predicting  
444 physical phenomena and internal information within the ocean. This study introduces a novel  
445 approach to marine SVFs prediction using the ST-LSTM-SA model, which leverages deep  
446 learning techniques. By treating the prediction of SVFs as a nonlinear time series prediction  
447 problem and adopting a data-driven approach, this method significantly enhances  
448 computational efficiency and reduces resource consumption. The ST-LSTM-SA model is  
449 designed to effectively integrate convolutional operations, LSTM, and self-attention  
450 mechanisms, allowing for the consideration of both spatial and temporal correlations in the  
451 SVF. This enables end-to-end prediction of the SVFs. During model training, transfer learning  
452 techniques are employed to train the model weights on different datasets. The SODA2.2.4  
453 reanalysis dataset assists in capturing simple variations in sound velocity over an extended time  
454 period, while the GDCSM\_Argo in-situ analysis data provides more realistic detailed  
455 characteristics of sound velocity, which further refines the model weights.

456 Through an analysis of the prediction results from January 2019 to September 2022, it  
457 is found that the ST-LSTM-SA model outperforms other models across all indicators,  
458 demonstrating better agreement with observed results. Temporally, the prediction results of the  
459 ST-LSTM-SA model exhibit stability over time, and the self-attention mechanism effectively  
460 handles long-term dependencies within the time series. Spatially, traditional ANN and LSTM  
461 models convert multi-dimensional data into one-dimensional data during input, disregarding  
462 the spatial and temporal correlations in the data, resulting in larger discrepancies in prediction  
463 accuracy across different locations. Conversely, the ST-LSTM-SA model demonstrates a more  
464 balanced spatial prediction capability, with prediction errors converging across different water  
465 depth layers.

466 Sound velocity in the ocean is influenced by a complex and dynamic environment,  
467 making it challenging to accurately describe and simulate its motion and underlying physical  
468 laws. In this study, we focus on investigating the spatiotemporal prediction of the sound  
469 velocity field. However, there is still room for further improvement in prediction accuracy. We  
470 plan to optimize the prediction model by refining its architecture and incorporating new feature  
471 data, which is expected to achieve better predictions of the sound velocity field in the future.

## 472 **Acknowledgments**

473 This work was supported by National Natural Science Foundation of China (42004030),  
474 Basic Scientific Fund for National Public Research Institutes of China (2022S03), Science and  
475 Technology Innovation Project (LSKJ202205102) Funded by Laoshan Laboratory, and  
476 National Key Research and Development Program of China (2020YFB0505805).

## 477 **Data Availability Statement**

478 We thank National Centers for Environmental Information/National Oceanic and  
479 Atmospheric Administration (NCEI/NOAA) for ETOPO1 surface topography data available  
480 at <https://www.ncei.noaa.gov/access/metadata/landing->

481 page/bin/iso?id=gov.noaa.ngdc.mgg.dem:316, IRI/LDEO Climate Data Library for the  
482 SODA version 2.2.4 data available at <https://www2.atmos.umd.edu/~ocean/>, and Shanghai  
483 Ocean University and China Argo Real-time Data Center for the GDCSM\_Argo gridded  
484 dataset available at <https://argo.ucsd.edu/data/argo-data-products/>.

## 485 **References**

486 Akyildiz, I. F., D. Pompili, and T. Melodia, 2005: Underwater acoustic sensor networks:  
487 research challenges. *Ad hoc networks*, **3**(3), 257-279,  
488 <https://doi.org/10.1016/j.adhoc.2005.01.004>

489 Amante, C., and B. W. Eakins, 2009: ETOPO1 1 Arc-Minute Global Relief Model: Procedures,  
490 Data Sources and Analysis. NOAA Technical Memorandum NESDIS NGDC-  
491 24[Dataset]. National Geophysical Data Center, NOAA,  
492 <http://dx.doi.org/10.7289/V5C8276M>

493 Andersson, T. R., J. S. Hosking, M. Pérez-Ortiz, B. Paige, A. Elliott, C. Russell, S. Law, D. C.  
494 Jones, J. Wilkinson, T. Phillips, J. Byrne, S. Tietsche, B. B. Sarojini, E. Blanchard-  
495 Wrigglesworth, Y. Aksenov, R. Downie, and E. Shuckburgh, 2021: Seasonal Arctic sea  
496 ice forecasting with probabilistic deep learning. *Nature Communications*, **12**(1), 5124,  
497 <https://doi.org/10.1038/s41467-021-25257-4>

498 Bengio, Y., P. Simard, and P. Frasconi, 1994: Learning long-term dependencies with gradient  
499 descent is difficult. *IEEE transactions on neural networks*, **5**(2), 157-166,  
500 <https://doi.org/10.1109/72.279181>

