



ParNeC

模式识别与神经计算研究组
PAttern Recognition and NEural Computing

Physics-Coupled Spatio-Temporal

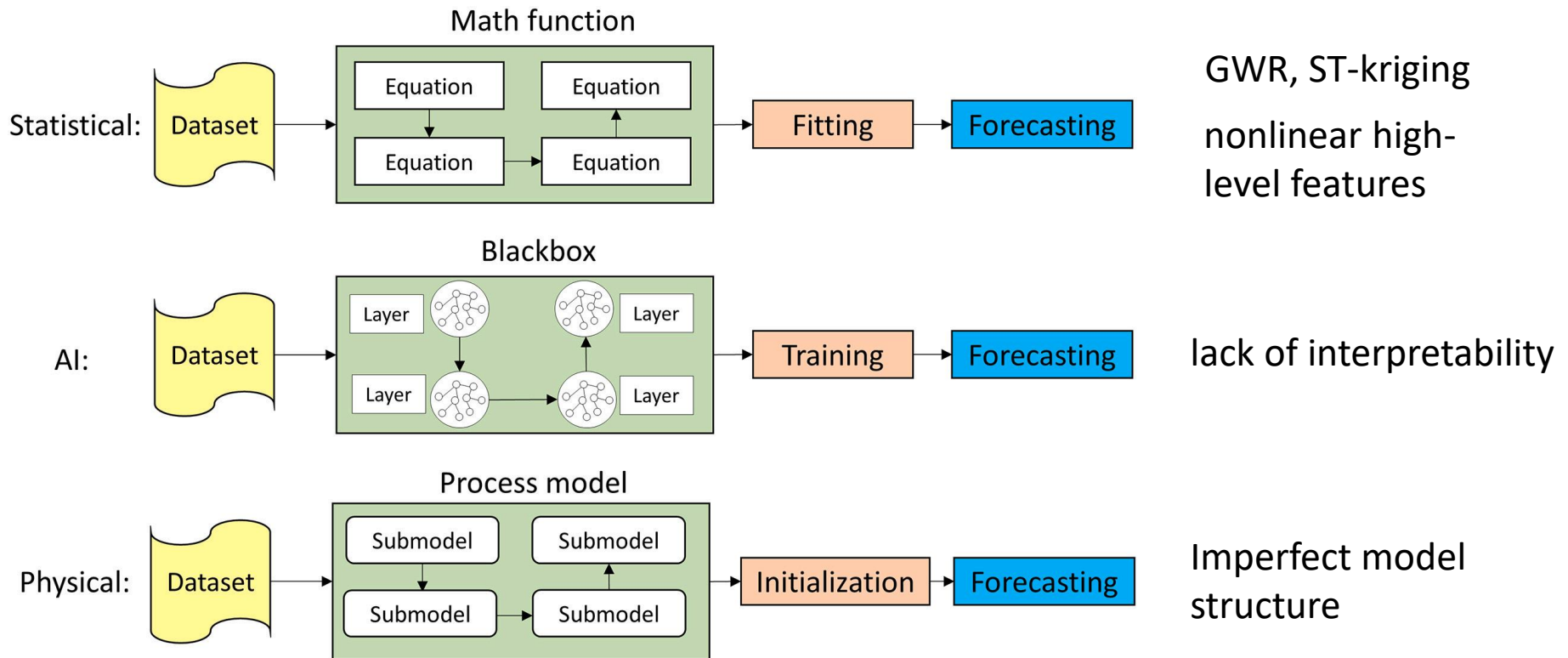
Active Learning for Dynamical Systems

Yu Huang, Yufei Tang, Xingquan Zhu, Min Shi, Ali Muhamed Ali, Hanqi Zhuang, Laurent Cherubin

{yhwang2018,tangy,xzhu3,mshi2018,amuhamedali2014,zhuang,lcherubin}@fau.edu

Florida Atlantic University
Boca Raton, Florida, USA

extends traditional time series forecasting or spatial interpolation problem to space and time dimensions.



- high cost of data acquisition

ocean data:

- need trained scientists
- cost limits the total number of sensors
- increasing costs does not necessarily mean improved model predictions



where and when to query
the data within constraints



AL

- infer the underlying causes

how data is generated

how data is propagated



homogeneity and heterogeneity



driven by the physical laws



impacted by the spatial
and temporal regions

- propose an active learning algorithm for spatial-temporal dynamical systems.
- propose ST-PCNN, which captures **heterogeneity** from all spatial locations.
- propose embedding a physics learning module to learn the inherited **homogeneity**.
- conduct extensive experiments and comparative studies on both synthetic and real-world datasets.

