

Reconstructing and Forecasting Marine Dynamic Variable Fields across Space and Time Globally and Gaplessly

Zhixi Xiong^{1†}, Yukang Jiang^{1†}, Wenfang Lu³, Xueqin Wang², and Ting Tian^{1*}

¹School of Mathematics, Sun Yat-sen University, Guangzhou, 510275, China

²School of Management, University of Science and Technology of China, Hefei, 230026, China

³School of Marine Sciences, Sun Yat-sen University, Guangzhou, 510275, China

*Corresponding author. E-mail: tiant55@mail.sysu.edu.cn

Contributing authors: xiongzhx5@mail2.sysu.edu.cn; jiangyk3@mail2.sysu.edu.cn; luwf6@sysu.edu.cn; wangxq20@ustc.edu.cn;

[†]These authors contributed equally to this work.

August 6, 2024

Abstract

Spatiotemporal projections in marine science are essential for understanding ocean systems and their impact on Earth’s climate. However, existing AI-based and statistics-based inversion methods face challenges in leveraging ocean data, generating continuous outputs, and incorporating physical constraints. We propose the Marine Dynamic Reconstruction and Forecast Neural Networks (MDRF-Net), which integrates marine physical mechanisms and observed data to reconstruct and forecast continuous ocean temperature-salinity and dynamic fields. MDRF-Net leverages statistical theories and techniques, incorporating parallel neural network sharing initial layer, two-step training strategy, and ensemble methodology, facilitating in exploring challenging marine areas like the Arctic zone. We have theoretically justified the efficacy of our ensemble method and the rationality of it by providing an upper bound on its generalization error. The effectiveness of MDRF-Net’s is validated through a comprehensive simulation study, which highlights its capability to reliably estimate unknown parameters. Comparison with other inversion methods and reanalysis data are also conducted, and the global test error is 0.455°C for temperature and 0.0714psu for salinity. Overall,

MDRF-Net effectively learns the ocean dynamics system using physical mechanisms and statistical insights, contributing to a deeper understanding of marine systems and their impact on the environment and human use of the ocean.

Keywords: Fields inversion, Global ocean dynamics, Primitive equations, Uncollected marine variables.

1 Introduction

Ocean changes quantified by the marine variable fields have significant implications for the economic and social systems that rely on them (Tittensor et al., 2021; Harley et al., 2006). For example, temperature, salinity, and current behaviors play a crucial role in the ocean state and the performance and sustainability of many organisms (Ashton et al., 2022; Smyth and Elliott, 2016). A core challenge in modern AI is the efficiency, exactness, and explainability of analyzing vast, multi-source ocean data, where statistical methods show great promise for extracting valuable insights relevant to both marine science and human activities in this realm. (Salcedo-Sanz et al., 2020). However, the collection of anthropogenically gathered ocean big data from various sources, such as buoys equipped with sensors, satellite remote sensing, and ship-based or airborne radars, has historically been sporadic and limited, resulting in difficulties quantifying continuous marine variability from in situ measurements (Johnson et al., 2022). Improving the comprehension of the oceans’ climate impact and facilitating enhanced decision-making across numerous marine-related industries can be achieved by utilizing statistics and AI, which are instrumental in deciphering the immense and varied data available (Bauer et al., 2015; Fox-Kemper et al., 2019).

Recent advancements in predicting the spatial and temporal variability of marine variable fields have demonstrated the potential of AI-based (Xiao et al., 2019; Xie et al., 2019; Song et al., 2020; Su et al., 2021; Song et al., 2022; Zhang et al., 2023) and statistical

inversion methods (Lee et al., 2019; Yarger et al., 2022). However, AI-based methods often encounter challenges, such as the requirement for time-series gridded data and the inability to provide continuous spatiotemporal predictions. On the other hand, the statistical methods are structurally complex, difficult to implement, and may struggle with handling large-scale data. Furthermore, these methods are purely data-driven, lacking interpretability, and do not fully exploit the relationships between different variable fields. Employing statistical methodologies and AI techniques plays a critical role in developing an integrated, robust, and interpretable framework for analyzing marine data and predicting ocean states (Lou et al., 2023).

To address the aforementioned issues, we propose a mesh-free and easy-to-implement model, Marine Dynamic Reconstruction and Forecast Neural Networks (MDRF-Net), for the reconstruction and forecasting of ocean temperature-salinity and dynamic fields. It offers a seamless integration of informative mechanistic equations with diverse data sources and types, grounded in statistical rigor, to achieve both high training efficiency and enhanced accuracy in practical AI applications. On one hand, it integrates fundamental physical laws of fluid dynamics, particularly the primitive equations, to ensure predictions are grounded in reality, accounting for temperature-salinity interactions in global fluid movement and mass conservation (Hieber et al., 2016; Lions et al., 1992). On the other hand, MDRF-Net adeptly merges relatively granular temperature and salinity data derived from Argo (Wong et al., 2020), with the meticulous and comprehensive currents reanalysis data provided by European Union-Copernicus Marine Service (2020). Thus, MDRF-Net enables the reconstruction and forecasting of complete continuous variable fields based on partial and incomplete observations.

MDRF-Net is built upon the fundamental framework of physics-informed neural net-

works as developed by Raissi et al. (Raissi et al., 2019), and three distinct design features of MDRF-Net contribute to its enhanced performance. Firstly, employing an ensemble method to rotate the Earth’s coordinate system facilitates multiple modeling and weighted averaging of results from diverse sub-learners, allowing for multiple models to work synergistically, compensating for individual weaknesses and improving overall accuracy and robustness of the predictions. Secondly, employing parallel neural networks that share the initial layer enables the sharing of fundamental information across different variable fields, while the parallel structural design of networks enhances fitting and predictive performance for variable fields in ultra-large-scale data. Finally, adopting a two-step training method reduces the redundant calculation of partial derivatives during the initial training stages, thereby easing the computational burden and speeds up the training process without compromising the model’s convergence or final performance.

Theoretically, we show the effectiveness of ensemble method and give an upper bound on the generalization error of MDRF-Net, proving the convergence of the solution provided that the model is adequately trained and the assumptions on the conditional stability estimate of the system of primitive equations hold. At the same time, when the sample domain is significantly smaller than the whole domain, namely, there are a large number of sea areas that cannot be reached by the detectors in reality, other inversion methods without adding mechanisms do not have this convergence, that is, they are unable to give the upper bound of the generalization error.

MDRF-Net’s validation via simulations reveals its superior performance in reconstructing missing data, enhancing field estimation precision, and stably identifying unknown equation parameters compared to similar methods. It also demonstrates consistent performance with real datasets, adapting to various spatial-temporal prediction tasks and data

volumes. Globally, MDRF-Net provides continuous space-time inference of ocean changes and temporal extrapolation, surpassing reanalysis data in certain aspects. It effectively identifies key ocean phenomena, such as the Mediterranean Salinity Crisis and the North Atlantic Warm Current, and monitors hard-to-observe regions like the Arctic. An interactive R Shiny platform (<https://tikitakatikitaka.shinyapps.io/mdrf-net-shiny/>) consolidates variable predictions for user exploration across spaces and times, and the code is accessible at <https://github.com/tikitaka0243/mdrf-net>.

2 Methodologies

2.1 Primitive equations

The goal of our research is to reconstruct and predict the ocean dynamics field, and the primitive equations (Hieber et al., 2016; Lions et al., 1992), which describe the ocean dynamics system including current motion and thermohaline diffusion effects, bring important physical information to the table. It is given by Equation 1 and contains the momentum (Equation 1a), hydrostatic balance (Equation 1b) and continuity (Equation 1c) equations, the equations of temperature (Equation 1d, thermodynamical equation) and

salinity (Equation 1e), and the equation of state (Equation 1f).

$$\frac{\partial \mathbf{v}}{\partial t} + \nabla_{\mathbf{v}} \mathbf{v} + w \frac{\partial \mathbf{v}}{\partial r_a} + \frac{1}{\rho_0} \nabla_h p + 2\boldsymbol{\omega}_e(\mathbf{e}_r \times \mathbf{v}) \cos \theta - \zeta \Delta \mathbf{v} - \eta \frac{\partial^2 \mathbf{v}}{\partial r_a^2} = 0, \quad (1a)$$

$$\frac{\partial p}{\partial r_a} = -\rho g, \quad (1b)$$

$$\text{div } \mathbf{v} + \frac{\partial w}{\partial r_a} = 0, \quad (1c)$$

$$\frac{\partial \tau}{\partial t} + \nabla_{\mathbf{v}} \tau + w \frac{\partial \tau}{\partial r_a} - \zeta_{\tau} \Delta \tau - \eta_{\tau} \frac{\partial^2 \tau}{\partial r_a^2} = 0, \quad (1d)$$

$$\frac{\partial \sigma}{\partial t} + \nabla_{\mathbf{v}} \sigma + w \frac{\partial \sigma}{\partial r_a} - \zeta_{\sigma} \Delta \sigma - \eta_{\sigma} \frac{\partial^2 \sigma}{\partial r_a^2} = 0, \quad (1e)$$

$$\rho_0[1 - \beta_{\tau}(\tau - \tau_0) + \beta_{\sigma}(\sigma - \sigma_0)] = \rho. \quad (1f)$$

Here we use $\tau, \sigma, w, v_{\theta}, v_{\phi}, p$, and ρ to denote temperature, salinity, vertical velocity, northward velocity, eastward velocity, pressure, and density respectively, they are all functions of r, θ, φ, t , namely radial distance, polar angle (transformed latitude), azimuthal angle (transformed longitude), and time. $\mathbf{v} = (v_{\theta}, v_{\phi})$ represents the horizontal velocity. $r = r_a + r_e$, r_e is the radius of the Earth and $r_a \leq 0$ the vertical coordinate with regard to the sea surface. $\rho_0 > 0$, $\tau_0 > 0$, $\sigma_0 > 0$ are the reference values of the density, temperature and salinity. $\boldsymbol{\omega}_e$ is the vector angular velocity of the Earth and \mathbf{e}_r the vertical unit vector. ∇_h represents the horizontal gradient operator and \times the cross product operator. η, ζ are eddy viscosity coefficients and $\eta_{\tau}, \zeta_{\tau}$ and $\eta_{\sigma}, \zeta_{\sigma}$ are eddy diffusivity coefficients for temperature and salinity respectively. Positive expansion coefficients $\beta_{\tau}, \beta_{\sigma}$, are used, along with the Laplacian Δ , gradient $\nabla_{\mathbf{v}}$ with respect to velocity, and divergence operator div .

And there are initial conditions (Equation 2) and boundary conditions (Equation 3), namely space-time boundary conditions:

$$\mathbf{v} = \mathbf{i}, \quad \tau = i_{\tau}, \quad \sigma = i_{\sigma}, \quad (2)$$

$$\frac{\partial \mathbf{v}}{\partial r_a} = \boldsymbol{\delta}_v, \quad w = 0, \quad \frac{\partial \tau}{\partial r_a} + \alpha(\tau - \tau_a) = 0, \quad \frac{\partial \sigma}{\partial r_a} = 0, \quad \text{on } \Gamma_u, \quad (3a)$$

$$\mathbf{v} = \mathbf{0}, \quad w = 0, \quad \tau = b_\tau, \quad \sigma = b_\sigma, \quad \text{on } \Gamma_b, \quad (3b)$$

$$\mathbf{v} = \mathbf{0}, \quad w = 0, \quad \frac{\partial \tau}{\partial \boldsymbol{\psi}} = 0, \quad \frac{\partial \sigma}{\partial \boldsymbol{\psi}} = 0, \quad \text{on } \Gamma_l, \quad (3c)$$

where i, i_τ, i_σ are the initial values. Γ_u, Γ_b , and Γ_l are the upper, bottom, and lateral boundaries of the ocean respectively. The wind stress, denoted as $\boldsymbol{\delta}_v$, is contingent upon the velocity of the atmosphere. α is a positive constant associated with the turbulent heating on the ocean's surface, τ_a is the apparent atmospheric equilibrium temperature, and b_τ as well as b_σ , which are functions of the latitude θ and longitude φ , represent the temperature and salinity of the sea at the ocean's bottom. On the boundary, the normal vector is $\boldsymbol{\psi}$. More details about the primitive equations are in Appendix A.