501 Bianco, M., and P. Gerstoft, 2016: Compressive acoustic sound speed profile estimation. *The*  
502 *Journal of the Acoustical Society of America*, **139**(3), EL90-EL94,  
503 <https://doi.org/10.1121/1.4943784>

504 Candy, J. V., and E. J. Sullivan, 1993: Sound velocity profile estimation: A system theoretic  
505 approach. *IEEE journal of oceanic engineering*, **18**(3), 240-252,  
506 <https://doi.org/10.1109/JOE.1993.236362>

507 Carrière, O., J. P. Hermand, and J. V. Candy, 2009: Inversion for time-evolving sound-speed  
508 field in a shallow ocean by ensemble Kalman filtering. *IEEE Journal of Oceanic*  
509 *Engineering*, **34**(4), 586-602, <https://doi.org/10.1109/JOE.2009.2033954>

510 Chen, C. T., and F. J. Millero, 1977: Speed of sound in seawater at high pressures. *The Journal*  
511 *of the Acoustical Society of America*, **62**(5), 1129-1135, <https://doi.org/10.1121/1.381646>

512 Chen, C., B. Lei, Y. L. Ma, and R. Duan, 2016: Investigating sound speed profile assimilation:  
513 An experiment in the Philippine Sea. *Ocean Engineering*, **124**, 135-140,  
514 <https://doi.org/10.1016/j.oceaneng.2016.07.062>

515 Choo, Y., and W. Seong, 2018: Compressive sound speed profile inversion using beamforming  
516 results. *Remote Sensing*, **10**(5), 704, <https://doi.org/10.3390/rs10050704>

517 Cummings, J. A., and O. M. Smedstad, 2013: Variational data assimilation for the global ocean.

- 518 Data Assimilation for Atmospheric, Oceanic and Hydrologic Applications (Vol. II), 303-  
519 343, [https://doi.org/10.1007/978-3-642-35088-7\\_13](https://doi.org/10.1007/978-3-642-35088-7_13)
- 520 Dai, M., Y. Li, and K. Yang, 2019: Joint inversion for sound speed field and moving source  
521 localization in shallow water. *Journal of Marine Science and Engineering*, **7**(9), 295,  
522 <https://doi.org/10.3390/jmse7090295>
- 523 Del Grosso, V. A., 1974: New equation for the speed of sound in natural waters (with  
524 comparisons to other equations). *The Journal of the Acoustical Society of America*, **56**(4),  
525 1084-1091, <https://doi.org/10.1121/1.1903388>
- 526 Espenholt, L., S. Agrawal, C. Sønderby, M. Kumar, J. Heek, C. Bromberg, C. Gazen, R. Carver,  
527 M. Andrychowicz, J. Hickey, A. Bell, and N. Kalchbrenner, 2022: Deep learning for  
528 twelve hour precipitation forecasts. *Nature communications*, **13**(1), 5145,  
529 <https://doi.org/10.1038/s41467-022-32483-x>
- 530 Gaillard, F., T. Reynaud, V. Thierry, N. Kolodziejczyk, and K. von Schuckmann, 2016: In  
531 situ-based reanalysis of the global ocean temperature and salinity with ISAS: Variability  
532 of the heat content and steric height. *Journal of Climate*, **29**(4), 1305-1323,  
533 <https://doi.org/10.1175/JCLI-D-15-0028.1>
- 534 Gerstoft, P., C. F. Mecklenbräuker, W. Seong, and M. Bianco, 2018: Introduction to  
535 compressive sensing in acoustics. *The Journal of the Acoustical Society of America*,  
536 **143**(6), 3731-3736, <https://doi.org/10.1121/1.5043089>
- 537 Giese, B. S., and S. Ray, 2011: El Niño variability in simple ocean data assimilation (SODA),  
538 1871–2008. *Journal of Geophysical Research: Oceans*, **116**(C2),  
539 <https://doi.org/10.1029/2010JC006695>
- 540 Goncharov, V. V., and A. G. Voronovich, 1993: An experiment on matched-field acoustic  
541 tomography with continuous wave signals in the Norway Sea. *The Journal of the*  
542 *Acoustical Society of America*, **93**(4), 1873-1881, <https://doi.org/10.1121/1.406702>
- 543 Good, S. A., M. J. Martin, and N. A. Rayner, 2013: EN4: Quality controlled ocean temperature  
544 and salinity profiles and monthly objective analyses with uncertainty estimates. *Journal*  
545 *of Geophysical Research: Oceans*, **118**(12), 6704-6716,  
546 <https://doi.org/10.1002/2013JC009067>
- 547 Ham, Y. G., J. H. Kim, and J. J. Luo, 2019: Deep learning for multi-year ENSO forecasts.  
548 *Nature*, **573**(7775), 568-572, <https://doi.org/10.1038/s41586-019-1559-7>
- 549 Heidemann, J., M. Stojanovic, and M. Zorzi, 2012: Underwater sensor networks: applications,  
550 advances and challenges. *Philosophical Transactions of the Royal Society A:*  
551 *Mathematical, Physical and Engineering Sciences*, **370**(1958), 158-175,  
552 <https://doi.org/10.1098/rsta.2011.0214>
- 553 Hochreiter, S., and J. Schmidhuber, 1997: Long short-term memory. *Neural computation*, **9**(8),  
554 1735-1780, <https://doi.org/10.1162/neco.1997.9.8.1735>
- 555 Huang, J., Y. Luo, J. Shi, X. Ma, Q. Q. Li, and Y. Y. Li, 2021: Rapid Modeling of the Sound