while *Termination Conditions NOT satisfied do*

Learn Physics: $[\lambda] \leftarrow$ learn the physics, *i.e.* coefficients λ of PDE, from the existing training data \mathbb{D} by minimizing Eq. (32);

GP

ST-PCNN Training: train the ST-PCNN with the learned physics λ from training data \mathbb{D} , see Algorithm 2;

Pos + LSTM

Prediction: $[\hat{\mathcal{S}}] \leftarrow$ make prediction at all locations in the network, $\forall \mathbf{x} \in \Omega$;

Kriging: $\Omega_{kriging}^n \leftarrow$ use Kriging to identify n query locations with the largest estimated errors for next active learning step;

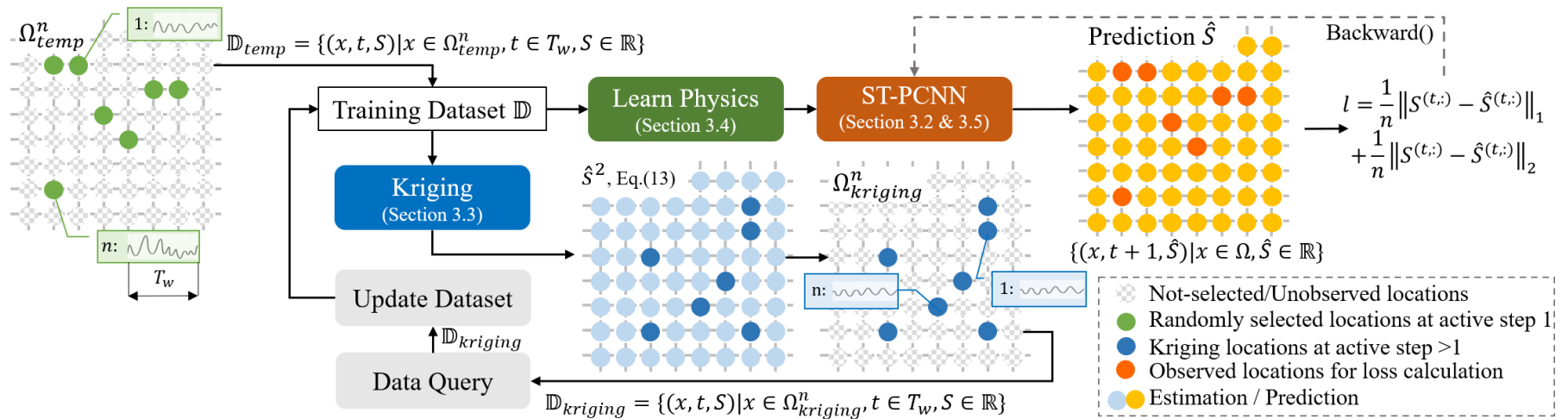
Kriging(GPR)

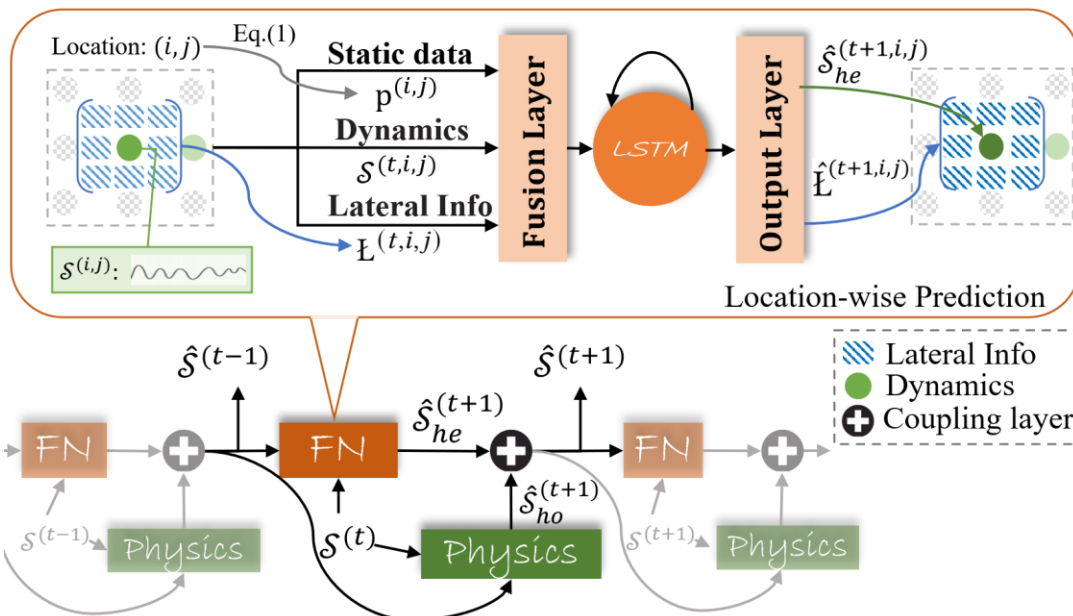
Data Query: obtain consecutive observations of window size T_w from the newly selected locations