The underlying domain of the primitive equations is $\Omega = [\{r_a\}_{min} + r_e, r_e] \times [0, \pi) \times [0, 2\pi) \in \mathbb{R}$ and the boundary with continuous first order derivatives is $\Gamma = \Gamma_u \cup \Gamma_b \cup \Gamma_l$. We also have the space-time domain and boundary $\bar{\Omega} = \Omega \times [0, T] \subset \mathbb{R}^4$ and $\bar{\Gamma} = \Gamma \times [0, T] \cup \Omega \times \{t = 0\}$. Then the primitive equations (Equation 1) can be generally represented as

$$\mathcal{F}_\beta(\mathbf{u}) = \mathbf{g} \quad (4)$$

where $\mathcal{F} : L^2(\bar{\Omega}, \mathbb{R}^6) \mapsto L^2(\bar{\Omega}, \mathbb{R}^6)$ is a general operator for the primitive equations and $L^p(X, Y)$ represents the p -norm finite function space mapping from X to Y ; $\beta = (\beta_\tau, \beta_\sigma)$ represents the unknown parameters of the equations; $\mathbf{u} = (\tau, \sigma, w, v_\theta, v_\varphi, p) \in L^2(\bar{\Omega}, \mathbb{R}^6)$ denotes the ocean variable fields and notice that the density field ρ can be derived from temperature field τ and salinity field σ ; $\mathbf{g} \in L^2(\bar{\Omega}, \mathbb{R}^6)$ is the source term and here we only consider the effect of gravity. We also assume that,

$$\|\mathcal{F}_\beta(\mathbf{u})\|_{\mathbb{R}^6} < +\infty, \quad \forall \mathbf{u} \in \mathbb{R}^6, \quad \text{with } \|\mathbf{u}\|_{\mathbb{R}^6} < +\infty, \quad \text{and} \quad \|\mathbf{g}\|_{\mathbb{R}^6} < +\infty. \quad (5)$$

where $\|\cdot\|_{\mathbb{R}^d}$ represents the norm on the d -dimensional real space.

The general form of the space-time boundary conditions (Equation 2, 3) is

$$\mathcal{B}(\text{tr}(\mathbf{u})) = \mathbf{b}, \quad (6)$$

with the general boundary operator $\mathcal{B} : L^2(\bar{\Gamma}, \mathbb{R}^6) \mapsto L^2(\bar{\Gamma}, \mathbb{R}^6)$, the trace operator $\text{tr} : L^2(\bar{\Omega}, \mathbb{R}^6) \mapsto L^2(\bar{\Gamma}, \mathbb{R}^6)$ and $\mathbf{b} \in L^2(\bar{\Gamma}, \mathbb{R}^6)$. All of the operators are bounded under the corresponding norm, that is

$$\|\mathcal{B}_\beta(\text{tr}(\mathbf{u}))\|_{\mathbb{R}^6} < +\infty, \quad \forall \mathbf{u} \in \mathbb{R}^6, \quad \text{with } \|\mathbf{u}\|_{\mathbb{R}^6} < +\infty, \quad \text{and} \quad \|\mathbf{b}\|_{\mathbb{R}^6} < +\infty. \quad (7)$$

It is known that the forward problem of the primitive equations (Equation 4, 6) is well-posed (Hieber et al., 2016; Charve, 2008), namely, given $\mathbf{g} \in L^2(\bar{\Omega}, \mathbb{R}^6)$ and $\mathbf{b} \in L^2(\bar{\Gamma}, \mathbb{R}^6)$, there exists a unique solution $\mathbf{u} \in L^2(\bar{\Omega}, \mathbb{R}^6)$ satisfying the equations (Equation 4, 6) and continuously depending on changes in boundary conditions.

In ocean dynamic systems, we do not know the complete boundary conditions and the equations' parameters. Thus, the forward problem for the primitive equations will be ill-posed, that is, no unique solution can be guaranteed. However, some observations \mathbf{u}_{obs} of the underlying solution \mathbf{u} in a sub-domain $\bar{\Omega}' \subset \bar{\Omega}$ (observation domain) are available, such as temperature and salinity data from Argo (Wong et al., 2020) and three-dimensional velocity data from European Union-Copernicus Marine Service (2020):

$$\mathbf{u} = \mathbf{u}_{\text{obs}} + \boldsymbol{\epsilon}, \quad \forall \mathbf{x} \in \bar{\Omega}', \quad (8)$$

where $\boldsymbol{\epsilon}$ is a noise term and $\mathbf{x} = (r, \theta, \varphi, t)$ represents the space-time coordinate. We assume that

$$\|\mathbf{u}_{\text{obs}}\|_{L^2(\bar{\Omega}')} < +\infty \quad (9)$$

Equations 4, 6, and 8 form the inverse problem of the primitive equations, and we assume that it has conditional stability estimate:

Assumption 2.1 (Conditional stability estimate). *For any $\mathbf{u}_1, \mathbf{u}_2 \in L^2(\bar{\Omega}, \mathbb{R}^6)$, the primitive equations satisfy*

$$\begin{aligned} \|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2(S)} \leq C \left(\|\mathbf{u}_1\|_{L^2(\bar{\Omega})}, \|\mathbf{u}_2\|_{L^2(\bar{\Omega})} \right) \cdot \left(\|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2(\bar{\Omega}')}^{\gamma'} + \right. \\ \left. \|\mathcal{F}_\beta(\mathbf{u}_1) - \mathcal{F}_\beta(\mathbf{u}_2)\|_{L^2(\bar{\Omega})}^{\gamma_1} + \|\mathcal{B}(\text{tr}(\mathbf{u}_1)) - \mathcal{B}(\text{tr}(\mathbf{u}_2))\|_{L^2(\bar{\Gamma})}^{\gamma_2} \right) \end{aligned} \quad (10)$$

for some $0 < \gamma', \gamma_1, \gamma_2 \leq 1$ and for any subset $\bar{\Omega}' \subset S \subset \bar{\Omega}$.

2.2 Marine Dynamic Reconstruction and Forecast Neural Networks (MDRF-Net)

In order to enable the reconstruct and forecast ocean dynamic fields, including temperature, salinity, vertical, northward and eastward velocities, and pressure fields, we constructed the Marine Dynamic Reconstruction and Forecast Neural Networks (MDRF-Net), which seamlessly integrates information from the data and the primitive equations in a concise structure by embedding the equations into the loss function. And since Copernicus' current reanalysis data is already highly refined and comprehensive, our focus is on using it as a bridge to connect the primitive equations and enabling MDRF-Net to reconstruct the temperature and salinity fields with more accurate and comprehensive information about the ocean current system.

Inside MDRF-Net (Figure 1), there is a parallel fully connected neural network that shares the first layer. This design allows for the adaptation of the neural networks to different sizes of marine variable fields with varying complexities. Additionally, we propose a two-step training strategy where parameters are first optimized using observed data, and then using the loss of physics and the observed data simultaneously. The second step adds the primitive equations into the loss function to provide spatial-temporal physics information. This approach eliminates the need for the model to calculate a large number

of partial derivatives when its outputs deviate from expectations, resulting in faster and more efficient training. To improve the model’s performance at high latitudes, we also use an ensemble method by rotating the Earth’s coordinate system multiple times along the 0° longitude (prime meridian), thus modeling the model multiple times and taking a weighted average of the results. More detailed information about MDRF-Net can be found in the Appendix C.

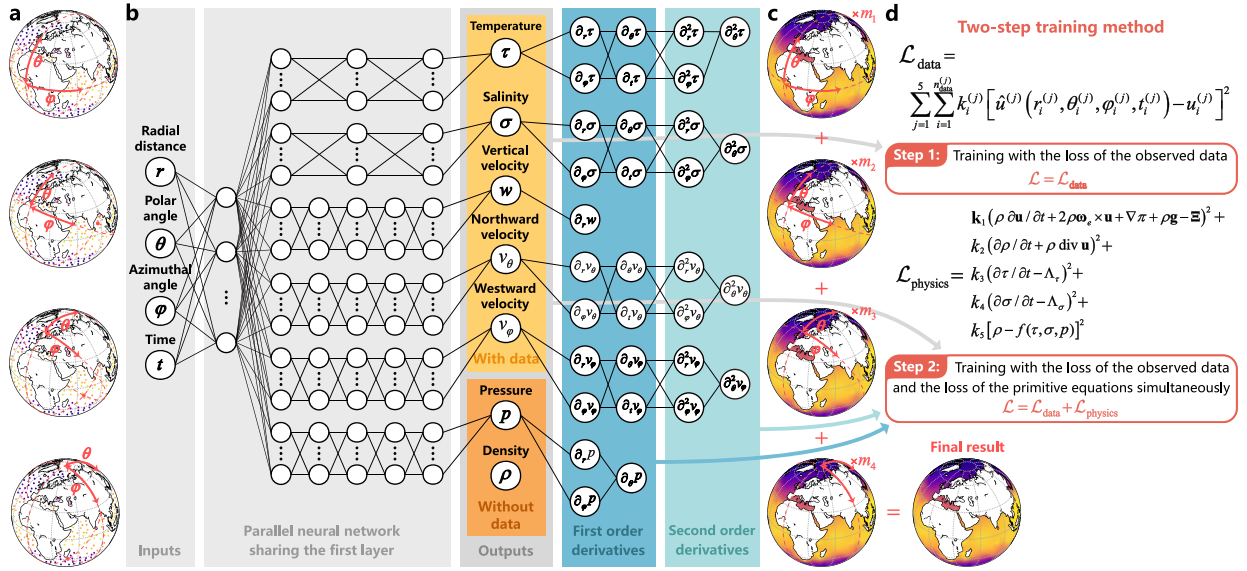


Figure 1: Marine Dynamic Reconstruction and Forecast Neural Networks (MDRF-Net). Before the training data is fed into the neural network, its coordinate system is rotated multiple times along the 0-degree meridian (a). The main body of MDRF-Net is a parallel neural network that shares the first layer (b), with inputs of coordinates and time, and outputs the variables that need to be differentiated in the primitive equation. There are three layers of the fully connected neural networks for the temperature and salinity fields, and five layers for the other four fields. The first and second-order partial derivatives of each ocean variable are calculated to compute the loss of equations. Finally, the results derived in each coordinate system are weighted and averaged to obtain the final outcome (c). In the two-step training strategy (d), the parameters of neural networks are optimized based on the loss of data in the first step and on the loss of data and physics together in the second step. The formula for $\mathcal{L}_{\text{physics}}$ here is a simplified version, see the Materials and Methods section for the precise implementation.

