

# Statistical Methods for Spatio-Temporal Tensor Data

Hu Sun

Department of Statistics, University of Michigan

Doctoral Dissertation Defense

July 15, 2024

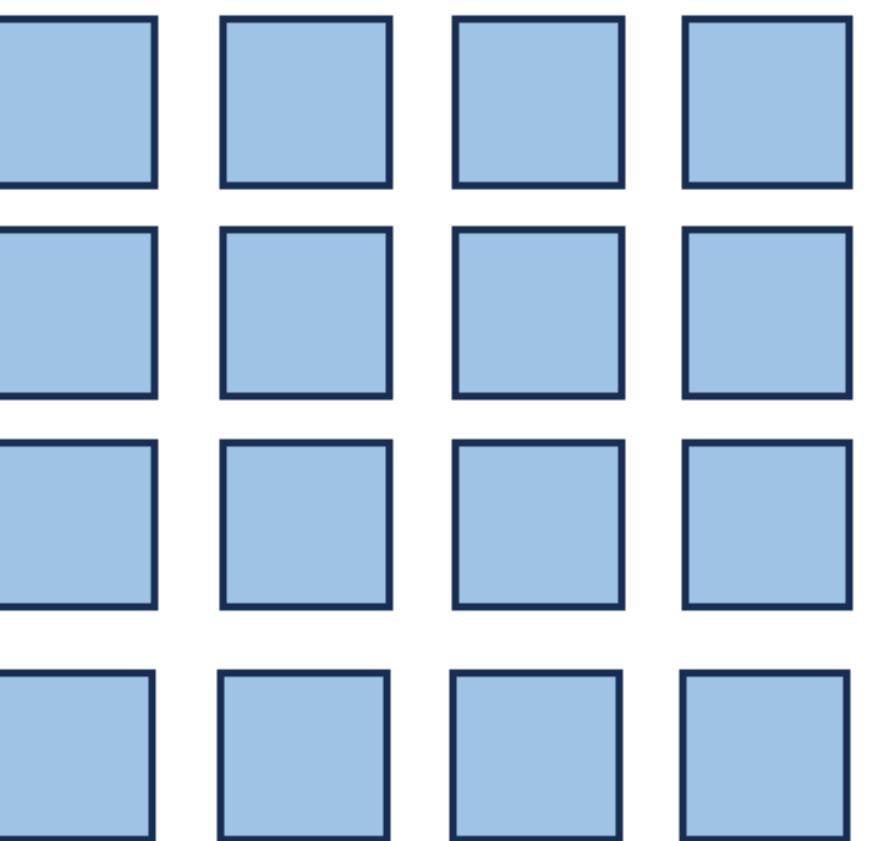
Advisor & Committee Chair: Prof. Yang Chen

Committee Members: Prof. Jian Kang, Prof. Kean Ming Tan, Prof. Ward B. Manchester

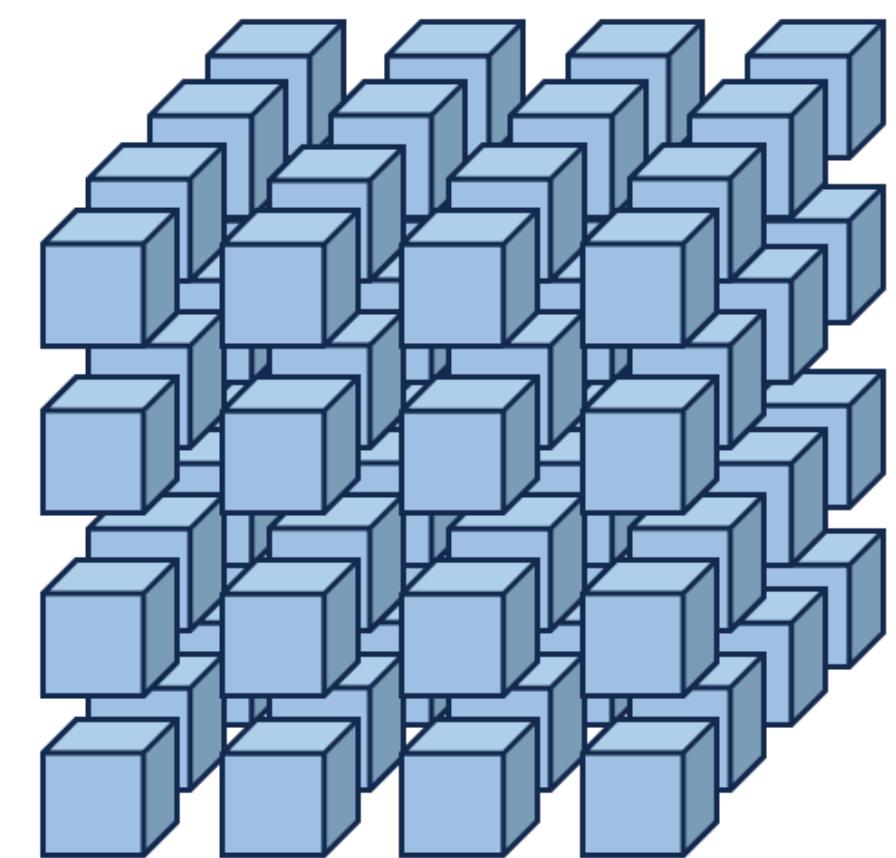
# Introduction to Tensor



order-1 tensor  $\mathbf{x}$   
(vector)



order-2 tensor  $\mathbf{X}$   
(matrix)



order-3 tensor  $\mathcal{X}$   
(tensor)

# Examples of Tensor Data

# Examples of Tensor Data



rating tensor  $\mathcal{X}$  in recommender system

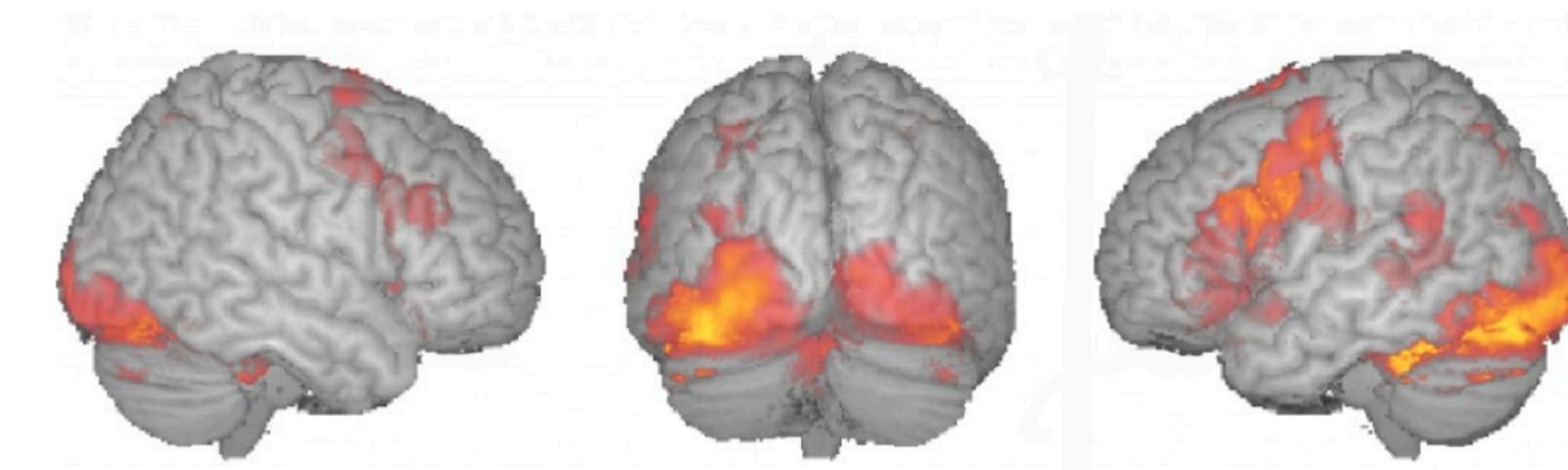
order-3 tensor:  $\mathcal{X}$  (user, item, context)