$\mathbb{D}_{kriging} = \{(\mathbf{x}, t, \mathcal{S}) \mid \mathbf{x} \in \Omega_{kriging}^n, t \in T_w, \mathcal{S} \in \mathbb{R}\}$;

Update Dataset: add $\mathbb{D}_{kriging}$ to the existing training data \mathbb{D} ;

end

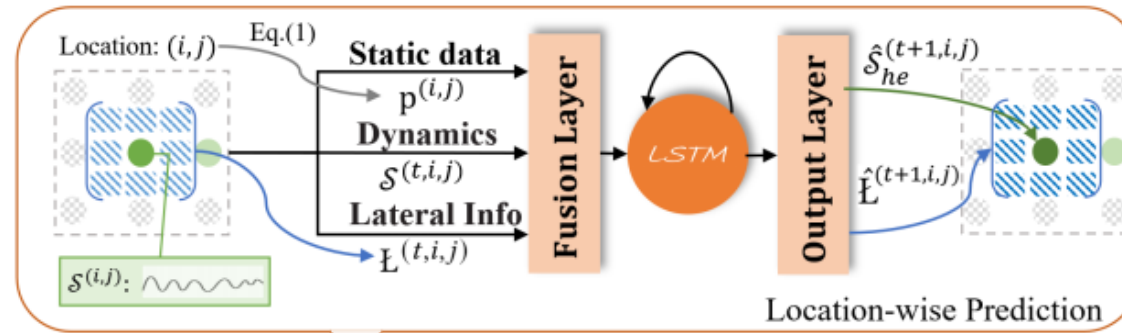




bi-network architecture:

- forecasting network (FN): produces heterogeneous prediction leveraging its own specific local attributes only
 - physics network (PN): generates homogeneous solution of the dynamics regularized by the overall underlain governing physics
- produce the final prediction by synthesizing heterogeneous prediction and homogeneous prediction

$$\hat{\mathcal{S}}^{(t+1,i,j)} = Relu([\hat{\mathcal{S}}_{he}^{(t+1,i,j)}, \hat{\mathcal{S}}_{ho}^{(t+1,i,j)}] \mathcal{W}_C^T + b_C)$$



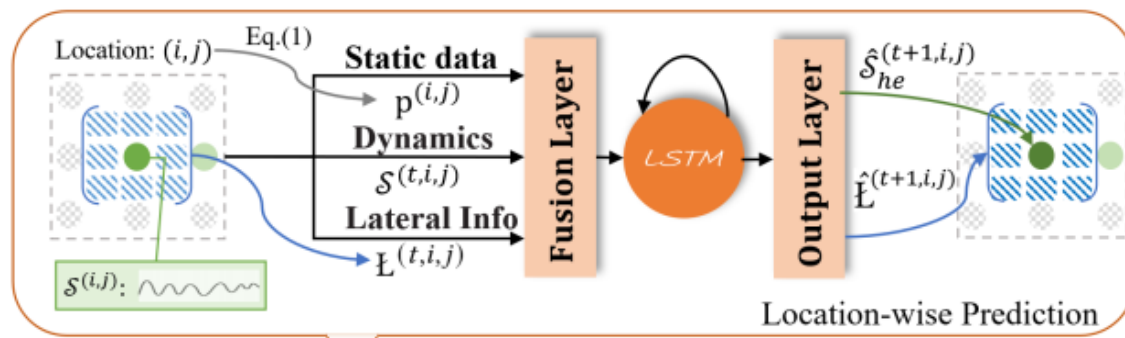
Static data: representing node location

Dynamics: representing dynamics of node
(temperature, conductivity and velocity)

Lateral info: capturing interaction between each
node and its neighbors

Predicted dynamics

Additional lateral information



FN encodes each view using a fusion layer



features are then fed into an LSTM to model the node-specific interactions



An output layer is stacked to transform the LSTM output into the expected dynamic prediction and additional lateral information