2.3 Implementation of MDRF-Net

Consider an input $\mathbf{x} \in \bar{\Omega}$, we represent the neural network of MDRF-Net by the following:

$$\begin{aligned} \mathbf{u}_{\Theta}^{(j)}(\mathbf{x}) &= C_{K^{(j)}}^{(j)} \circ \sigma_{\tanh} \circ C_{K^{(j)}-1}^{(j)} \circ \cdots \circ \sigma_{\tanh} \circ C_2^{(j)} \circ \sigma_{\tanh} \circ C_1(\mathbf{x}), \\ \hat{\mathbf{u}}(\mathbf{x}) &= [\mathbf{u}_{\Theta}^{(1)}(\mathbf{x}), \mathbf{u}_{\Theta}^{(2)}(\mathbf{x}), \dots, \mathbf{u}_{\Theta}^{(6)}(\mathbf{x})]. \end{aligned} \quad (11)$$

where $\mathbf{u}_{\Theta}(\mathbf{x})$ is the variable fields represented by a parallel neural network sharing the first layer with parameters Θ ; $\mathbf{u}_{\Theta}^{(j)}(\mathbf{x})$ is the j^{th} sub-network for the j^{th} variable field; \circ is the composition of the functions and σ_{\tanh} is the tanh activation function. $C_k^{(j)}$ represents the linear transformation of the k^{th} layer of the j^{th} sub-network, and

$$C_k^{(j)}(\mathbf{x}_k^{(j)}) = \mathbf{W}_k^{(j)} \mathbf{x}_k^{(j)} + \mathbf{b}_k^{(j)}, \quad 1 \leq j \leq 6, \quad 2 \leq k \leq K^{(j)}, \quad (12)$$

where $\mathbf{W}_k^{(j)} \in \mathbb{R}^{d_{k+1}^{(j)} \times d_k^{(j)}}$ is the weights matrix, $\mathbf{b}_k^{(j)} \in \mathbb{R}^{d_{k+1}^{(j)}}$ is the bias term and $\mathbf{x}_k^{(j)} \in \mathbb{R}^{d_k^{(j)}}$ is the output of the $(k-1)^{\text{th}}$ layer. Here, we set $K^{(1)} = K^{(2)} = 3$ and $K^{(3)} = K^{(4)} = \dots = K^{(6)} = 5$; $d_2^{(j)} = d_3^{(j)} = \dots = d_{K^{(j)}}^{(j)} = 128$, $d_{K^{(j)}+1}^{(j)} = 1$, $1 \leq j \leq 6$. For the first shared layer, similarly, $C_1(\mathbf{x}_1) = \mathbf{W}_1 \mathbf{x}_1 + \mathbf{b}_1$, $\mathbf{W}_1 \in \mathbb{R}^{d_2^{(j)} \times 4}$, $\mathbf{x}_1 \in \mathbb{R}^4$, $\mathbf{b}_1 \in \mathbb{R}^{d_2^{(j)}}$, $1 \leq j \leq 6$. Then the parameters of the neural network are $\Theta = \{\mathbf{W}_k^{(j)}, \mathbf{b}_k^{(j)}, \mathbf{W}_1, \mathbf{b}_1\}$, $\forall 1 \leq k \leq K^{(j)}$, $1 \leq j \leq 6$.

Given the primitive equations with initial and boundary conditions and the observed data, the residuals of the data and physics information for MDRF-Net \mathbf{u}_{Θ} and the equations' unknown parameters β are $\mathcal{L}_{\text{data}}(\mathbf{u}_{\Theta}) := \mathbf{u}_{\Theta} - \mathbf{u}_{\text{obs}} - \epsilon$, $\mathcal{L}_{\text{pde}}(\mathbf{u}_{\Theta}, \beta) := \mathcal{F}_{\beta}(\mathbf{u}_{\Theta}) - \mathbf{g}$, and $\mathcal{L}_{\text{icbc}}(\mathbf{u}_{\Theta}) := \mathcal{B}(\text{tr}(\mathbf{u}_{\Theta})) - \mathbf{b}$. By the first three assumptions (Equation 5, 7, 9), we know that $\mathcal{L}_{\text{data}} \in L^2(\bar{\Omega}', \mathbb{R}^6)$, $\mathcal{L}_{\text{pde}} \in L^2(\bar{\Omega}, \mathbb{R}^6)$, $\mathcal{L}_{\text{icbc}} \in L^2(\bar{\Gamma}, \mathbb{R}^6)$, and $\|\mathcal{L}_{\text{data}}\|_{L^2(\bar{\Omega}')} , \|\mathcal{L}_{\text{pde}}\|_{L^2(\bar{\Omega})} \|\mathcal{L}_{\text{icbc}}\|_{L^2(\bar{\Gamma})} < +\infty$, for all $\Theta \in \Theta$, $\beta \in B$, Θ and B are the parameter spaces of the neural network and the primitive equations respectively.

Therefore, the optimization problem of MDRF-Net is

$$\Theta^*, \beta^* = \underset{\substack{\Theta \in \Theta \\ \beta \in B}}{\text{argmin}} \left\{ \|\mathcal{L}_{\text{data}}\|_{L^2(\bar{\Omega}')} + \lambda'_1 \|\mathcal{L}_{\text{pde}}\|_{L^2(\bar{\Omega})} + \lambda'_2 \|\mathcal{L}_{\text{icbc}}\|_{L^2(\bar{\Gamma})} \right\}. \quad (13)$$

Equivalently,

$$\Theta^*, \beta^* = \underset{\substack{\Theta \in \Theta \\ \beta \in B}}{\operatorname{argmin}} \left\{ \int_{\bar{\Omega}'} |\mathcal{L}_{\text{data}}|^2 d\mathbf{x} + \lambda_1 \int_{\bar{\Omega}} |\mathcal{L}_{\text{pde}}|^2 d\mathbf{x} + \lambda_2 \int_{\bar{\Gamma}} |\mathcal{L}_{\text{icbc}}|^2 d\mathbf{x} \right\}. \quad (14)$$

where λ'_1, λ_1 and λ'_2, λ_2 are additional regularization hyper-parameters, and the physics residual terms can be considered as penalization terms. Additionally, the problem is able to be approximated using the quadratic rules, with the coordinates of the observations $\{\mathbf{x}'_i\}$ with all $1 \leq i \leq N_{\text{data}}$, combined with the selected sampling points inside the space-time domain $\{\mathbf{x}_{1,i}\}$, $\mathbf{x}_{1,i} \in \bar{\Omega}$, $1 \leq i \leq N_{\text{pde}}$ and points at the space-time boundary $\{\mathbf{x}_{2,i}\}$, $\mathbf{x}_{2,i} \in \bar{\Gamma}$, $1 \leq i \leq N_{\text{icbc}}$. Then we have the following loss function:

$$J(\Theta, \beta) := \sum_{i=1}^{N_{\text{data}}} k'_i |\mathcal{L}_{\text{data}}(\mathbf{x}'_i)|^2 + \lambda_1 \sum_{i=1}^{N_{\text{pde}}} k_{1,i} |\mathcal{L}_{\text{pde}}(\mathbf{x}_{1,i})|^2 + \lambda_2 \sum_{i=1}^{N_{\text{icbc}}} k_{2,i} |\mathcal{L}_{\text{icbc}}(\mathbf{x}_{2,i})|^2 \quad (15)$$

where k'_i , $k_{1,i}$ and $k_{2,i}$ are weights. When calculating the errors of the observed data ($\mathcal{L}_{\text{data}}$), the sampling points are naturally the positions in space and time of the samples in the training data set $(r_i^{(j)}, \theta_i^{(j)}, \varphi_i^{(j)}, t_i^{(j)}), i = 1, 2, \dots, N_{\text{data}}^{(j)}, j = 1, 2, \dots, 5$; and when computing the loss of the primitive equations (\mathcal{L}_{pde}), the sampling points can come from the location of the observed data, or can be generated additionally within the equations' underlying domain $\bar{\Omega} = \Omega \times [0, T]$. In this study, additional gridded sampling points are generated for calculating the equations' loss in the global scenario due to the non-uniform distribution of observed data locations. It is important to note that the number of observations strictly at the initial time or definition domain boundaries is usually small. Therefore, separate sampling points need to be generated for calculating the errors in the initial conditions and boundary conditions ($\mathcal{L}_{\text{icbc}}$) based on the definition domain.

During the training of MDRE-Net, the partial derivatives of first and second order of the outputs \mathbf{u}_{Θ} with respect to the inputs \mathbf{x} are cleverly computed by the automatic differentiation algorithm required in the optimization process of neural networks. As in

classical neural networks, MDRF-Net uses the gradient descent method based on the back propagation algorithm to optimize the parameters. For the unknown parameters of the primitive equations, MDRF-Net optimizes them together with the parameters of the neural networks, thus solving the ‘inverse problem’ of the primitive equations.

Regarding the ensemble method, suppose we have inversion results for both the original and rotated variable fields, represented as $\hat{\mathbf{u}}^{(r)}, r = 1, \dots, N_{ro}$. Then the final weighted result at the polar angle θ is

$$\hat{\mathbf{u}}(\theta) = \frac{\sum_{r=1}^{N_{ro}} m_r(\theta_r) \hat{\mathbf{u}}^{(r)}(\theta_r)}{\sum_{r=1}^{N_{ro}} m_r(\theta_r)}, \quad (16)$$

where $m_r(\theta_r) = \text{logistic}(10(\theta_r/\pi - 0.5))$, θ_r is the polar angle of the rotated coordinates.

2.4 Estimation of Generalization Error

The generalization error is the error of MDRF-Net on unseen data. We set $\bar{\Omega}' \subset S \subset \bar{\Omega}$ and define the corresponding generalization error as

$$\mathcal{E}(S) = \mathcal{E} \left(S; \Theta^*, \boldsymbol{\beta}^*, \{\mathbf{x}'_i\}_{i=1}^{N_{\text{data}}}, \{\mathbf{x}_{1,i}\}_{i=1}^{N_{\text{pde}}}, \{\mathbf{x}_{2,i}\}_{i=1}^{N_{\text{icbc}}} \right) = \|\mathbf{u} - \hat{\mathbf{u}}\|_{L^2(S)}, \quad (17)$$

which depends on the train sampling points, data and the optimized MDRF-Net’s parameters Θ^* and equations’ unknown parameters $\boldsymbol{\beta}^*$; $\hat{\mathbf{u}} = \mathbf{u}_{\Theta^*}$ represents the trained MDRF-Net with Θ^* as its parameters.

Due to neural networks’ inability to well capture the variation pattern of longitudinal density with latitude, we design an ensemble method involving multiple rotations in spherical coordinates, and its effectiveness can be given by the following theorem.