# Examples of Tensor Data



rating tensor  $\mathcal{X}$  in recommender system

order-3 tensor:  $\mathcal{X} (\text{user}, \text{item}, \text{context})$

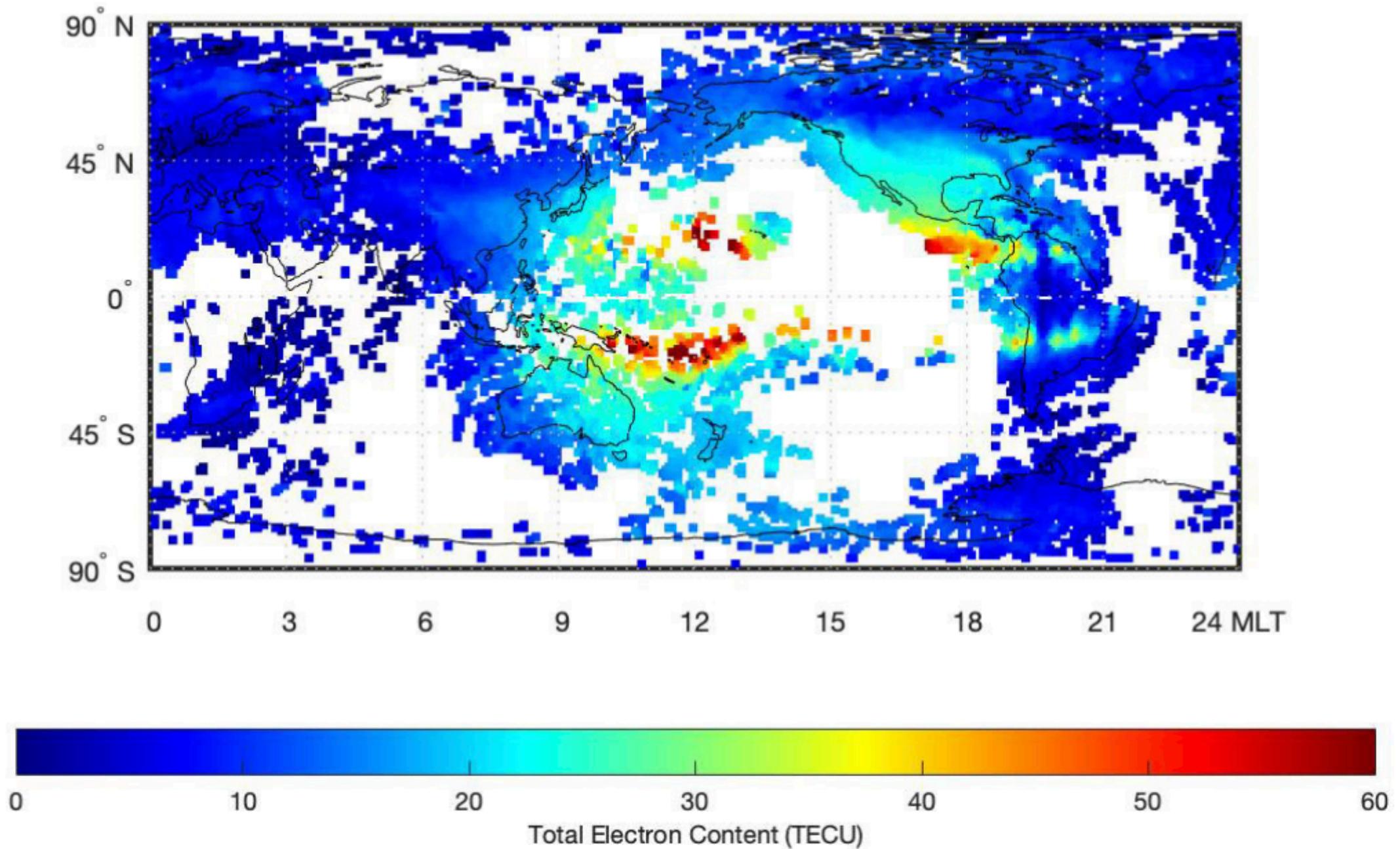


functional MRI 3-D brain imaging tensor  $\mathcal{Y}$

order-3 tensor:  $\mathcal{Y} (x, y, z)$

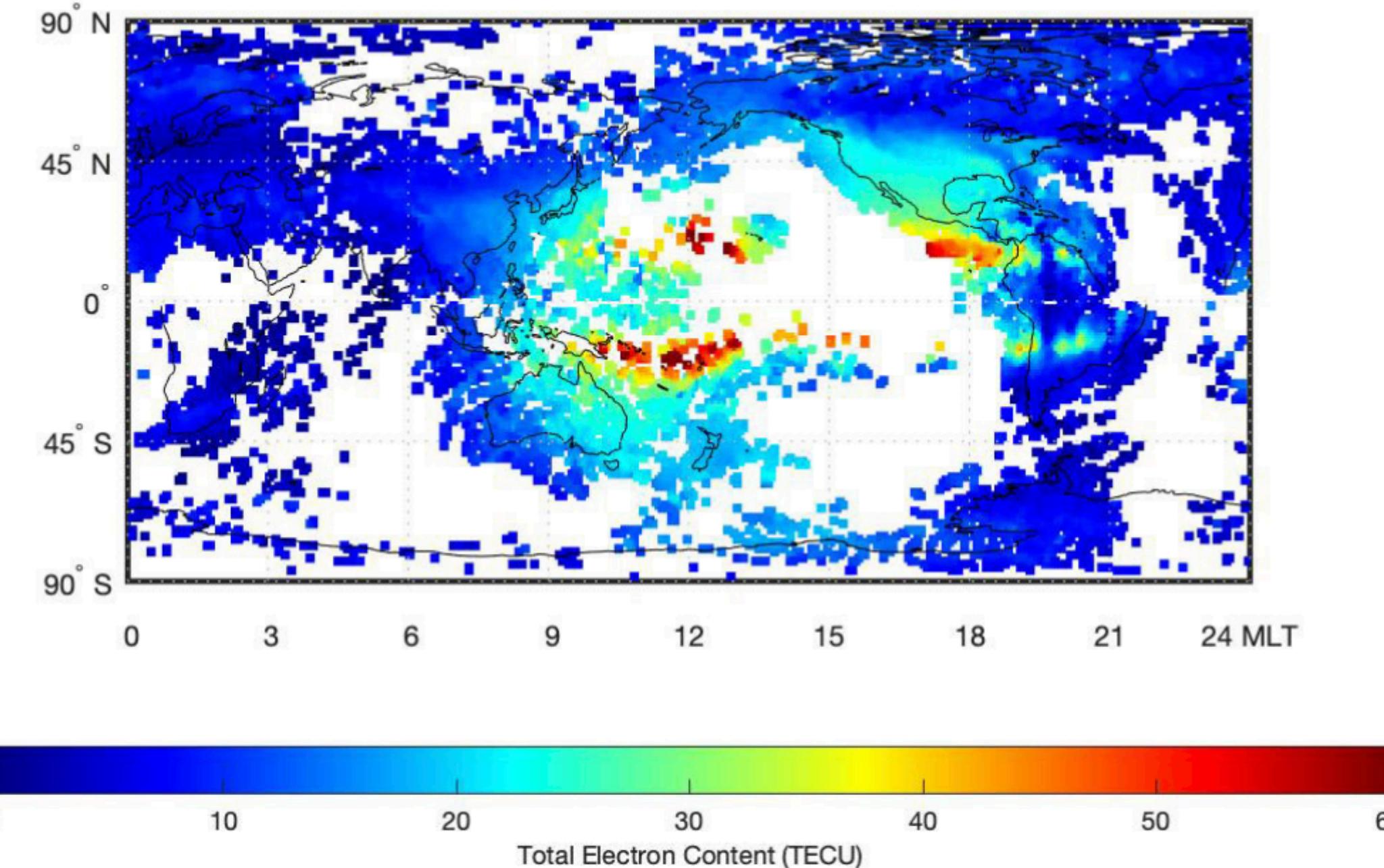
# Spatio-Temporal Tensor

Motivating Dataset: Global Total Electron Content (TEC)