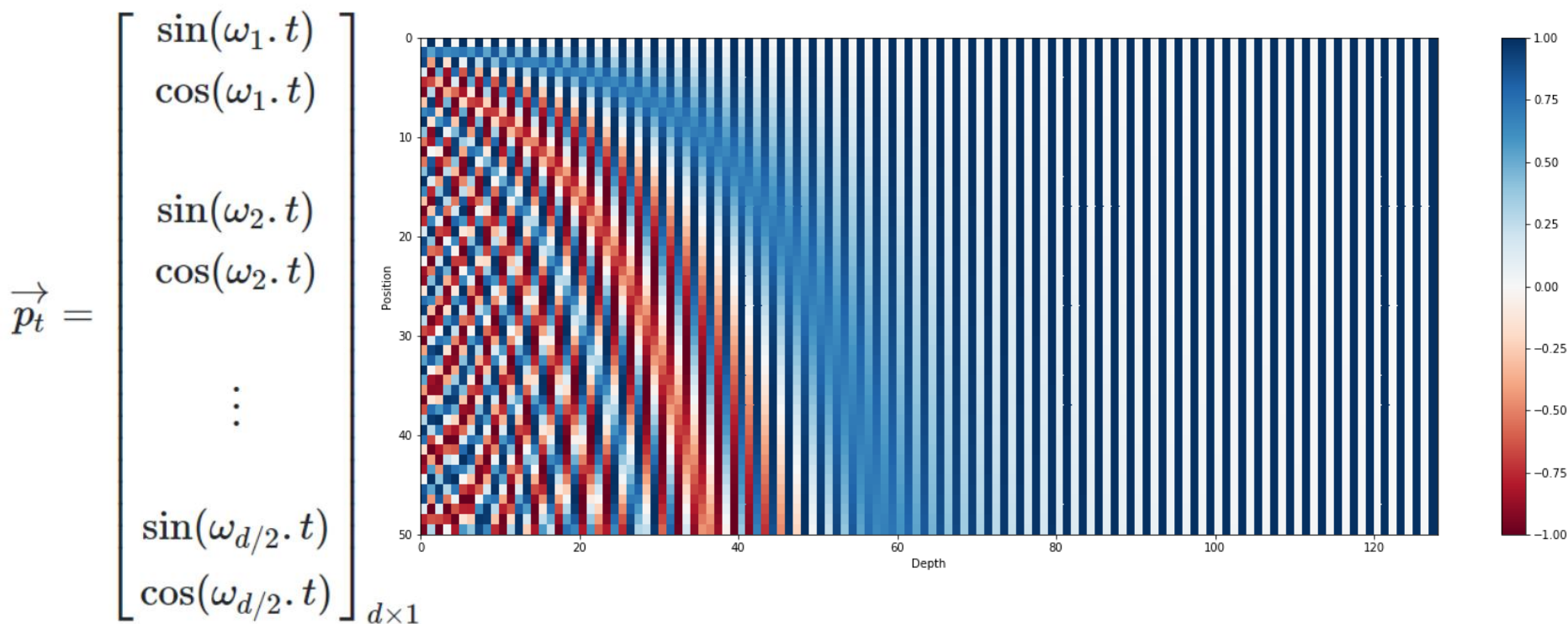
$$p^{(i,j)} := \begin{cases} \sin(\omega_k, i), \sin(\omega_k, j) & \text{if } i, j = 2k \\ \cos(\omega_k, i), \cos(\omega_k, j) & \text{if } i, j = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10,000^{2k/d}}, k \in \mathbb{N}_{\leq \lfloor \frac{d}{2} \rfloor}.$$

Positional Encoding

- Why? {
- Position and order of words define the grammar and actual semantics of a sentence.
 - Transformer ditched the recurrence mechanism, word order information is lost
- How? {
- assign a number within the [0, 1] range in which 0 means the first word and 1 is the last
 - assign a number to each time-step linearly
- can't figure out how many words are present within a specific range
- values get quite large; model may face sentences longer than the ones in training

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

output a unique encoding for each time-step (word's position in a sentence)



allows the model to attend relative positions effortlessly

$$M \cdot \begin{bmatrix} \sin(\omega_k \cdot t) \\ \cos(\omega_k \cdot t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k \cdot (t + \phi)) \\ \cos(\omega_k \cdot (t + \phi)) \end{bmatrix}$$

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k \cdot t) \\ \cos(\omega_k \cdot t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k \cdot (t + \phi)) \\ \cos(\omega_k \cdot (t + \phi)) \end{bmatrix}$$

$$M_{\phi,k} = \begin{bmatrix} \cos(\omega_k \cdot \phi) & \sin(\omega_k \cdot \phi) \\ -\sin(\omega_k \cdot \phi) & \cos(\omega_k \cdot \phi) \end{bmatrix}$$

fit GP model \longrightarrow get coefficients of PDE \longrightarrow homogeneous prediction