Theorem 2.1 (The effectiveness of ensemble method). *Let $\mathbf{u} \in L^2(\bar{\Omega}, \mathbb{R}^6)$ be the solution of the inverse problem of the primitive equations, and $\hat{\mathbf{u}}, \hat{\mathbf{u}}^{(r)}$ are the trained MDRF-Net*

and the r^{th} sub-learner. Then the generalization error for MDRF-Net with ensemble method satisfies

$$\|\mathbf{u} - \hat{\mathbf{u}}\| \leq \|\mathbf{u} - \hat{\mathbf{u}}^{(r)}\|, \quad \forall 1 \leq r \leq N_{ro} \quad (18)$$

We can estimate the generalization error in terms of the training error $\mathcal{E}_{\text{train}} = \mathcal{E}_{\text{train}}(\Theta^*, \{\mathbf{x}'_i\}_{i=1}^{N_{\text{data}}}, \{\mathbf{x}_{1,i}\}_{i=1}^{N_{\text{pde}}}, \{\mathbf{x}_{2,i}\}_{i=1}^{N_{\text{icbc}}})$, defined by:

$$\mathcal{E}_{\text{train}} = (\mathcal{E}_{\text{data}}^2 + \lambda_1 \mathcal{E}_{\text{pde}}^2 + \lambda_2 \mathcal{E}_{\text{icbc}}^2)^{\frac{1}{2}}, \quad (19)$$

where $\mathcal{E}_{\text{data}} := (\sum_{i=1}^{N_{\text{data}}} k'_i |\mathcal{L}_{\text{data}}(\mathbf{x}'_i)|^2)^{\frac{1}{2}}$, $\mathcal{E}_{\text{pde}} := (\sum_{i=1}^{N_{\text{pde}}} k_{1,i} |\mathcal{L}_{\text{pde}}(\mathbf{x}_{1,i})|^2)^{\frac{1}{2}}$, and $\mathcal{E}_{\text{icbc}} := (\sum_{i=1}^{N_{\text{icbc}}} k_{2,i} |\mathcal{L}_{\text{icbc}}(\mathbf{x}_{2,i})|^2)^{\frac{1}{2}}$. The bound on generalization error in terms of training error $\mathcal{E}_{\text{train}}$ is given by the following estimate (Mishra and Molinaro, 2021):

Theorem 2.2 (Upper bound of generalization error). *Let $\mathbf{u} \in L^2(\bar{\Omega}, \mathbb{R}^6)$ be the solution of the inverse problem of the primitive equations. Assume the conditional stability estimate assumption (Equation 10) holds for any S , s.t. $\bar{\Omega}' \subset S \subset \bar{\Omega}$. Let $\hat{\mathbf{u}}$ be the trained MDRF-Net, based on the training sampling points $\{\mathbf{x}_{1,i}\}_{i=1}^{N_{\text{pde}}}$, $\{\mathbf{x}_{2,i}\}_{i=1}^{N_{\text{icbc}}}$ and the observed data coordinates $\{\mathbf{x}'_i\}_{i=1}^{N_{\text{data}}}$. Additionally, assume the residuals \mathcal{L}_{pde} , $\mathcal{L}_{\text{icbc}}$ and $\mathcal{L}_{\text{data}}$ be such that $\mathcal{L}_{\text{data}} \in L^2(\bar{\Omega}', \mathbb{R}^6)$, $\mathcal{L}_{\text{pde}} \in L^2(\bar{\Omega}, \mathbb{R}^6)$ and $\mathcal{L}_{\text{icbc}} \in L^2(\bar{\Gamma}, \mathbb{R}^6)$ and the quadrature error is bounded. Then the following estimate on the generalization error holds:*

$$\begin{aligned} \mathcal{E}(S) \leq C & \left(\mathcal{E}_{\text{data}}^{\gamma'} + \mathcal{E}_{\text{pde}}^{\gamma_1} + \mathcal{E}_{\text{icbc}}^{\gamma_2} + \|\epsilon\|_{L^2(\bar{\Omega}')}^{\gamma'} + \right. \\ & \left. (C'_q)^{\frac{\gamma'}{2}} N_{\text{data}}^{-\frac{a'\gamma'}{2}} + (C_{1,q})^{\frac{\gamma_1}{2}} N_{\text{pde}}^{-\frac{a_1\gamma_1}{2}} + (C_{2,q})^{\frac{\gamma_2}{2}} N_{\text{icbc}}^{-\frac{a_2\gamma_2}{2}} \right), \end{aligned} \quad (20)$$

with constants $C = C(\|\mathbf{u}\|_{L^2(\bar{\Omega})}, \|\hat{\mathbf{u}}\|_{L^2(\bar{\Omega})})$, $C'_q = C'_q(\|\mathcal{L}_{\text{data}}\|_{L^2(\bar{\Omega}')}^2)$, $C_{1,q} = C_q(\|\mathcal{L}_{\text{pde}}\|_{L^2(\bar{\Omega})}^2)$, $C_{2,q} = C_q(\|\mathcal{L}_{\text{icbc}}\|_{L^2(\bar{\Gamma})}^2)$.

Notice that Theorem 2.2 is satisfied by the precondition that the conditional stability estimate (Equation 10) of the primitive equations hold. And when MDRF-Net has no

constraints from the ocean dynamics mechanism, that is, the primitive equations are not embedded in the loss function, the model degenerates into a purely data-driven neural network. At this point we can view the training of the model as a special ‘inverse problem’, i.e., $\mathcal{F}_\beta(\mathbf{u}) = c_1$ and $\mathcal{B}(\mathbf{u}) = c_2$ for any $\mathbf{u} \in L^2(\bar{\Omega}, \mathbb{R}^6)$, where c_1 and c_2 are constants and the observations (Equation 8) remains the same. Now assume that the measure of the integration domain of the observation space $\bar{\Omega}'$ is ϱ_i times of the whole space $\bar{\Omega}$, $0 < \varrho_i < 1$ for all $1 \leq i \leq 4$. Then the norm

$$\begin{aligned} \|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2(\bar{\Omega}')}^{\gamma'} &= \left(\int_{\bar{\Omega}'} |\mathbf{u}_1 - \mathbf{u}_2|^2 d\mathbf{x} \right)^{\frac{\gamma'}{2}} \\ &= \left(\prod_{i=1}^4 \varrho_i \int_{\bar{\Omega}} |\mathbf{u}_1 - \mathbf{u}_2|^2 d\mathbf{x} \right)^{\frac{\gamma'}{2}} = \left(\prod_{i=1}^4 \varrho_i \right)^{\frac{\gamma'}{2}} \|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2(\bar{\Omega})}^{\gamma'}, \end{aligned} \quad (21)$$

which can not be control by the constant $C(\|\mathbf{u}_1\|_{L^2(\bar{\Omega})}, \|\mathbf{u}_2\|_{L^2(\bar{\Omega})})$, since ϱ_i is determined by the observation space $\bar{\Omega}'$. Hence, there exists a space S' satisfying $\bar{\Omega}' \subset S' \subset \bar{\Omega}$ that

$$\|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2(S')} > C(\|\mathbf{u}_1\|_{L^2(\bar{\Omega})}, \|\mathbf{u}_2\|_{L^2(\bar{\Omega})}) \cdot \|\mathbf{u}_1 - \mathbf{u}_2\|_{L^2(\bar{\Omega}')}^{\gamma'} \quad (22)$$

for any $0 < \gamma' \leq 1$. Then the conditional stability estimate cannot be satisfied and an upper bound on the generalization error cannot be given. In fact, this derivation applies to all methods that do not incorporate mechanism, such as Gaussian Process Regression and Regression Kriging, which will be discussed later. This implies that when the data domain is not sufficiently large, these methods lack an upper bound on their generalization error.

3 Results

3.1 Simulation Study

We first explore the capabilities of MDRF-Net in a simulated system by considering a 2D simplified version (Equation 23) of the primitive equations (Equation 1) which has only

one dimension in the horizontal direction and does not include the diffusion equation for salinity (Equation 1e) as well as the equation of state (Equation 1f), so it has only four variables temperature τ , horizontal velocity v , vertical velocity w and pressure p . Unlike Equation 1, the underlying domain of this primitive equations is based on a Cartesian coordinate system rather than a spherical coordinate system, and it is dimensionless, that is, all variables take on values without units.

$$\begin{aligned}
\frac{\partial v}{\partial t} + v \frac{\partial v}{\partial x} + w \frac{\partial v}{\partial z} - \eta \frac{\partial^2 v}{\partial x^2} - \zeta \frac{\partial^2 v}{\partial z^2} + \frac{\partial p}{\partial x} &= 0, \\
\frac{\partial p}{\partial z} &= -\tau, \\
\frac{\partial v}{\partial x} + \frac{\partial w}{\partial z} &= 0, \\
\frac{\partial \tau}{\partial t} + v \frac{\partial \tau}{\partial x} + w \frac{\partial \tau}{\partial z} - \eta_\tau \frac{\partial^2 \tau}{\partial x^2} - \zeta_\tau \frac{\partial^2 \tau}{\partial z^2} &= Q,
\end{aligned} \tag{23}$$

This system of equations has a specific Taylor-Green vortex solution for a specific periodic source term Q (Hu et al., 2023),

$$\begin{aligned}
v &= -\sin(2\pi x) \cos(2\pi z) \exp[-4\pi^2(\eta + \zeta)t], \\
w &= \cos(2\pi x) \sin(2\pi z) \exp[-4\pi^2(\eta + \zeta)t], \\
p &= \frac{1}{4} \cos(4\pi x) \exp[-8\pi^2(\eta + \zeta)t] + \frac{1}{2\pi} \cos(2\pi z) \exp(-4\pi^2\zeta_\tau t), \\
\tau &= \sin(2\pi z) \exp(-4\pi^2\zeta_\tau t), \\
Q &= \pi \cos(2\pi x) \sin(4\pi z) \exp[-4\pi^2(\eta + \zeta + \zeta_\tau)t].
\end{aligned} \tag{24}$$

Accordingly, we set $\eta = \zeta = 0.01$, $\zeta_\tau = 0.02$, and generated 1000 samples randomly in the data domain and without pressure variables to mimic the reality of the ocean data. Note that the Taylor-Green vortex is independent of the value of ζ_τ . In reconstructing these variable fields using MDRF-Net, we set the parameters ζ and ζ_τ to be unknown with initial value 0 and the value of η is correlated with ζ , which would be unrecognizable if they were both set to be unknown. We also examine the performance of MDRF-Net without mechanism (N-MDRF-Net), and the commonly used marine variable field interpolation

methods Gaussian Process Regression (GPR) and Regression Kriging (R-Kriging), which support ungridded data and provide continuous inversion results, on simulated datasets.

MDRF-Net perfectly reconstructs the real variable fields, even in domain where there is no data and for the variable without data. The other models can only reconstruct the variable fields fairly well in the domain with data, while it is basically a failure in the domain without data. For pressure fields with no data at all, GPR and R-Kriging fail to give results completely, whereas N-MDRF-Net still gives an output because it uses a parallel neural network that shares the first layer (Figure 2 **a**).

In terms of accuracy, the overall root mean square error (RMSE) of MDRF-Net is below 10^{-2} , which is far superior to the other methods over the whole domain, even in domain where data is available, which is consistent with the expectations derived from theoretical analyses. Notice that even though R-Kriging is not as continuous as the N-MDRF-Net and GPR results, the three methods that do not include the mechanism are at the same level of overall RMSE. In the data domain, N-MDRF-Net slightly outperforms GPR and R-Kriging, however, its predicted temperature fields exhibit a tendency to accumulate substantially larger errors over extended temporal horizons. In addition, the RMSEs of all the models show more or less a decreasing and then increasing pattern throughout the time span (Figure 2 **b**).

During the training process of MDRF-Net, the two unknown parameters, ζ and ζ_τ , steadily approach their true values of 0.01 and 0.02 from below and above, respectively. Notably, in the early stages of the solution procedure, considerable fluctuations are observed in the values of these unknowns, with ζ , for instance, momentarily reaching around 0.06 (Figure 2 **c**).