# Spatio-Temporal Tensor

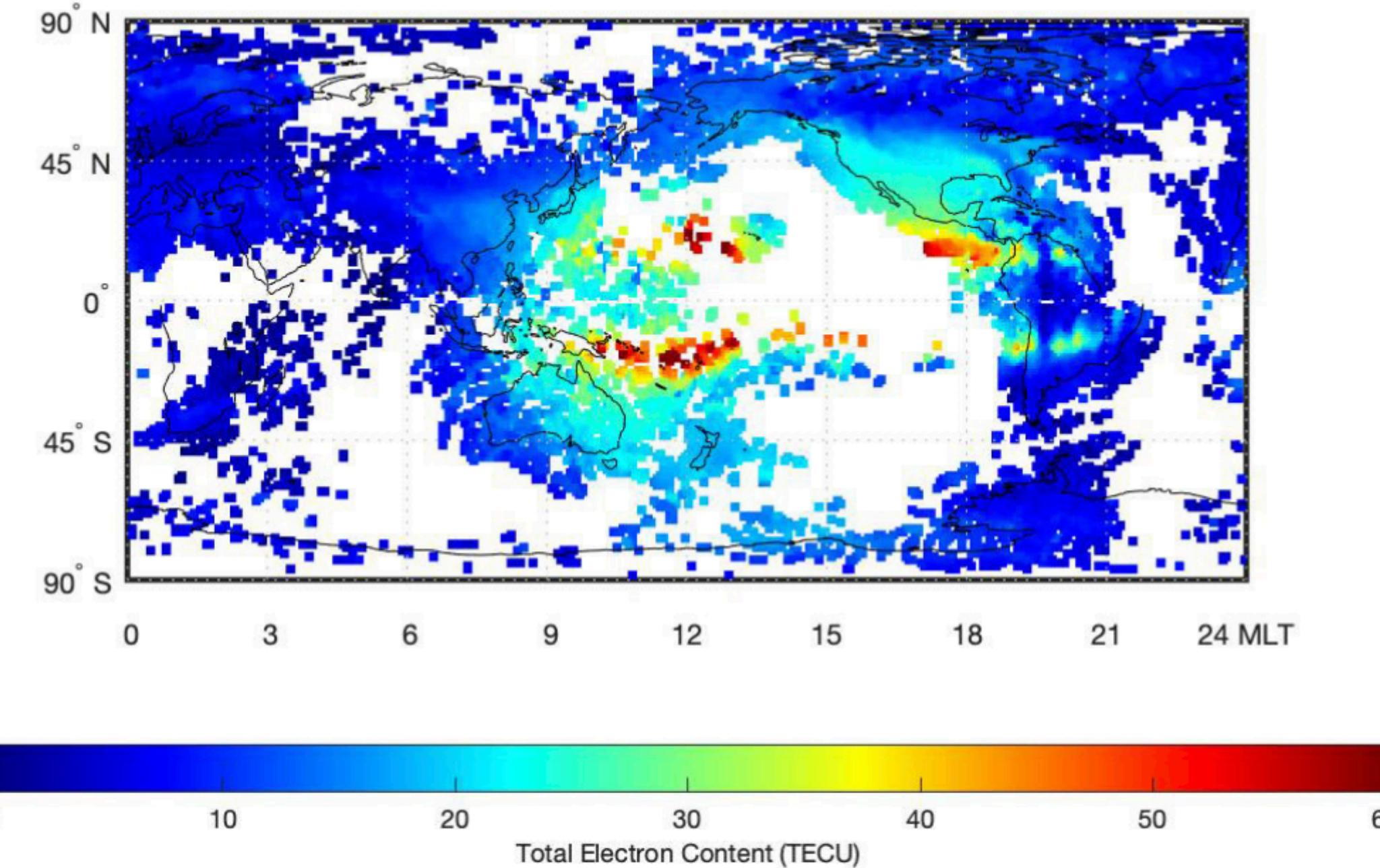
## Motivating Dataset: Global Total Electron Content (TEC)



- Monitoring the global TEC is critical since TEC affect satellite communication and satellite navigation (e.g. GPS positioning).

# Spatio-Temporal Tensor

Motivating Dataset: Global Total Electron Content (TEC)