prior: $\mathbf{y} \sim \mathcal{N}(\bar{\mathbf{V}}, \mathbf{C})$

estimation: $\mathbf{y}^* \sim \mathcal{N}(\bar{\mathbf{V}}^*, \mathbf{C}^{**})$

joint distribution: $\begin{bmatrix} \mathbf{y} \\ \mathbf{y}^* \end{bmatrix} = \mathcal{N} \left(\begin{bmatrix} \bar{\mathbf{V}} \\ \bar{\mathbf{V}}^* \end{bmatrix}, \begin{bmatrix} \mathbf{C} + \sigma_y^2 \mathbf{I} & \mathbf{C}^* \\ \mathbf{C}^{*T} & \mathbf{C}^{**} \end{bmatrix} \right)$

conditional predictive distribution: $p(\mathbf{y}^* | \mathbf{V}^*, \phi, \sigma_y, \mathbf{V}, \mathbf{y}) = \mathcal{N}(\mathbf{y}^* | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{y}^* | \mathbf{V}^*, \phi, \sigma_y, \mathbf{V}, \mathbf{y}] = \bar{\mathbf{V}}^* + \mathbf{C}^{*T} (\mathbf{C} + \sigma_y^2 \mathbf{I})^{-1} (\mathbf{y} - \bar{\mathbf{V}})$$

$$\boldsymbol{\Sigma} = \text{cov}[\mathbf{y}^* | \mathbf{V}^*, \phi, \sigma_y, \mathbf{V}, \mathbf{y}] = \mathbf{C}^{**} - \mathbf{C}^{*T} (\mathbf{C} + \sigma_y^2 \mathbf{I})^{-1} \mathbf{C}^* + \sigma_y^2 \mathbf{I}$$

Estimating the PDE coefficients requires derivatives of the state variable (fitted by a GP model) with respect to independent variables.

for constant mean GP, the first order derivative:

$$p(\dot{\mathbf{y}}_{gp} | \mathbf{y}, \phi) = \mathcal{N}(\dot{\boldsymbol{\mu}}, \dot{\boldsymbol{\Sigma}}_{gp}). \quad \dot{\boldsymbol{\mu}} = \frac{\partial}{\partial v_j^*} \mathbb{E}[\mathbf{y}^* | V^*, \phi, \sigma_y, V, \mathbf{y}] = \dot{C}^{*T} (C + \sigma_y^2 I)^{-1} (\mathbf{y} - \bar{\mathbf{V}})$$

$$\dot{\boldsymbol{\Sigma}}_{gp} = \dot{C}^{**} - \dot{C}^{*T} (C + \sigma_y^2 I)^{-1} \dot{C}^*$$

$$\text{cov}(\mathbf{y}, \frac{\partial}{\partial v_j^*} \mathbf{y}^*) = \frac{\partial}{\partial v_j^*} k(\mathbf{v}, \mathbf{v}^*)$$

$$\text{cov}(\frac{\partial}{\partial v_j} \mathbf{y}, \frac{\partial}{\partial v_j^*} \mathbf{y}^*) = \frac{\partial^2}{\partial v_j \partial v_j^*} k(\mathbf{v}, \mathbf{v}^*)$$

the second-derivative:

$$\ddot{\mathbf{y}}_{gp} = \frac{\partial^2}{\partial v_j^{*2}} \mathbb{E}[\mathbf{y}^* | V^*, \phi, \sigma_y, V, \mathbf{y}] = \ddot{C}^{*T} (C + \sigma_y^2 I)^{-1} (\mathbf{y} - \bar{\mathbf{V}})$$

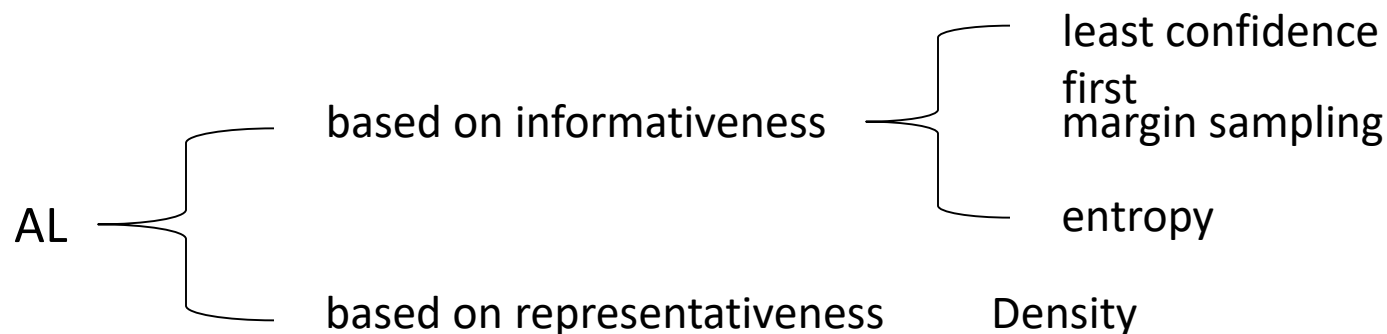
$$\text{cov}(y, \frac{\partial^2}{\partial v_j^{*2}} \mathbf{y}^*) = \frac{\partial^2}{\partial v_j^{*2}} k(\mathbf{v}, \mathbf{v}^*)$$