- 556 Speed Field in the South China Sea Based on a Comprehensive Optimal LM-BP Artificial  
557 Neural Network. *Journal of Marine Science and Engineering*, **9**(5), 488,  
558 <https://doi.org/10.3390/jmse9050488>
- 559 Jain, S., and M. M. Ali, 2006: Estimation of sound speed profiles using artificial neural  
560 networks. *IEEE Geoscience and Remote Sensing Letters*, **3**(4), 467-470,  
561 <https://doi.org/10.1109/LGRS.2006.876221>
- 562 Kingma, D. P., and J. Ba, 2014: Adam: a method for stochastic optimization. In: Proceedings  
563 of the 3rd International Conference on Learning Representations, ICLR 2015,  
564 <https://doi.org/10.48550/arXiv.1412.6980>
- 565 Kinsler, L. E., A. R. Frey, A. B. Coppens, and J. V. Sanders, 2000: Fundamentals of acoustics.  
566 3rd ed., John Wiley and Sons, 480pp.
- 567 LeCun, Y., Y. Bengio, and G. Hinton, 2015: Deep learning. *Nature*, **521**(7553), 436-444,  
568 <https://doi.org/10.1038/nature14539>
- 569 Li, B., and J. Zhai, 2022: A Novel Sound Speed Profile Prediction Method Based on the  
570 Convolutional Long-Short Term Memory Network. *Journal of Marine Science and  
571 Engineering*, **10**(5), 572, <https://doi.org/10.3390/jmse10050572>
- 572 Liu, Y., Y. Chen, Z. Meng, and W. Chen, 2023: Performance of single empirical orthogonal  
573 function regression method in global sound speed profile inversion and sound field  
574 prediction. *Applied Ocean Research*, **136**, 103598,  
575 <https://doi.org/10.1016/j.apor.2023.103598>
- 576 Mackenzie, K. V., 1981: Nine-term equation for sound speed in the oceans. *The Journal of the  
577 Acoustical Society of America*, **70**(3), 807-812, <https://doi.org/10.1121/1.386920>
- 578 Munk, W., and C. Wunsch, 1979: Ocean acoustic tomography: A scheme for large scale  
579 monitoring. *Deep Sea Research Part A. Oceanographic Research Papers*, **26**(2), 123-161,  
580 [https://doi.org/10.1016/0198-0149\(79\)90073-6](https://doi.org/10.1016/0198-0149(79)90073-6)
- 581 Pan, S. J., and Q. Yang, 2010: A survey on transfer learning. *IEEE Transactions on knowledge  
582 and data engineering*, **22**(10), 1345-1359, <https://doi.org/10.1109/TKDE.2009.191>
- 583 Park, J. C., and R. M. Kennedy, 1996: Remote sensing of ocean sound speed profiles by a  
584 perceptron neural network. *IEEE journal of oceanic engineering*, **21**(2), 216-224,  
585 <https://doi.org/10.1109/48.486796>
- 586 Johnson, G. C., S. Hosoda, S. R. Jayne, P. R. Oke, S. C. Riser, D. Roemmich, T. Suga, V.  
587 Thierry, S. E. Wijffels, and J. P. Xu, 2022: Argo—Two decades: Global oceanography,  
588 revolutionized. *Annual review of marine science*, **14**, 379-403,  
589 <https://doi.org/10.1146/annurev-marine-022521-102008>
- 590 Saunders, P. M., 1981: Practical conversion of pressure to depth. *Journal of Physical  
591 Oceanography*, **11**(4), 573-574, [https://doi.org/10.1175/1520-  
592 0485\(1981\)011%3C0573:PCOPTD%3E2.0.CO;2](https://doi.org/10.1175/1520-0485(1981)011%3C0573:PCOPTD%3E2.0.CO;2)
- 593 Shao, Q., W. Li, G. Han, G. Hou, S. Liu, Y. Gong, and P. Qu, 2021: A deep learning model for