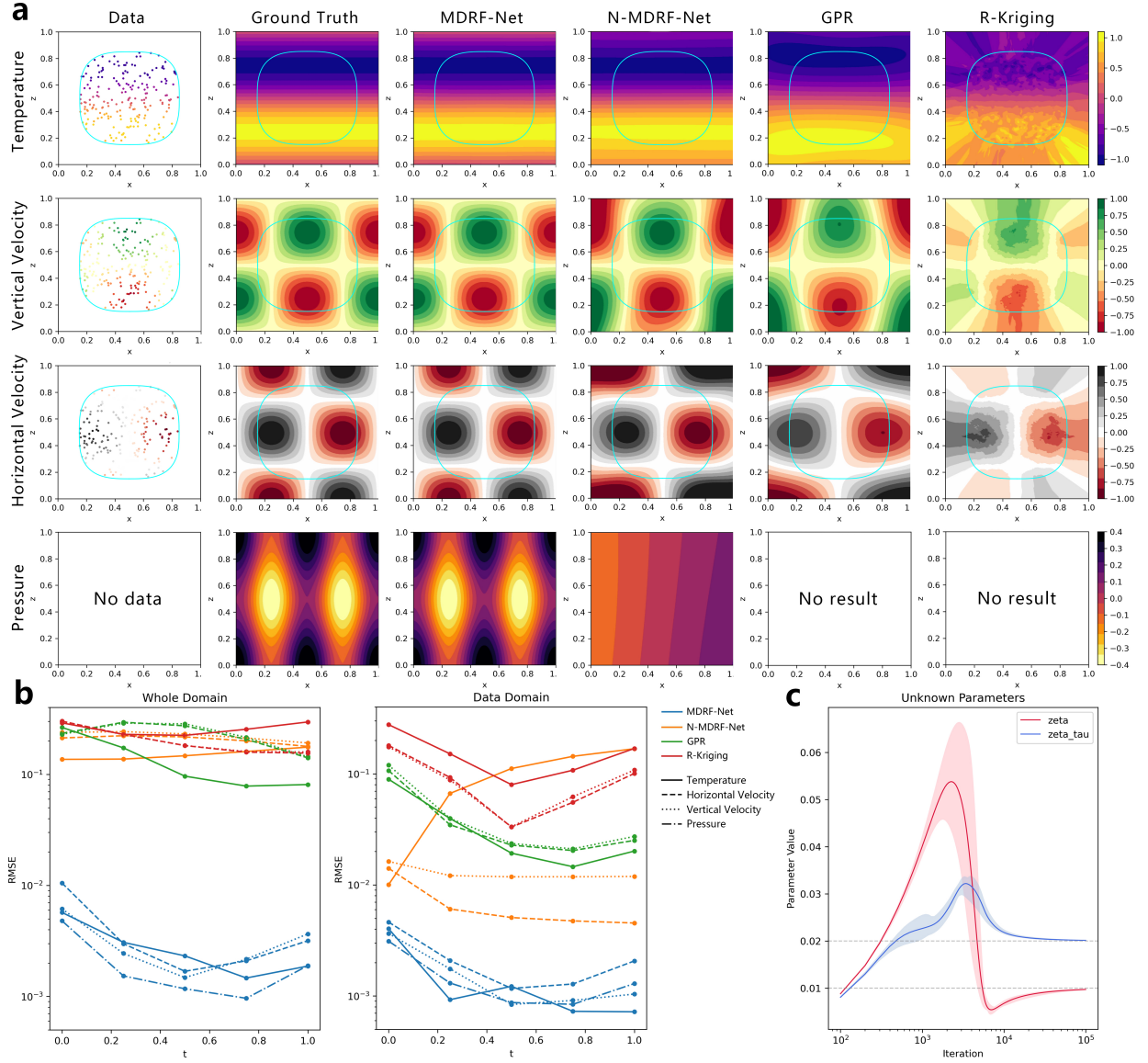


Figure 2: **Simulation study.** MDRF-Net, MDRF-Net without mechanism (N-MDRF-Net) are compared with Gaussian Process Regression (GPR) and Regression Kriging (R-Kriging) at $t=0$ (a, more results can be found in Appendix E). Comparison of the accuracy of competing methods over time in the whole domain and the data domain, the accuracies of the pressure fields except for MDRF-Net are not plotted (b). The data are generated within rounded rectangles and do not contain the pressure variables. Changes of the unknown parameters during training, the solid lines are the means of the results of 100 repetitions of the experiment, while the light-colored areas in the background are the point-wise 95% confidence intervals (c).

3.2 Assessment of MDRF-Net with Real Data

We assess the performance of MDRF-Net in the equatorial Pacific (Appendix B) and analyze error trends over space and time (Figures 3 **a,b,c**; additional errors in Appendix F). These trends prove steady, suggesting MDRF-Net effectively captures spatial and temporal patterns. Errors near islands are marginally higher but show a distinct peak at depths of 250-350m for all variables except vertical velocity, which follows a separate pattern (see Appendix F). When tested against Argo and Copernicus data, MDRF-Net yields satisfactory RMSEs: 0.358°C for temperature, 0.0474 psu for salinity, and for velocities- 1.955×10^{-5} m/s (vertical), 0.0465 m/s (northward), and 0.0513 m/s (eastward), demonstrating its competence in predicting ocean conditions.

MDRF-Net is compared with the same competing models as the simulation study. We compute the RMSEs of all models utilizing training sets of different sizes on the test sets of the original spatiotemporal domain, two extended spatial domains (outwardly stretching and inwardly filling), and an extended temporal domain (forecasting). MDRF-Net excels in all evaluations, particularly in 3-month forecasting (Figure 3 **d iv**), outperforming traditional statistical inversing methods. It significantly outmatches GPR and R-Kriging for temperature and salinity, and dominates in flow velocity predictions, especially with larger datasets. As data volume increases, MDRF-Net’s errors decrease sharply. Its incorporation of primitive equations further boosts accuracy in inferring unobserved pressure and density fields, highlighting the equations’ critical role in augmenting the model’s performance with additional physics-driven insights.

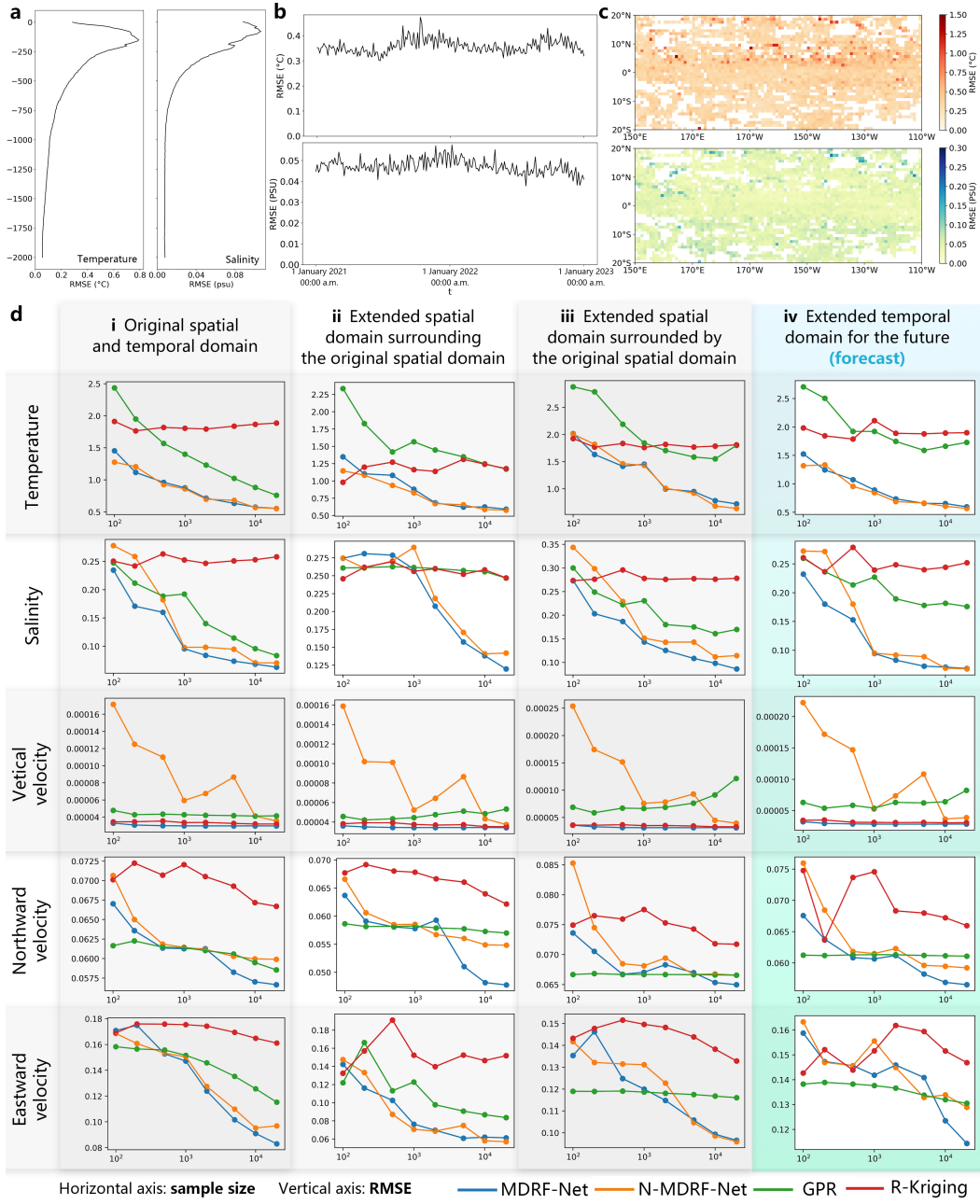


Figure 3: **The trends of root-mean-square errors (RMSEs) with the test set.** The variability of RMSEs is shown over depth (a), time (b), and coordinates (c) for temperature and salinity. The blanks in (c) represent islands. The RMSE charts for all marine variables over depth, time and coordinates can be found in Appendix F. Comparisons (d) are made to evaluate the performance of MDRF-Net in relation to its version that does not include the primitive equations (N-MDRF-Net), as well as the Gaussian process regression (GPR) and Regression Kriging algorithms (R-Kriging). Panel (i) represents the inversion accuracy of the models on the original spatiotemporal scale, while panels (ii), (iii) depict the prediction effects on the extended spatial scales (outwardly stretching and inwardly filling), respectively. Panel (iv) shows the 3-month average forecasting errors. scale All error values presented are averaged over five replicated experiments.

3.3 Reconstructing Global Oceanic Variations over Space Continuously

In order to characterize global marine activity, we have expanded the application of MDRF-Net to the global ocean, and local scenario of equatorial Pacific is presented in Appendix G. To enhance MDRF-Net’s performance along the coastlines, we have implemented Neumann boundary conditions for temperature and salinity, along with Dirichlet boundary conditions for current flow (Equation 3c). Despite slightly lower accuracy compared to smaller sea areas, MDRF-Net delivers impressive results, even enabling the inversion of two variable fields without collected data. The overall RMSE values for temperature, salinity, vertical velocity, northward velocity, and eastward velocity fields are 0.455 °C, 0.0714 psu, 4.254×10^{-6} m/s, 0.0777 m/s, and 0.0825 m/s respectively.