the residual error in the PDE at an observed point is given by:

$$\epsilon = f\left(v_1, \dots, v_m, y, \lambda, \mathcal{GP}\left\{\frac{\partial y}{\partial v_1}, \dots, \frac{\partial y}{\partial v_m}, \frac{\partial^2 y}{\partial v_1 \partial v_1}, \dots, \frac{\partial^2 y}{\partial v_1 \partial v_m}, \dots, \right\}\right)$$

minimizing the sum of square of residual error (SSRE)

$$SSRE = \epsilon^T \epsilon = \sum_{\forall v, y \in \mathbb{D}} f\left(v_1, \dots, v_m, y, \lambda, \mathcal{GP}\left\{\frac{\partial y}{\partial v_1}, \dots, \frac{\partial y}{\partial v_m}, \frac{\partial^2 y}{\partial v_1 \partial v_1}, \dots, \frac{\partial^2 y}{\partial v_1 \partial v_m}, \dots, \right\}\right)^2$$



$$\mathcal{Y}(x) \sim \mathcal{GP}(\mu(x), k(x, x'))$$

$$\hat{y}(x^*) = \mu(x^*) + c^T C^{-1}(\mathbf{y} - \mu)$$

mean squared error (MSE) of this prediction:

$$\hat{s}^2(x^*) = \mathbb{E}\{(\hat{y}(x^*) - \mathcal{Y}(x^*))^2\}$$

$$\hat{s}^2(x^*) = \sigma^2(x^*) - c^T C^{-1} c$$

Advantages:

Nonparametric

obtain model's uncertainty directly

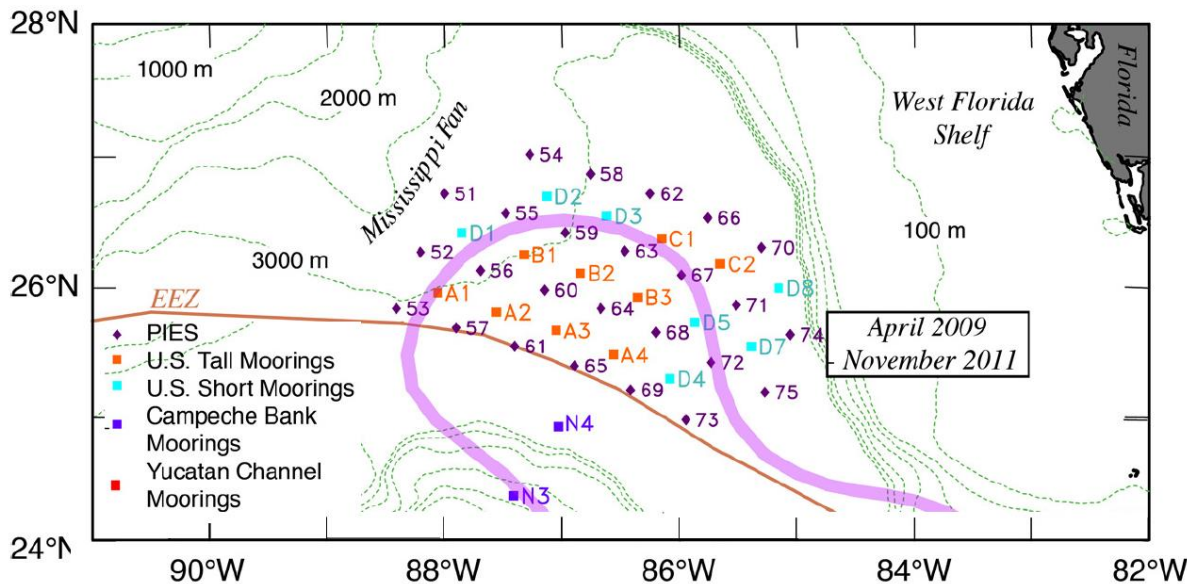
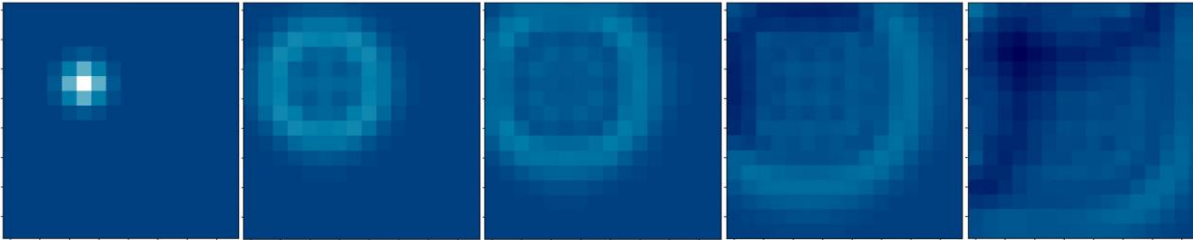
fit for nonlinear data