- 594 forecasting sea surface height anomalies and temperatures in the South China Sea. *Journal*  
595 *of Geophysical Research: Oceans*, **126**(7), e2021JC017515,  
596 <https://doi.org/10.1029/2021JC017515>
- 597 Shi, X., Z. Chen, H. Wang, D. Y. Yeung, W. K. Wong, and W. C. Woo, 2015: Convolutional  
598 LSTM network: A machine learning approach for precipitation nowcasting. *Advances in*  
599 *neural information processing systems*, **28**, <https://doi.org/10.48550/arXiv.1506.04214>
- 600 Shi, X., Z. Gao, L. Lausen, H. Wang, D. Y. Yeung, W. K. Wong, and W. C. Woo, 2017: Deep  
601 learning for precipitation nowcasting: A benchmark and a new model. *Advances in neural*  
602 *information processing systems*, **30**, <https://doi.org/10.48550/arXiv.1706.03458>
- 603 Skarsoulis, E. K., G. A. Athanassoulis, and U. Send, 1996: Ocean acoustic tomography based  
604 on peak arrivals. *The Journal of the Acoustical Society of America*, **100**(2), 797-813,  
605 <https://doi.org/10.1121/1.416212>
- 606 Stojanovic, M., and J. Preisig, 2009: Underwater acoustic communication channels:  
607 Propagation models and statistical characterization. *IEEE communications magazine*,  
608 **47**(1), 84-89, <https://doi.org/10.1109/MCOM.2009.4752682>
- 609 Storto, A., S. Falchetti, P. Oddo, Y. M. Jiang, and A. Tesei, 2020: Assessing the impact of  
610 different ocean analysis schemes on oceanic and underwater acoustic predictions. *Journal*  
611 *of Geophysical Research: Oceans*, **125**(7), e2019JC015636,  
612 <https://doi.org/10.1029/2019JC015636>
- 613 Tolstoy, A., O. Diachok, and L. N. Frazer, 1991: Acoustic tomography via matched field  
614 processing. *The Journal of the Acoustical Society of America*, **89**(3), 1119-1127,  
615 <https://doi.org/10.1121/1.400647>
- 616 Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I.  
617 Polosukhin, 2017: Attention is all you need. *Advances in neural information processing*  
618 *systems*, **30**, <https://doi.org/10.48550/arXiv.1706.03762>
- 619 Wang, J., T. Xu, W. Nie, X. Yu, 2020: The construction of sound speed field based on back  
620 propagation neural network in the global ocean. *Marine Geodesy*, **43**(6), 621-642,  
621 <https://doi.org/10.1080/01490419.2020.1815912>
- 622 Wang, Y., M. Long, J. Wang, Z. Gao, and P. S. Yu, 2017: Predrnn: Recurrent neural networks  
623 for predictive learning using spatiotemporal lstms. *Advances in neural information*  
624 *processing systems*, **30**, <https://doi.org/10.48550/arXiv.2103.09504>
- 625 Wang, Y., H. Wu, J. Zhang, Z. Gao, J. Wang, S. Y. Philip, and M. Long, 2022: Predrnn: A  
626 recurrent neural network for spatiotemporal predictive learning. *IEEE Transactions on*  
627 *Pattern Analysis and Machine Intelligence*, **45**(2), 2208-2225,  
628 <https://doi.org/10.1109/TPAMI.2022.3165153>
- 629 Xiao, C., N. Chen, C. Hu, K. Wang, J. Gong, and Z. Chen, 2019: Short and mid-term sea  
630 surface temperature prediction using time-series satellite data and LSTM-AdaBoost  
631 combination approach. *Remote Sensing of Environment*, **233**, 111358,  
632 <https://doi.org/10.1016/j.rse.2019.111358>

- 633 Zhang, C., D. Wang, Z. Liu, S. Lu, C. Sun, Y. Wei, and M. Zhang, 2022: Global Gridded Argo  
634 Dataset Based on Gradient-Dependent Optimal Interpolation. *Journal of Marine Science*  
635 *and Engineering*, **10**(5), 650, <https://doi.org/10.3390/jmse10050650>
- 636 Zhou, G., G. Han, W. Li, X. Wang, X. Wu, L. Cao, and C. Li, 2023: High-resolution gridded  
637 temperature and salinity fields from Argo floats based on a spatiotemporal  
638 four-dimensional multigrid analysis method. *Journal of Geophysical Research: Oceans*,  
639 e2022JC019386, <https://doi.org/10.1029/2022JC019386>  
640

in press