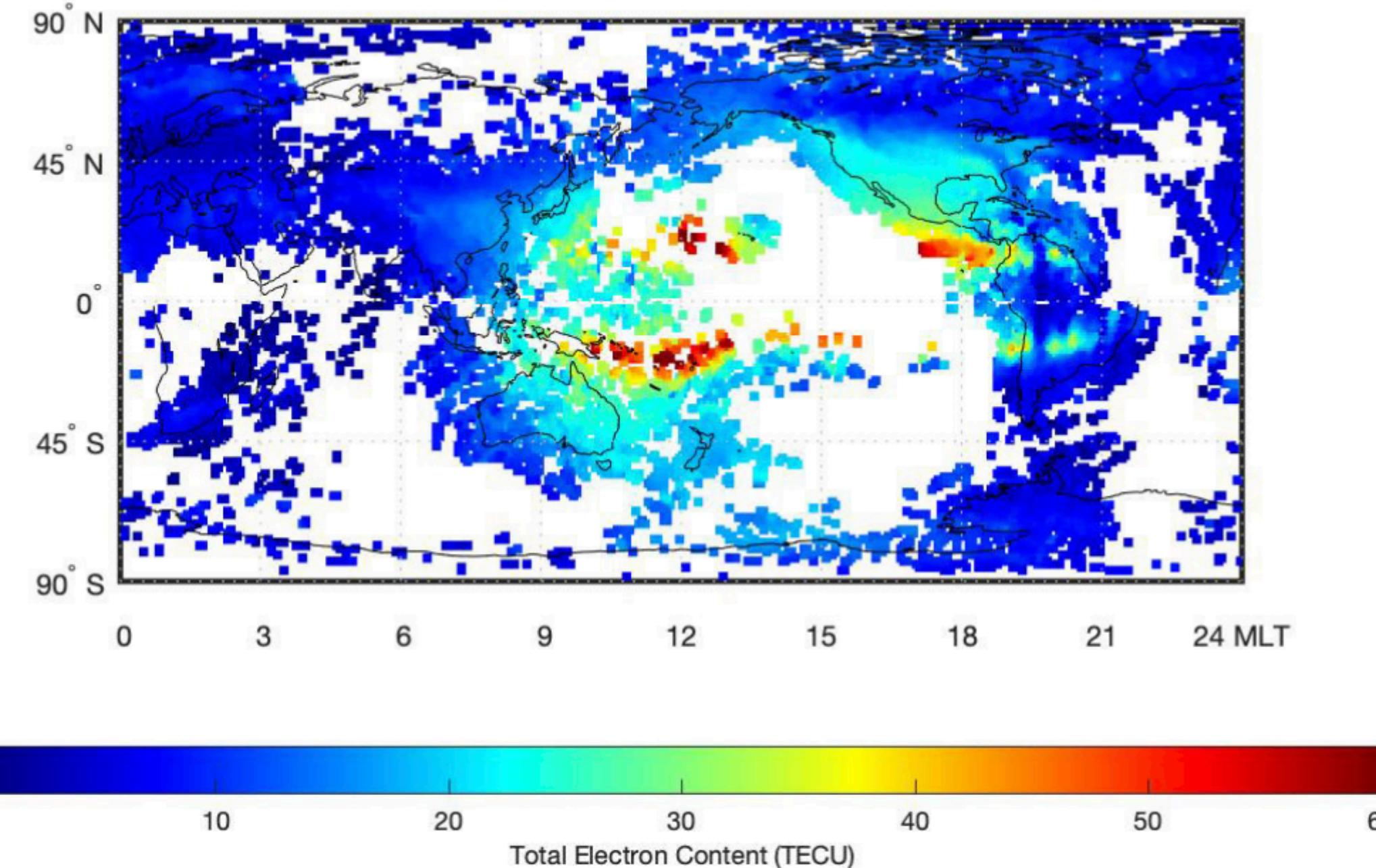


- Monitoring the global TEC is critical since TEC affect satellite communication and satellite navigation (e.g. GPS positioning).

**Challenges for such spatio-temporal tensors:**

# Spatio-Temporal Tensor

## Motivating Dataset: Global Total Electron Content (TEC)



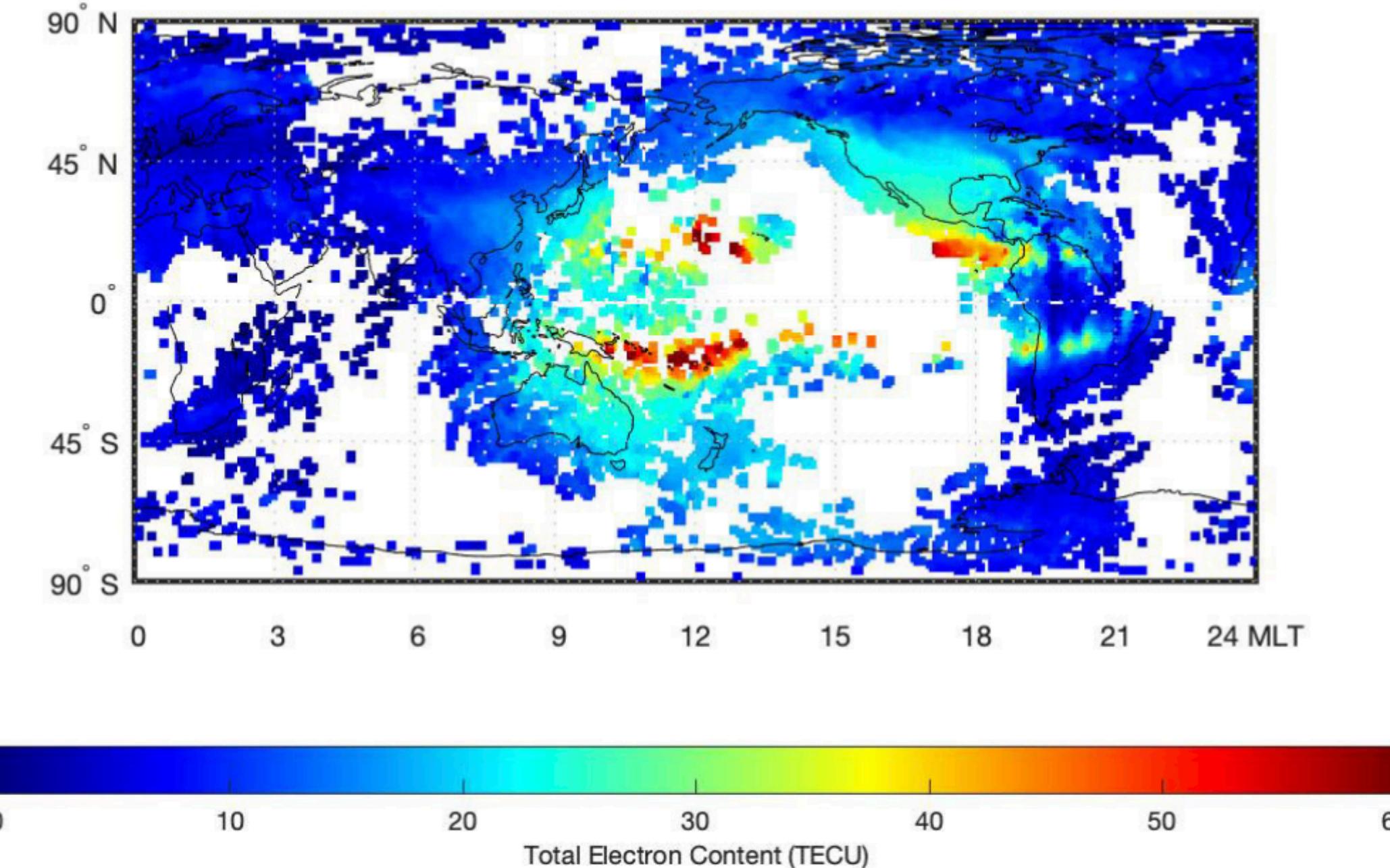
- Monitoring the global TEC is critical since TEC affect satellite communication and satellite navigation (e.g. GPS positioning).