Disdvantages:

need to calculate inverse matrix of

all training data points

- Reflected Wave Simulation Data
- Gulf of Mexico (GoM) Loop Current Data



Ablation study:
Random sampling
Without physics learning
All data available

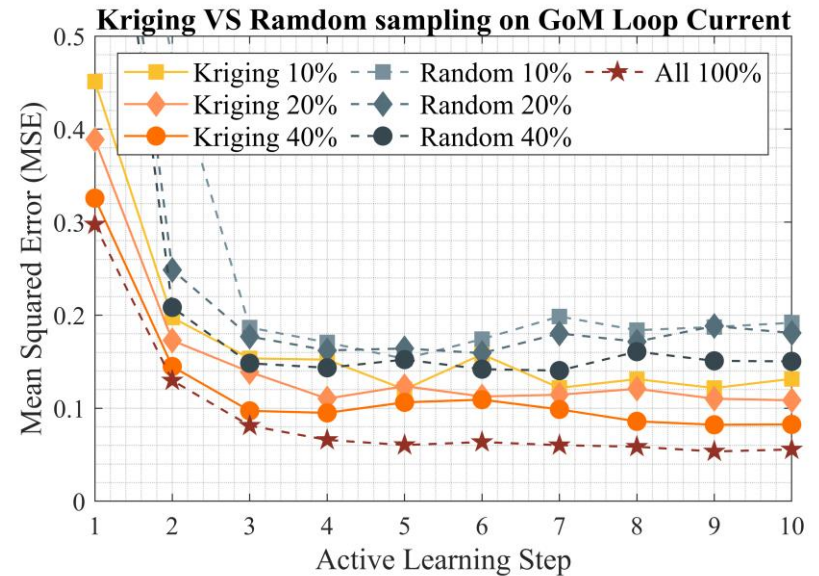
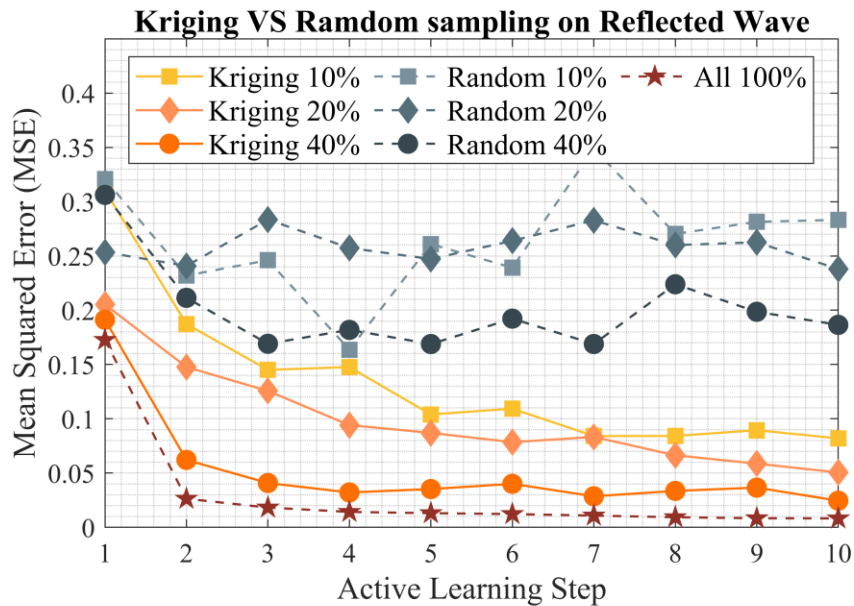


Figure 6: Explore Kriging and Random sampling effect on the performance of ST-PCNN on unseen data.