The inversion results for the five variable fields with available observations show good performance. The temperature, salinity, eastward velocity, and northward velocity fields exhibit clear patterns (Figure 4 **a,b,d,e**). For example, the temperature field decreases with increasing latitude at the surface and remains relatively constant at deeper levels (Figure 4 **a**). The salinity field shows extreme values in specific regions such as the Mediterranean Sea and the Black Sea, as well as subsurface occurrences in the Arctic Ocean (Figure 4 **b**). The three current fields display higher currents near the equator and changes corresponding to land formations (Figure 4 **c,d,e**). For instance, the northward currents (Figure 4 **d**) are stronger on the east side of the continent than on the west side. Regarding the pressure and density fields (Figure 4 **f,g**) without available observed values, they exhibit expected patterns consistent with local scenarios. The pressure field correlates primarily with water depth, increasing with greater depths. The density field, on the other hand, shows a maximum and minimum in a region similar to that of the salinity field.

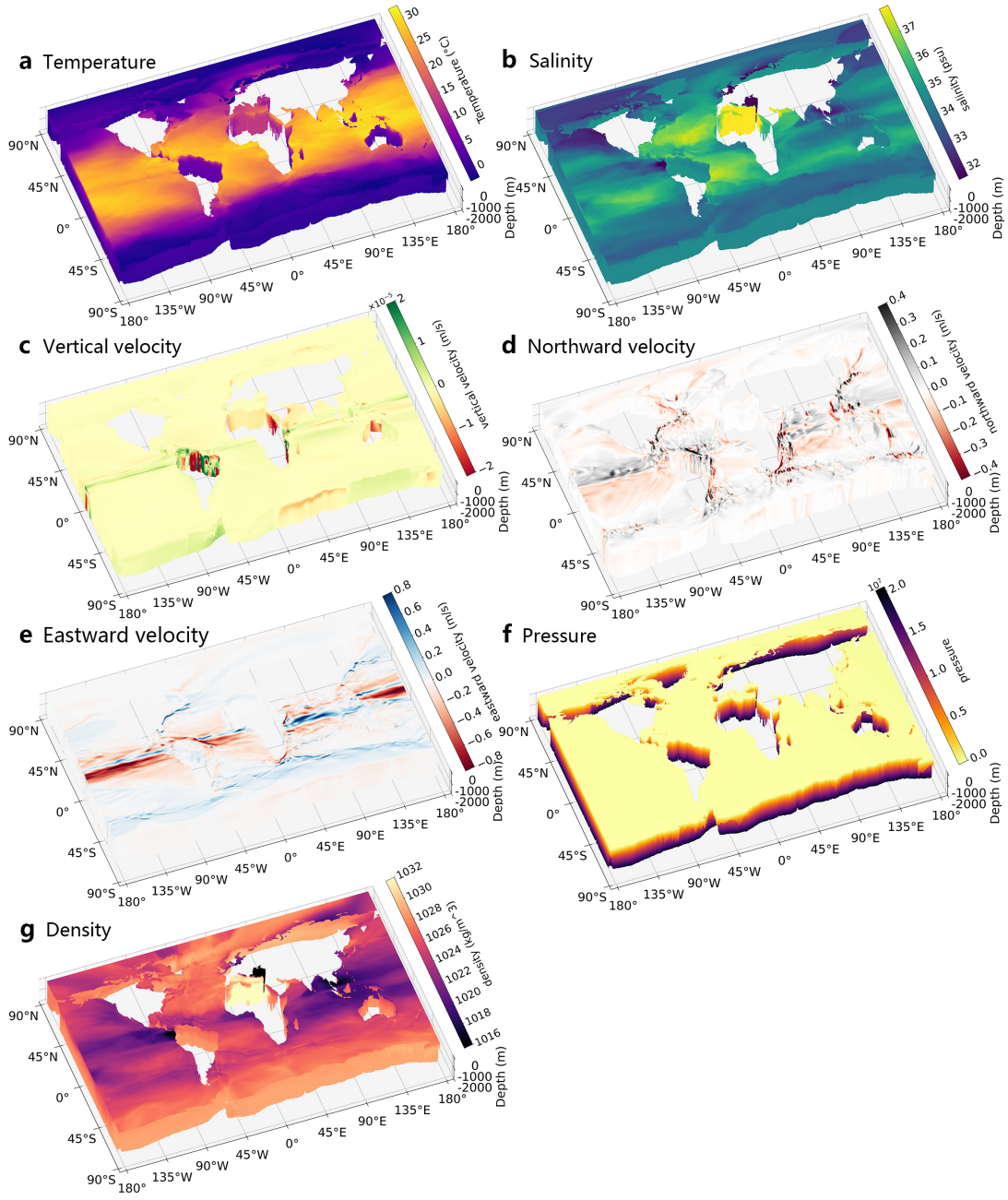


Figure 4: **Global inversion performances of temperature (a), salinity (b), vertical velocity (c), northward velocity (d), eastward velocity (e), pressure (f) and density (g) fields at 16 January 2021, 12:00 a.m.** The pressure (f) and density (g) fields are inverted without observed data. Results for more time spots can be found in Appendix G, while the full performances for 2021 and 2022 can be found in the R Shiny platform (<https://tikitakatikitaka.shinyapps.io/mdrf-net-shiny/>).

From the direct comparison with the reanalyzed fields, it is evident that MDRF-Net

is able to invert temperature and salinity fields that closely resemble their reanalyzed counterparts, despite not being directly trained using the reanalyzed data (Figure 5 **a,b,c**). This is achieved through the backfeeding of the exact current fields using the bridging of the primitive equations. Additionally, our variable fields demonstrate continuity in both space and time. For instance, in Figure 5 **b,c**, the reanalysis field shows a discontinuity in water depth and clear stratification. It should be noted that Copernicus’ reanalysis fields only include 40 layers up to a depth of 2000 meters.

The comparison of the density field illustrates the greater advantage of MDRF-Net (Figure 5 **d**), as it is not based on collected data like the pressure field. Instead, MDRF-Net inverts the density field using observations of other variable fields and the original equations exclusively, particularly through the equation of state. While we do not have an actual representation of the pressure field, and even the density field in the Copernicus reanalysis dataset only contains data from the ocean surface, MDRF-Net is capable of inverting the density field for the three-dimensional(3D) ocean and shows similarity to the reanalysis field at the ocean surface. The density field shows better results than the local scenario in the surface density versus the reanalyzed field (Figure 5 **e**), with a very elegant match of patterns. Furthermore, we showcase the inversion performances of the temperature and salinity fields of MDRF-Net at the Arctic zone (Figure 5 **f**), and there are still some differences in the Arctic Ocean in Northern Canada, but when compared to the reanalysis data, the pattern of the MDRF-Net inversion results is roughly the same.

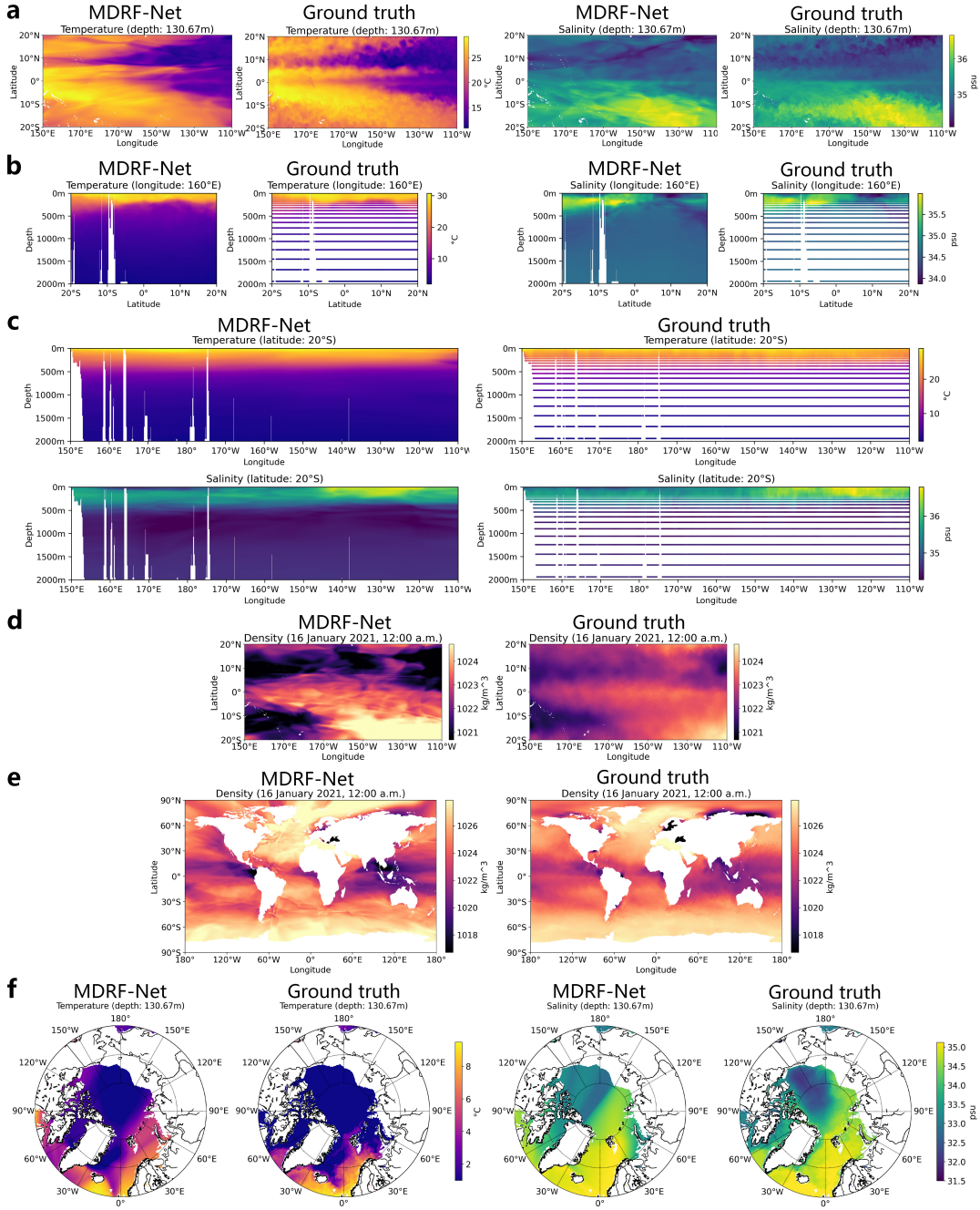


Figure 5: The comparison between the inversion performances of the temperature, salinity and density fields of MDRF-Net and the reanalyzed data from Copernicus, labeled as ground truth. Comparisons are conducted at the depth (a), longitude (b), and latitude cross-sections (c). The blank areas in the images represent islands. Note that the layering seen in the Copernicus images is a result of the dataset's sparse coverage of depth. Within the range of 2000 meters, only 40 layers are available for the reanalyzed data and there is only sea surface density data available. Panels (d) and (e) depict the density fields for the local and global scenarios respectively. The performance of MDRF-Net for temperature and salinity fields at the Arctic zone are shown (f). The more comparison can be found in Appendix G.