### Challenges for such spatio-temporal tensors:

- multi-way spatio-temporal dependency

# Spatio-Temporal Tensor

## Motivating Dataset: Global Total Electron Content (TEC)



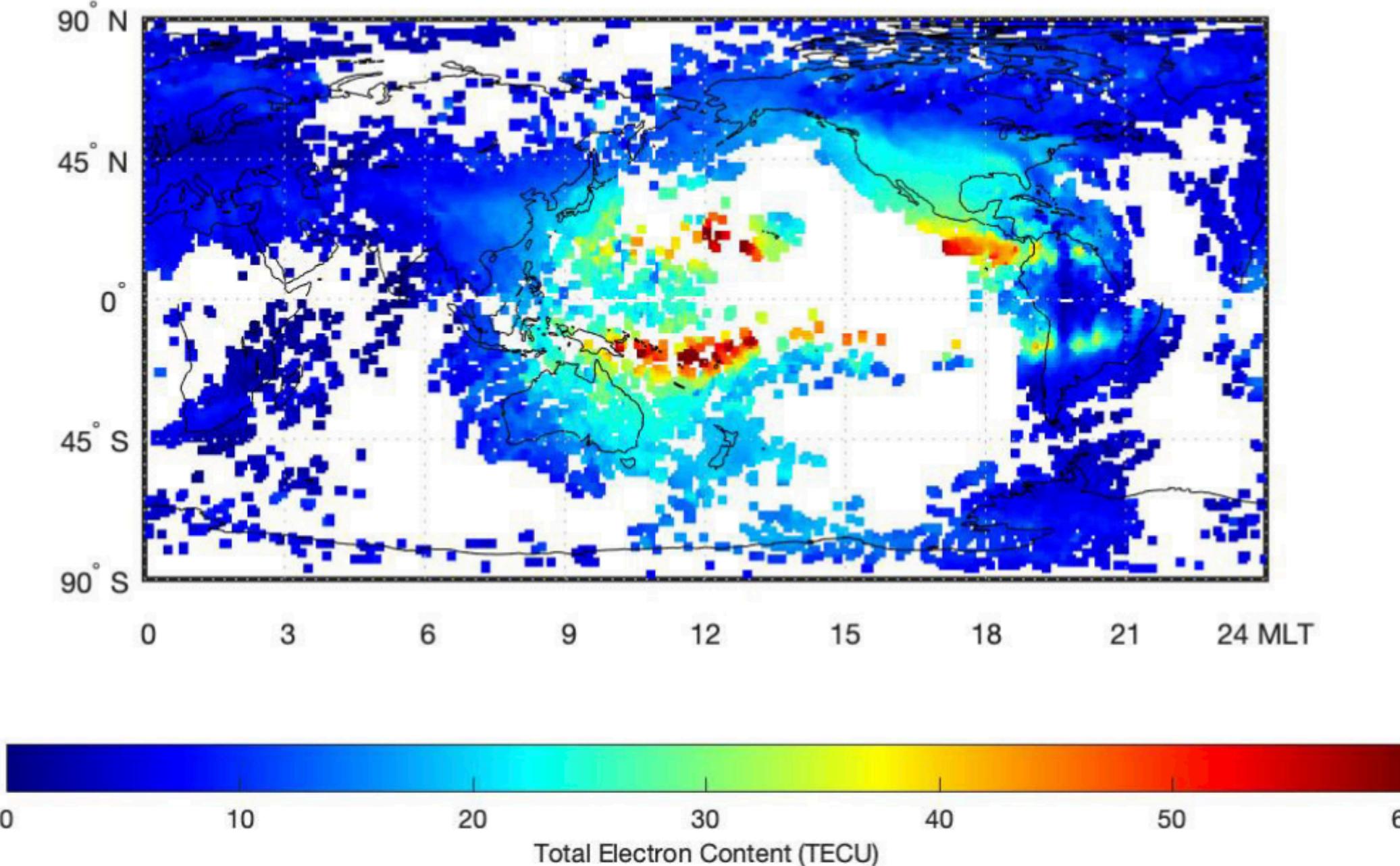
- Monitoring the global TEC is critical since TEC affect satellite communication and satellite navigation (e.g. GPS positioning).

### Challenges for such spatio-temporal tensors:

- multi-way spatio-temporal dependency
- dependent data missingness

# Spatio-Temporal Tensor

## Motivating Dataset: Global Total Electron Content (TEC)



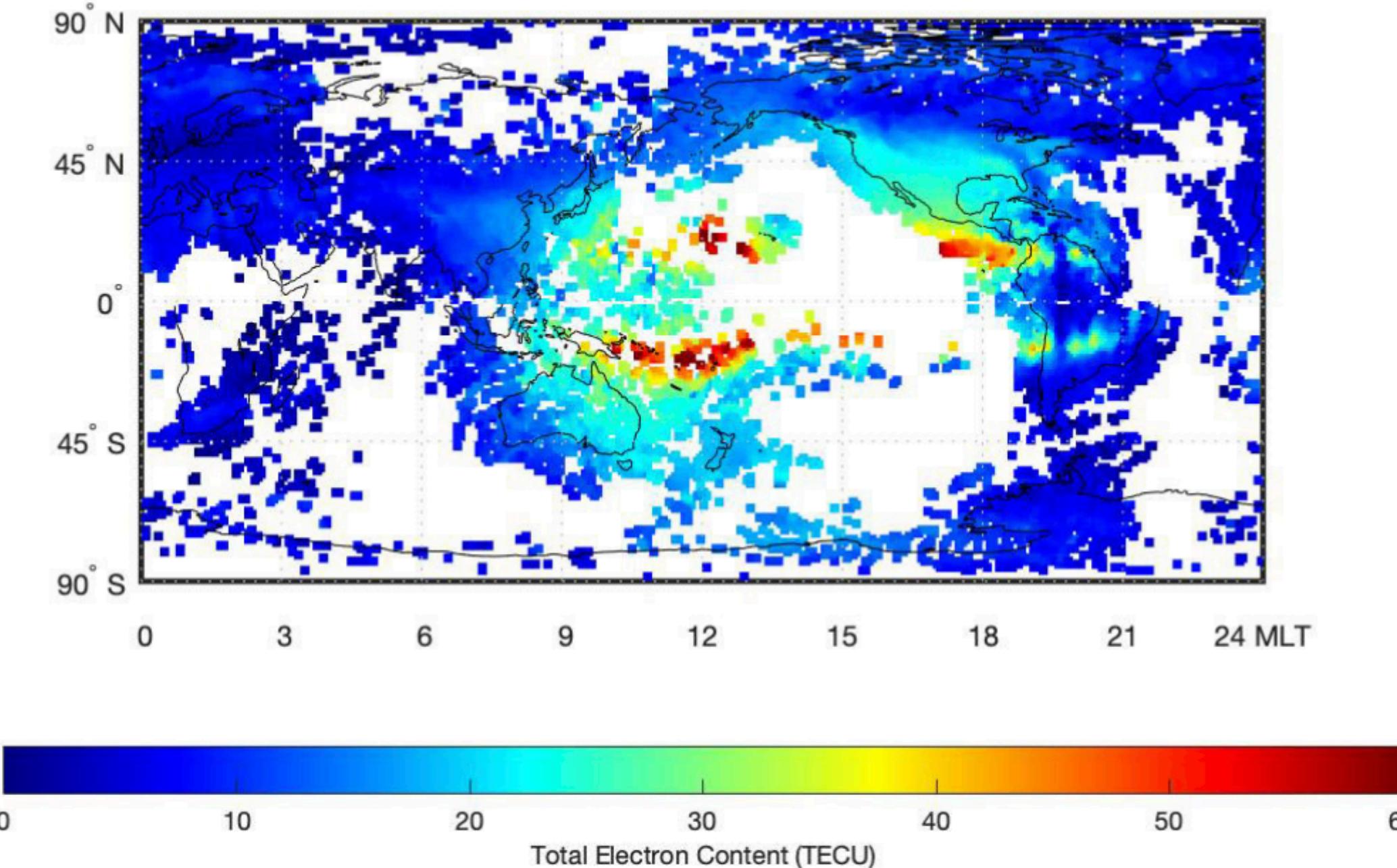
- Monitoring the global TEC is critical since TEC affect satellite communication and satellite navigation (e.g. GPS positioning).

### Challenges for such spatio-temporal tensors:

- multi-way spatio-temporal dependency
- dependent data missingness
- multiple data modalities

# Spatio-Temporal Tensor

## Motivating Dataset: Global Total Electron Content (TEC)



- Monitoring the global TEC is critical since TEC affect satellite communication and satellite navigation (e.g. GPS positioning).