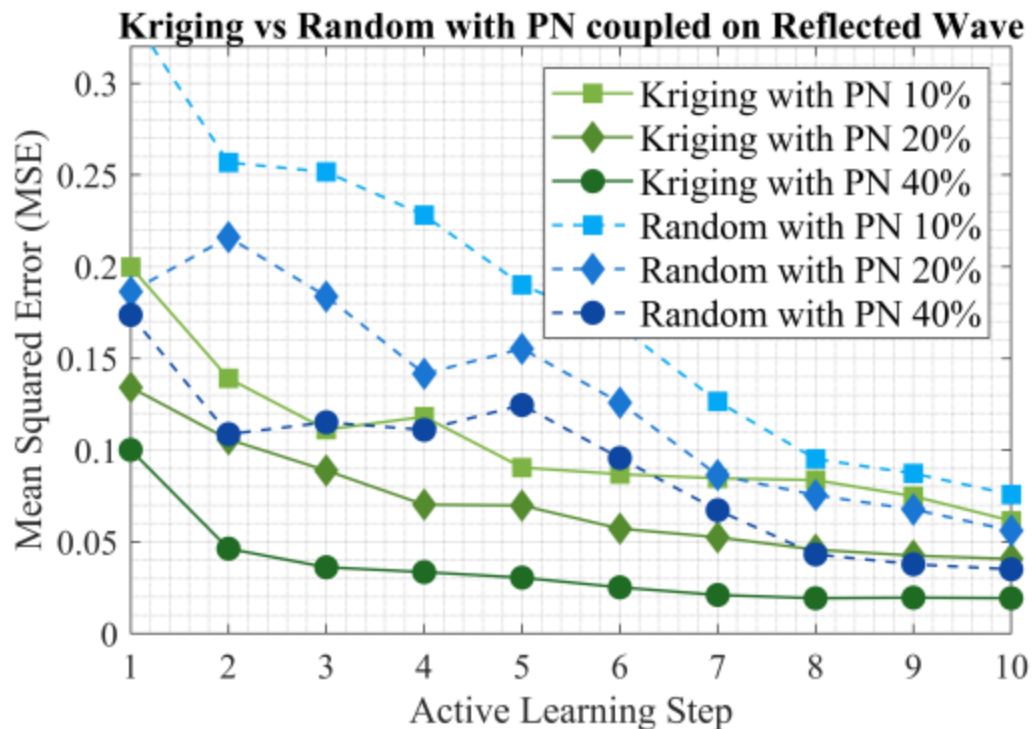


Figure 7: Explore Kriging and random sampling effect on performance of ST-PCNN on unseen data (the effect on GoM Loop Current refers to Figures 6 and 8).

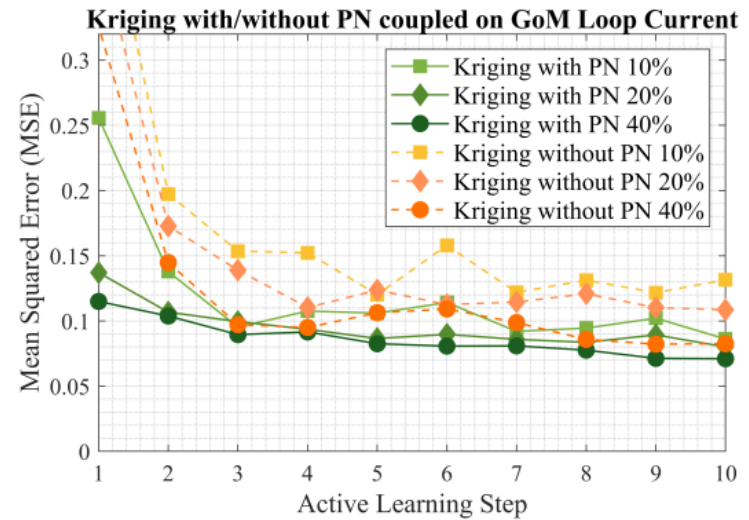
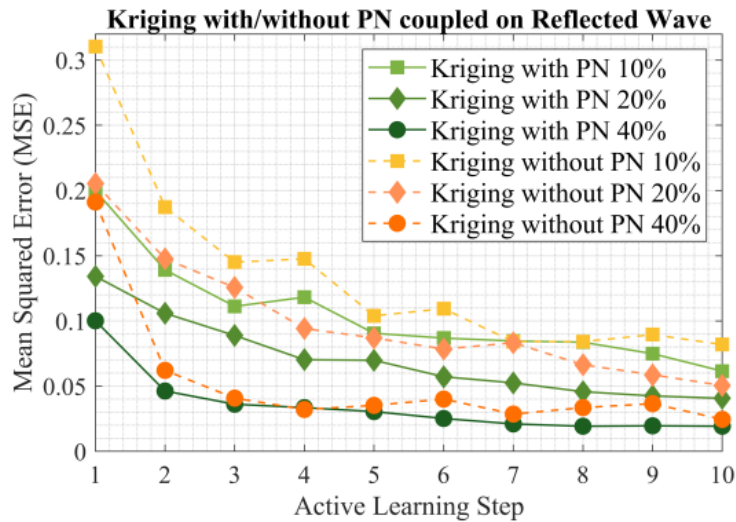


Figure 8: Explore physics effect on performance of ST-PCNN (with Kriging) on unseen data.