3.4 Forecasting Global Oceanic Variations over Time Gaplessly

In addition to projecting variable fields onto a gapless space and time scale, MDRF-Net can deliver uninterrupted forecasts. We give short- (1 and 7 days) and long-term (30 days) projections of MDRF-Net and compare them against reanalysis data for the same time-frame. In the short-term predictions (Appendix G), the MDRF-Net results show a high degree of agreement with the reanalysis variable fields, especially the temperature and salinity fields, while the 3D current field appears smoother than the reanalysis field, but the main current patterns are still well recognized by MDRF-Net. For long-term forecasting (Figure 6 **a**), it can be observed that most forecast results closely align with the reanalysis data across different depths with minor discrepancies observed in some land-marginal seas like the Black Sea and the Red Sea. Notably, MDRF-Net accurately captures features such as a warmer zone in the western Atlantic Ocean at a depth of 380.213 meters and the spread of high salinity from the Mediterranean Sea into the Atlantic Ocean at a depth of 1,062.440 meters.

Quantitatively, MDRF-Net demonstrates remarkable accuracy and consistent stability in forecasting, evident from the two-dimensional distributions of absolute errors (Figure 6 **b**). These visualizations highlight that a large portion of prediction errors in the test set cluster closely around zero. Furthermore, in trend analysis, MDRF-Net displays exceptional consistency. Over a one-month period, there is a slight observable upward trend, but this minor deviation does not detract from its overall robust and steady performance that remains consistent throughout the evaluation period.

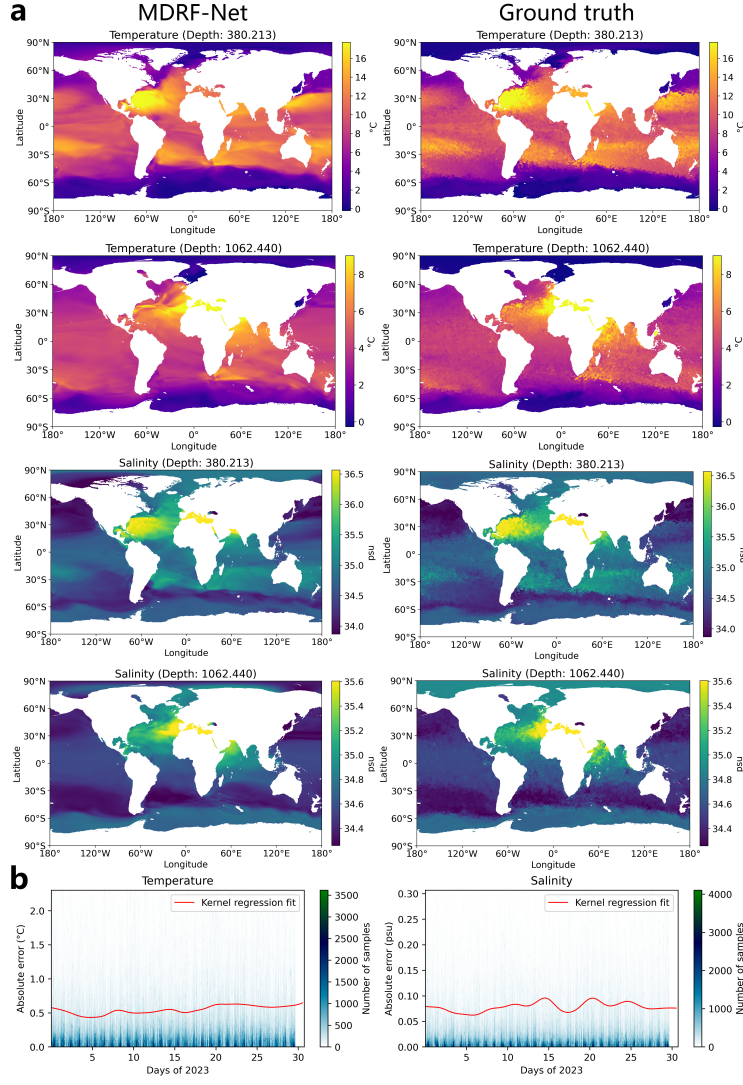


Figure 6: **Long-term forecasting.** 30 days forecasting (30 January 2023, 12:00 a.m.) results of temperature and salinity are compared with the corresponding reanalysis data (a). More forecast time spans are shown in Appendix G. The absolute errors of 1-month forecasting are calculated with the Argo data (b). The red line is the kernel regression result of 5000 samples from all the data, with bandwidth 1.

4 Discussion

MDRF-Net addresses a core AI problem in contemporary oceanography by integrating neural networks with ocean dynamics' fundamental laws. With its ensemble method, shared-layer full connectivity, and two-step training process, MDRF-Net contributes to statistical

methodology, theory, and application, demonstrating the potential of statistics in empowering AI tools to bridge data-driven models with physical oceanography. This innovative approach employs statistical theories and methodologies to offer unique insights into marine phenomena, advancing real-world scientific understanding.

MDRF-Net excels in integrating multiple datasets seamlessly, fusing Argo observations with Copernicus reanalysis data for a robust inversion base. It infers seven oceanic variables from five observed ones, deepening our insights. Its design promotes easy data source integration, adapting to varied sample sizes and ranges. Utilizing 2+ million Argo profiles and hundreds of millions of reanalysis points, MDRF-Net’s multi-source data fusion has the potential to significantly improve the generalization capability and accuracy of forecasts in large-scale ocean (Guillou et al., 2020; Adhikary and Banerjee, 2023).

MDRF-Net stands out for its remarkable extensibility, adeptly managing variable marine data at scale. It inverts two uncollected fields (density, pressure) from five known variables and extends this capability to poorly monitored areas like the Arctic and deep sea (Ura, 2013; Charrassin et al., 2008). By facilitating future change predictions, MDRF-Net supports proactive decision-making and resource management (Ellefmo et al., 2023), with its swift forecasting prowess making it a powerful tool for anticipating ocean dynamics.

MDRF-Net’s generalization error estimation reveals its effectiveness in inversely solving the primitive equation for ocean variables, meeting three key criteria: a well-trained model, indicated by low training error; appropriate regularization enabled by the smooth tanh activation, ensuring accurate residual approximation via quadrature; fulfillment of the conditional stability assumption for the inverse problem. Consequently, MDRF-Net successfully approximates the inverse problem and yields precise reconstructions and forecasts of ocean dynamics. As for those methods that do not include physics mechanisms

(N-MDRF-Net, GPR, R-Kriging), there are no upper bounds on their generalization errors, and their theoretical convergence is not guaranteed. We also show that the ensemble method of rotating the earth coordinates yields lower generalisation errors and improves the model’s performance in polar regions (Appendix Figure 5).

Simulation studies indicate that MDRF-Net excels in reconstructing both partial and complete variable fields with missing data, consistently delivering higher accuracy than other methods, even in data-rich domains. This aligns with the theoretical results, specifically, in regions outside the rounded rectangular area where data is absent (Figure 2), other methods lacking mechanistic insights fail to make effective predictions, thereby resulting in unbounded generalization errors in those locales. Meanwhile, MDRF-Net is stable and recognizable for the unknown parameter positional parameters in the dynamical system of the simulation experiment.

MDRF-Net distinguishes itself from both AI-based (Xiao et al., 2019; Xie et al., 2019; Song et al., 2020; Su et al., 2021; Song et al., 2022; Zhang et al., 2023) and statistical inversion models (Lee et al., 2019; Yarger et al., 2022) by effectively handling diverse sampling locations and offering continuous results. Its adaptability to sparse spatiotemporal data, coupled with natural interpretability and structural simplicity, facilitates seamless data integration, achieving higher accuracy and efficiency. MDRF-Net thus emerges as a superior solution, precisely capturing marine dynamics compared to other existing methodologies (Figure 3 **d**).

Few methods can efficiently integrate and invert crucial variables across the global ocean in a continuous spatiotemporal manner. Upon comparing our MDRF-Net results with re-analyzed data (Figure 5), it is evident that MDRF-Net’s consistently superior performance in global scenarios positions it as a formidable tool for conducting comprehensive marine

studies. The continuous provision of global inversion marine fields allows for a nuanced understanding of long-term trends and patterns, offering a unique advantage in the field. MDRF-Net excels in making reasonable inferences about ocean changes that exhibit continuity both spatially and temporally.

Forecasting over extended periods poses a major challenge, especially for global, high-resolution demands requiring substantial computational efficiency (Wolff et al., 2020; O’Donncha et al., 2015). MDRF-Net addresses this by delivering high-precision forecasts for multiple climate variables, excelling in both short- and long-term projections. Comparative analysis with reanalysis data reveals MDRF-Net’s exceptional performance in extrapolating forecasts over time. Notably, MDRF-Net exhibits remarkable accuracy in temperature and salinity fields, with extrapolation errors only slightly higher than those of interpolation. Ocean velocity forecasts benefit from incorporated mechanics, yielding smoother outputs than the turbulent reanalysis data (Appendix G). Crucially, MDRF-Net’s gapless predictions for various variables across the globe and time span do not heavily tax computational resources, advancing four-dimensional oceanic forecasting with enhanced resolution compared to current high-resolution forecasting methods (Wolff et al., 2020; Zhang et al., 2023).

MDRF-Net has broad real-world implications, uncovering previously unseen patterns. Its reconstruction of 4D oceanic fields enhances understanding and guides marine activities, vital for studying impacts on marine life, nutrients, and environment stability (Behrenfeld and Falkowski, 1997; Anderson et al., 2021; Sarmiento, 2006). By elucidating events like the Mediterranean Salinity Crisis (Hsü et al., 1977; Ryan, 2009) and the North Atlantic Warm Current (Schmitz Jr and McCartney, 1993; Jiang et al., 2021; Gunnarsson, 2021), MDRF-Net showcases historical and climatic insight. Its practical applications span fisheries management (Carnevale et al., 2022), maritime navigation (Melia et al., 2016; Kelly

et al., 2020), and offshore infrastructure resilience (El-Khoury et al., 2022), while also finely monitoring tropical phenomena such as El Niño (Cai et al., 2021), highlighting versatility and practical significance.

MDRF-Net’s mesh-free methodology addresses a pivotal challenge in marine scientific AI, amplifying its effectiveness through statistical insights. This model, which addresses limitations in existing methods and difficulties of variable field inversion, showcases the unique role of statistical theories in improving AI methods. While external forces’ impact on shallow coastal areas is acknowledged for future refinement, MDRF-Net has demonstrated excellent overall performance and generalization capabilities. MDRF-Net’s innovative application in sustainable marine resource management, underlines the significant contribution of statistical methodologies to AI, offering solutions for real-world science advancement.

Data and Code Availability

We select worldwide temperature and salinity data collected by the lifting floats from the Argo program (<https://argo.ucsd.edu/>), and eastward, northward and vertical reanalysis seawater velocity data from the EU Copernicus Marine Service (<https://marine.copernicus.eu/>) for the two years 2021 to 2022 and water depths up to 2000 m.

The code for this article has been uploaded to Github at the link (<https://github.com/tikitaka0243/mdrf-net>). We utilized the Python library DeepXDE (Lu et al., 2021) to construct the MDRF-Net.

Acknowledgments

This research was supported by the National Key R&D Program of China (No. 2022YFA1003800), the National Natural Science Foundation of China (No. 71991474; No. 12001554; No. 72171216), the Key Research and Development Program of Guangdong, China (No. 2019B020228001), the Science and Technology Program of Guangzhou, China (No. 202201011578), and the Natural Science Foundation of Guangdong Province, China (No. 2021A1515010205). The funding agencies had no role in the study design, data collection, analysis, decision to publish, or preparation of the manuscript.