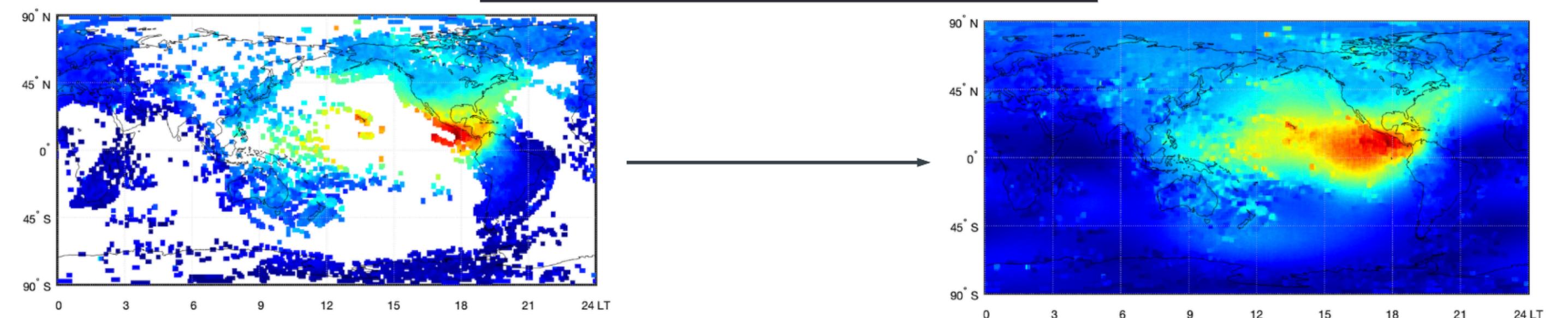
### Challenges for such spatio-temporal tensors:

- multi-way spatio-temporal dependency
- dependent data missingness
- multiple data modalities
- high data dimensionality

# Dissertation Summary

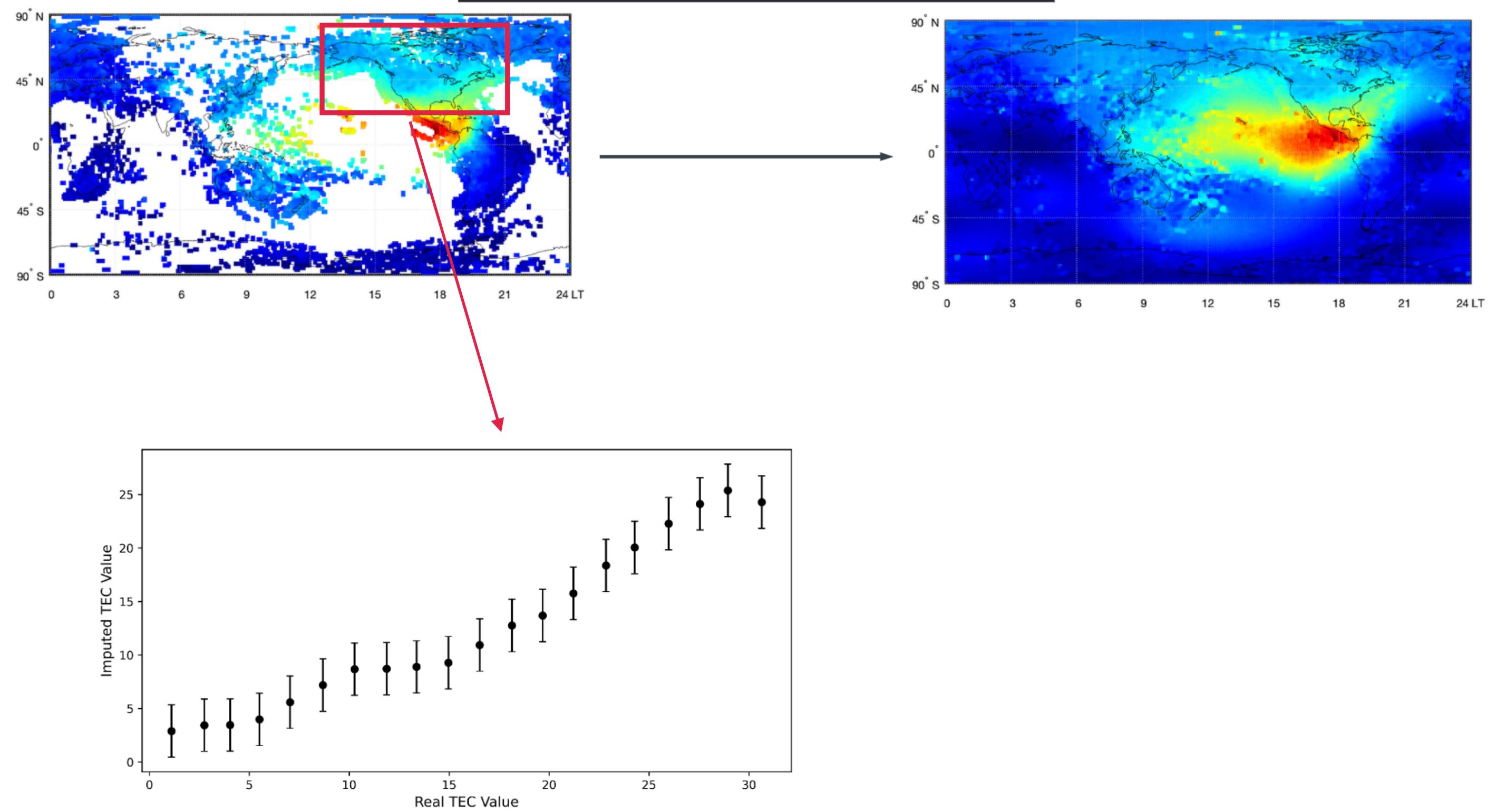
# Dissertation Summary

## I. Tensor Completion with Spatio-Temporal smoothing



# Dissertation Summary

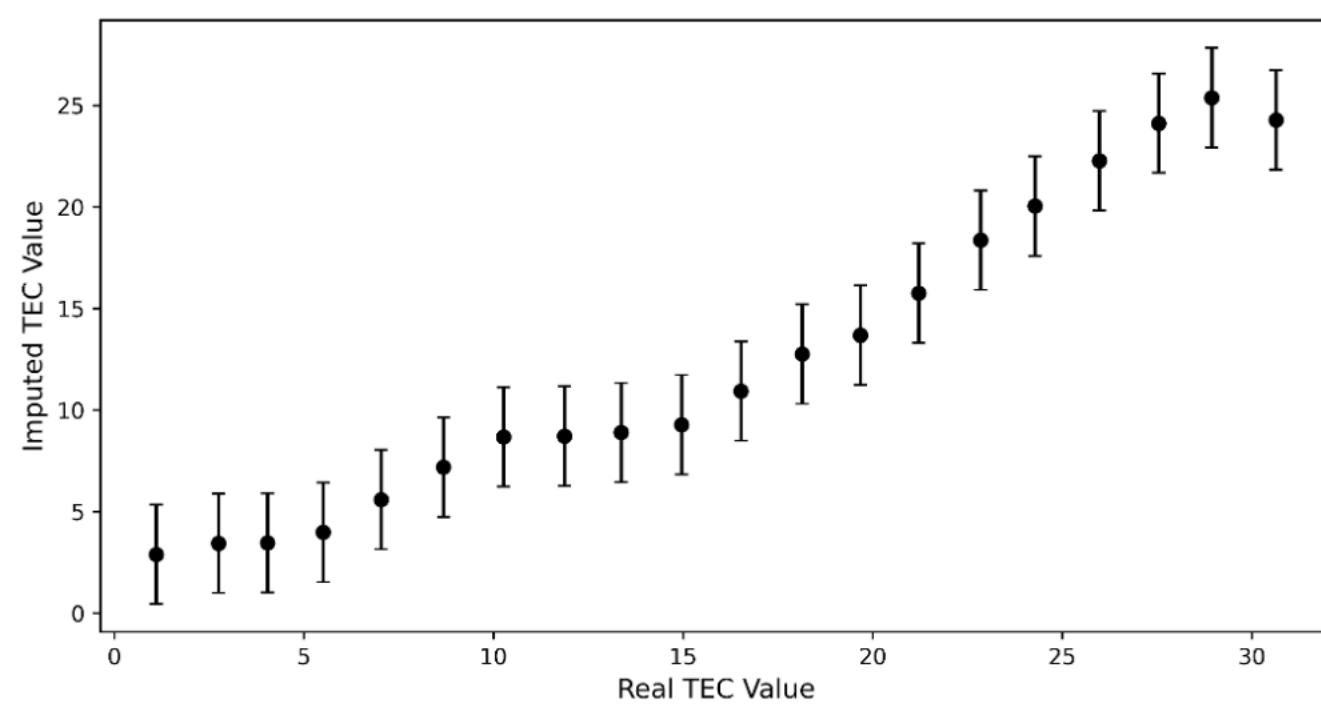
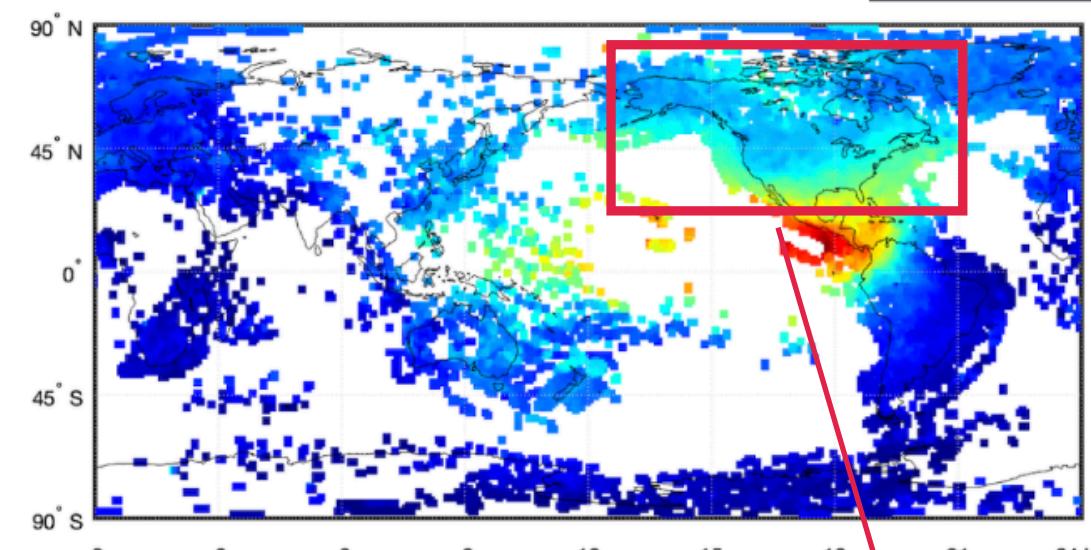
## I. Tensor Completion with Spatio-Temporal smoothing