Disclosure Statement

The authors report there are no competing interests to declare.

SUPPLEMENTARY MATERIAL

Title: Supplementary Material for “Reconstructing and Forecasting Marine Dynamic Variable Fields across Space and Time Globally and Gaplessly”

Overview: The derivation of the primitive equations is provided in **Appendix A**. More details about the Argo and Copernicus data sets are provided in **Appendix B**. More details about the Marine Dynamic Reconstruction and Forecast neural network (MDRF-Net) are provided in **Appendix C**. More derivations of Generalization Error Estimation are provided in **Appendix D**. More results of simulation study are provided in **Appendix E**. The error analysis of MDRF-Net is provided in **Appendix F**. Additional reconstruction and forecast results are provided in **Appendix G**.

References

- Adhikary, S. and S. Banerjee (2023). Improved large-scale ocean wave dynamics remote monitoring based on big data analytics and reanalyzed remote sensing. *Nature Environment & Pollution Technology* 22(1).
- Anderson, D. M., E. Fensin, C. J. Gobler, A. E. Hoeglund, K. A. Hubbard, D. M. Kulis, J. H. Landsberg, K. A. Lefebvre, P. Provoost, M. L. Richlen, et al. (2021). Marine harmful algal blooms (habs) in the united states: History, current status and future trends. *Harmful Algae* 102, 101975.
- Ashton, G. V., A. L. Freestone, J. E. Duffy, M. E. Torchin, B. J. Sewall, B. Tracy, M. Albano, A. H. Altieri, L. Altvater, R. Bastida-Zavala, et al. (2022). Predator control of marine communities increases with temperature across 115 degrees of latitude. *Science* 376(6598), 1215–1219.
- Bauer, P., A. Thorpe, and G. Brunet (2015). The quiet revolution of numerical weather prediction. *Nature* 525(7567), 47–55.
- Behrenfeld, M. J. and P. G. Falkowski (1997). Photosynthetic rates derived from satellite-based chlorophyll concentration. *Limnology and oceanography* 42(1), 1–20.
- Cai, W., A. Santoso, M. Collins, B. Dewitte, C. Karamperidou, J.-S. Kug, M. Lengaigne, M. J. McPhaden, M. F. Stuecker, A. S. Taschetto, et al. (2021). Changing el niño–southern oscillation in a warming climate. *Nature Reviews Earth & Environment* 2(9), 628–644.
- Carnevale, G., S. Werner, et al. (2022). Marine life in the mediterranean during the messinian salinity crisis: a paleoichthyological perspective. *Rivista Italiana di Paleontologia e Stratigrafia* 128, 283–324.
- Charrassin, J.-B., M. Hindell, S. R. Rintoul, F. Roquet, S. Sokolov, M. Biuw, D. Costa, L. Boehme, P. Lovell, R. Coleman, et al. (2008). Southern ocean frontal structure and sea-ice formation rates revealed by elephant seals. *Proceedings of the National Academy of Sciences* 105(33), 11634–11639.
- Charve, F. (2008). Global well-posedness for the primitive equations with less regular initial data. In *Annales de la Faculté des sciences de Toulouse: Mathématiques*, Volume 17, pp. 221–238.
- El-Khoury, M., E. Roziere, F. Grondin, R. Cortas, and F. H. Chehade (2022). Experimental evaluation of the effect of cement type and seawater salinity on concrete offshore structures. *Construction and Building Materials* 322, 126471.

- Ellefmio, S. L., N. Aberle, V. Hagspiel, M. Ingulstad, and K. Aasly (2023). Marine minerals’ role in future holistic mineral resource management. *Geological Society, London, Special Publications* 526(1), SP526–2022.
- European Union-Copernicus Marine Service (2020). Global ocean 1/12° physics analysis and forecast updated daily.
- Fox-Kemper, B., A. Adcroft, C. W. Böning, E. P. Chassignet, E. Curchitser, G. Danabasoglu, C. Eden, M. H. England, R. Gerdes, R. J. Greatbatch, et al. (2019). Challenges and prospects in ocean circulation models. *Frontiers in Marine Science* 6, 65.
- Guillou, N., G. Lavidas, and G. Chapalain (2020). Wave energy resource assessment for exploitation—a review. *Journal of Marine Science and Engineering* 8(9), 705.
- Gunnarsson, B. (2021). Recent ship traffic and developing shipping trends on the northern sea route—policy implications for future arctic shipping. *Marine Policy* 124, 104369.
- Harley, C. D., A. Randall Hughes, K. M. Hultgren, B. G. Miner, C. J. Sorte, C. S. Thornber, L. F. Rodriguez, L. Tomanek, and S. L. Williams (2006). The impacts of climate change in coastal marine systems. *Ecology letters* 9(2), 228–241.
- Hieber, M., A. Hussein, and T. Kashiwabara (2016). Global strong lp well-posedness of the 3d primitive equations with heat and salinity diffusion. *Journal of Differential Equations* 261(12), 6950–6981.
- Hsü, K. J., L. Montadert, D. Bernoulli, M. B. Cita, A. Erickson, R. E. Garrison, R. B. Kidd, F. Mèlierès, C. Müller, and R. Wright (1977). History of the mediterranean salinity crisis. *Nature* 267(5610), 399–403.
- Hu, R., Q. Lin, A. Raydan, and S. Tang (2023). Higher-order error estimates for physics-informed neural networks approximating the primitive equations.
- Jiang, W., G. Gastineau, and F. Codron (2021). Multicentennial variability driven by salinity exchanges between the atlantic and the arctic ocean in a coupled climate model. *Journal of Advances in Modeling Earth Systems* 13(3), e2020MS002366.
- Johnson, G. C., S. Hosoda, S. R. Jayne, P. R. Oke, S. C. Riser, D. Roemmich, T. Suga, V. Thierry, S. E. Wijffels, and J. Xu (2022). Argo—two decades: Global oceanography, revolutionized. *Annual review of marine science* 14, 379–403.
- Kelly, S., E. Popova, Y. Aksenov, R. Marsh, and A. Yool (2020). They came from the pacific: How changing arctic currents could contribute to an ecological regime shift in the atlantic ocean. *Earth’s Future* 8(4), e2019EF001394.

- Lee, K. M. B., C. Yoo, B. Hollings, S. Anstee, S. Huang, and R. Fitch (2019). Online estimation of ocean current from sparse gps data for underwater vehicles. In *2019 International conference on robotics and automation (ICRA)*, pp. 3443–3449. IEEE.
- Lions, J. L., R. Temam, and S. Wang (1992, sep). On the equations of the large-scale ocean. *Nonlinearity* 5(5), 1007.
- Lou, R., Z. Lv, S. Dang, T. Su, and X. Li (2023). Application of machine learning in ocean data. *Multimedia Systems* 29(3), 1815–1824.
- Lu, L., X. Meng, Z. Mao, and G. E. Karniadakis (2021). Deepxde: A deep learning library for solving differential equations. *SIAM review* 63(1), 208–228.
- Melia, N., K. Haines, and E. Hawkins (2016). Sea ice decline and 21st century trans-arctic shipping routes. *Geophysical Research Letters* 43(18), 9720–9728.
- Mishra, S. and R. Molinaro (2021, 06). Estimates on the generalization error of physics-informed neural networks for approximating a class of inverse problems for PDEs. *IMA Journal of Numerical Analysis* 42(2), 981–1022.
- O’Donncha, F., M. Hartnett, S. Nash, L. Ren, and E. Ragnoli (2015). Characterizing observed circulation patterns within a bay using hf radar and numerical model simulations. *Journal of Marine Systems* 142, 96–110.
- Raissi, M., P. Perdikaris, and G. E. Karniadakis (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics* 378, 686–707.
- Ryan, W. B. (2009). Decoding the mediterranean salinity crisis. *Sedimentology* 56(1), 95–136.
- Salcedo-Sanz, S., P. Ghamisi, M. Piles, M. Werner, L. Cuadra, A. Moreno-Martínez, E. Izquierdo-Verdiguier, J. Muñoz-Marí, A. Mosavi, and G. Camps-Valls (2020). Machine learning information fusion in earth observation: A comprehensive review of methods, applications and data sources. *Information Fusion* 63, 256–272.
- Sarmiento, J. L. (2006). *Ocean biogeochemical dynamics*. Princeton university press.
- Schmitz Jr, W. J. and M. S. McCartney (1993). On the north atlantic circulation. *Reviews of Geophysics* 31(1), 29–49.
- Smyth, K. and M. Elliott (2016). Effects of changing salinity on the ecology of the marine environment. *Stressors in the marine environment*, 161–174.

- Song, T., Z. Wang, P. Xie, N. Han, J. Jiang, and D. Xu (2020). A novel dual path gated recurrent unit model for sea surface salinity prediction. *Journal of Atmospheric and Oceanic Technology* 37(2), 317–325.
- Song, T., W. Wei, F. Meng, J. Wang, R. Han, and D. Xu (2022). Inversion of ocean subsurface temperature and salinity fields based on spatio-temporal correlation. *Remote Sensing* 14(11), 2587.
- Su, H., A. Wang, T. Zhang, T. Qin, X. Du, and X.-H. Yan (2021). Super-resolution of subsurface temperature field from remote sensing observations based on machine learning. *International Journal of Applied Earth Observation and Geoinformation* 102, 102440.
- Tittensor, D. P., C. Novaglio, C. S. Harrison, R. F. Heneghan, N. Barrier, D. Bianchi, L. Bopp, A. Bryndum-Buchholz, G. L. Britten, M. Büchner, et al. (2021). Next-generation ensemble projections reveal higher climate risks for marine ecosystems. *Nature Climate Change* 11(11), 973–981.
- Ura, T. (2013). Observation of deep seafloor by autonomous underwater vehicle.
- Wolff, S., F. O’Donncha, and B. Chen (2020). Statistical and machine learning ensemble modelling to forecast sea surface temperature. *Journal of Marine Systems* 208, 103347.
- Wong, A. P., S. E. Wijffels, S. C. Riser, S. Pouliquen, S. Hosoda, D. Roemmich, J. Gilson, G. C. Johnson, K. Martini, D. J. Murphy, et al. (2020). Argo data 1999–2019: Two million temperature-salinity profiles and subsurface velocity observations from a global array of profiling floats. *Frontiers in Marine Science* 7, 700.
- Xiao, C., N. Chen, C. Hu, K. Wang, Z. Xu, Y. Cai, L. Xu, Z. Chen, and J. Gong (2019). A spatiotemporal deep learning model for sea surface temperature field prediction using time-series satellite data. *Environmental Modelling & Software* 120, 104502.
- Xie, J., J. Zhang, J. Yu, and L. Xu (2019). An adaptive scale sea surface temperature predicting method based on deep learning with attention mechanism. *IEEE Geoscience and Remote Sensing Letters* 17(5), 740–744.
- Yarger, D., S. Stoev, and T. Hsing (2022). A functional-data approach to the argo data. *The Annals of Applied Statistics* 16(1), 216–246.
- Zhang, X., N. Zhao, and Z. Han (2023). A modified u-net model for predicting the sea surface salinity over the western pacific ocean. *Remote Sensing* 15(6), 1684.
- Zhang, Y., H. Mo, Z. Zu, and Y. Qin (2023). Preliminary validation for an eddy-resolving global ocean forecasting system–nmefc-nemo. In *Journal of Physics: Conference Series*, Volume 2486, pp. 012030. IOP Publishing.