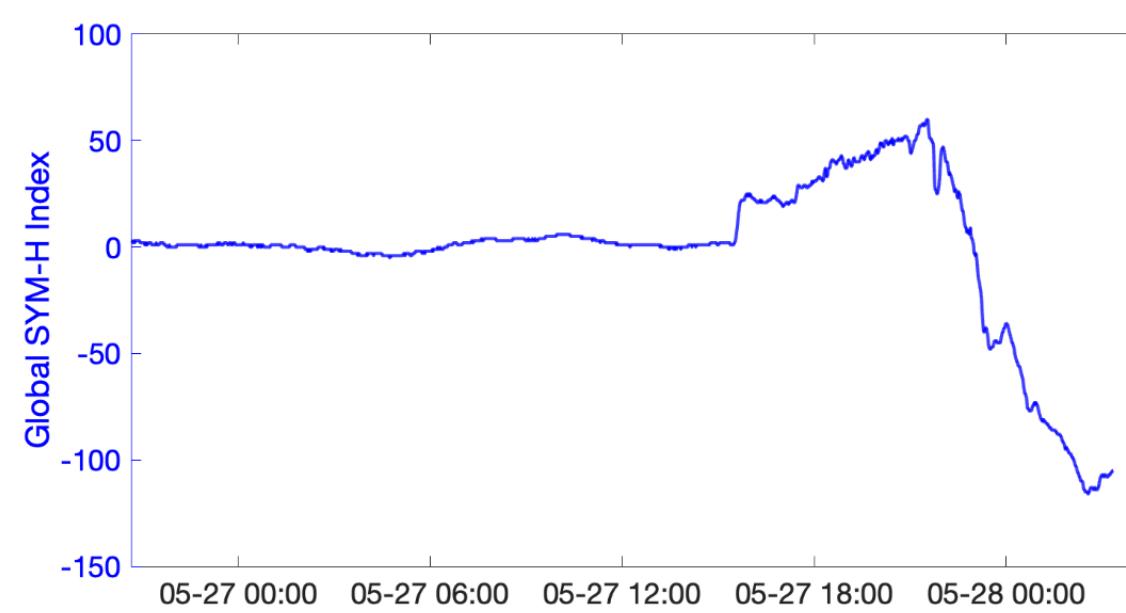
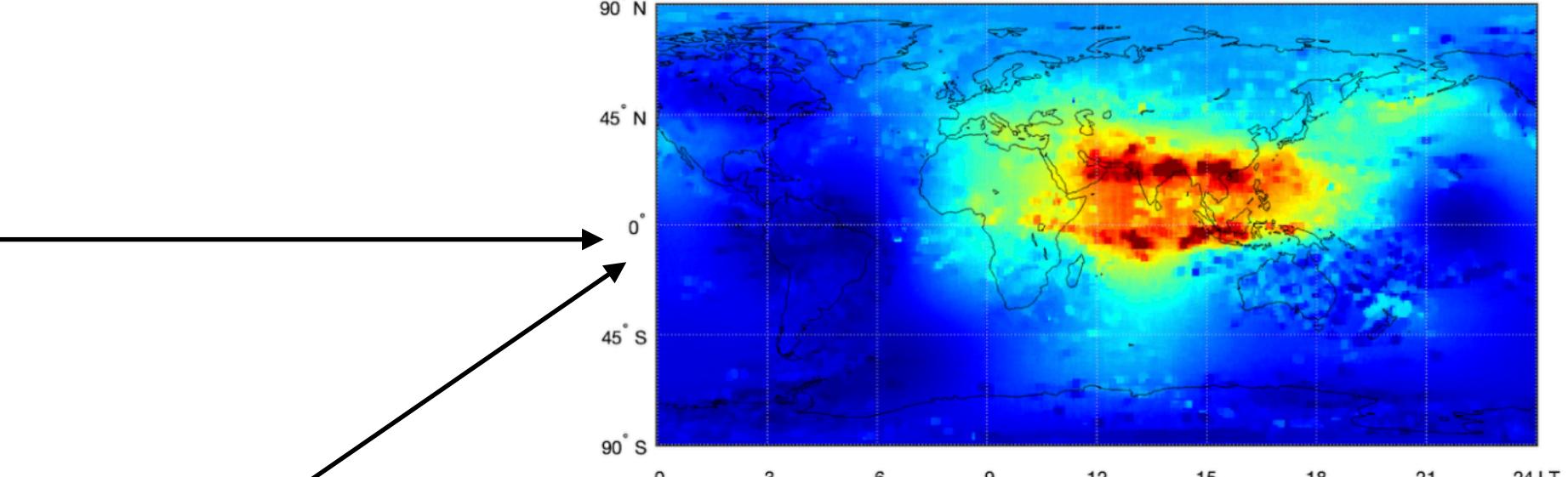
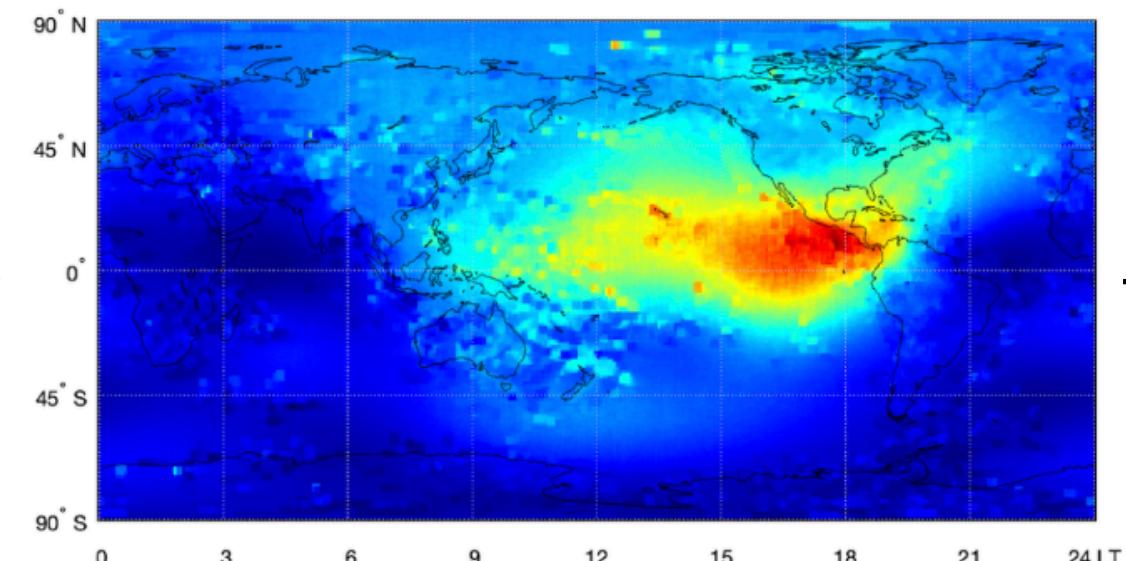
## II. Uncertainty Quantification for Tensor Completion

# Dissertation Summary

## I. Tensor Completion with Spatio-Temporal smoothing

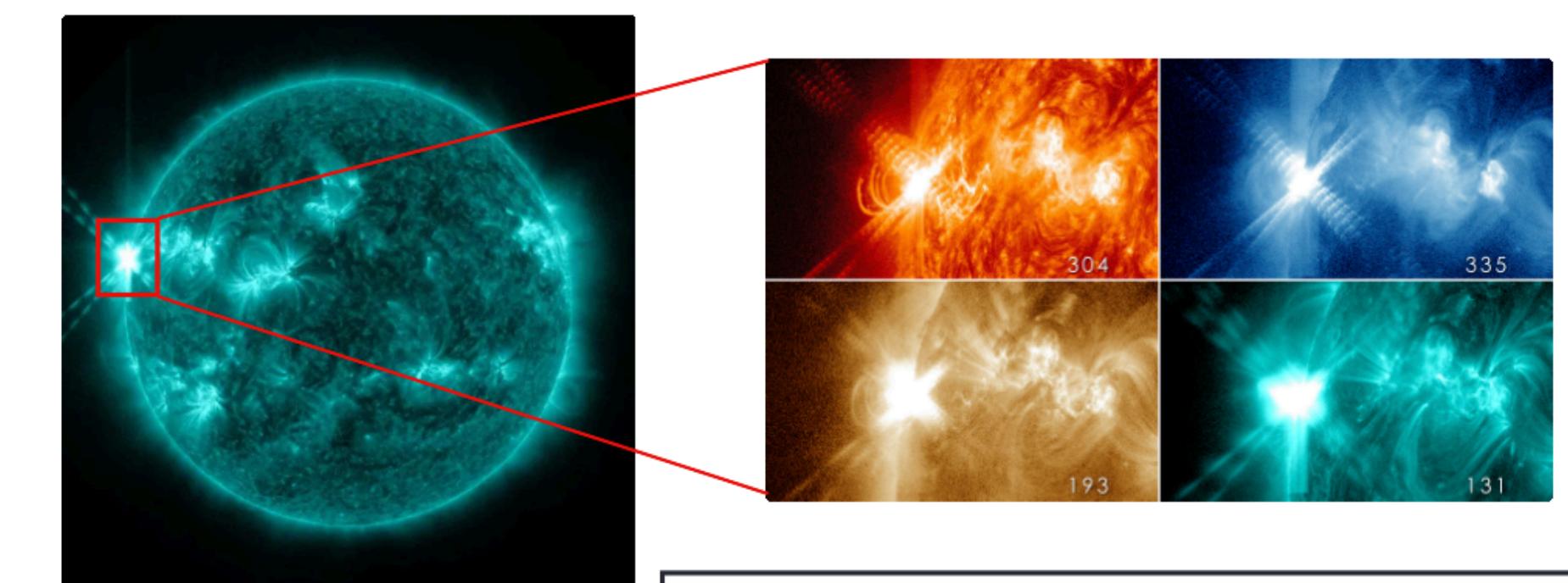
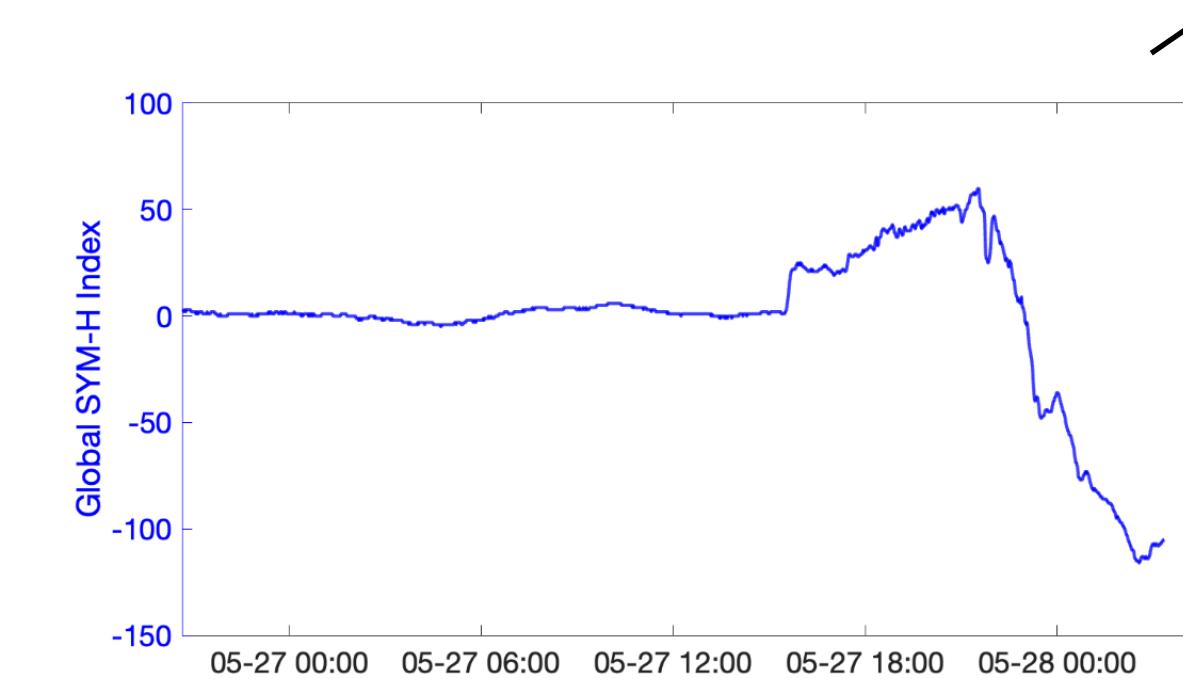
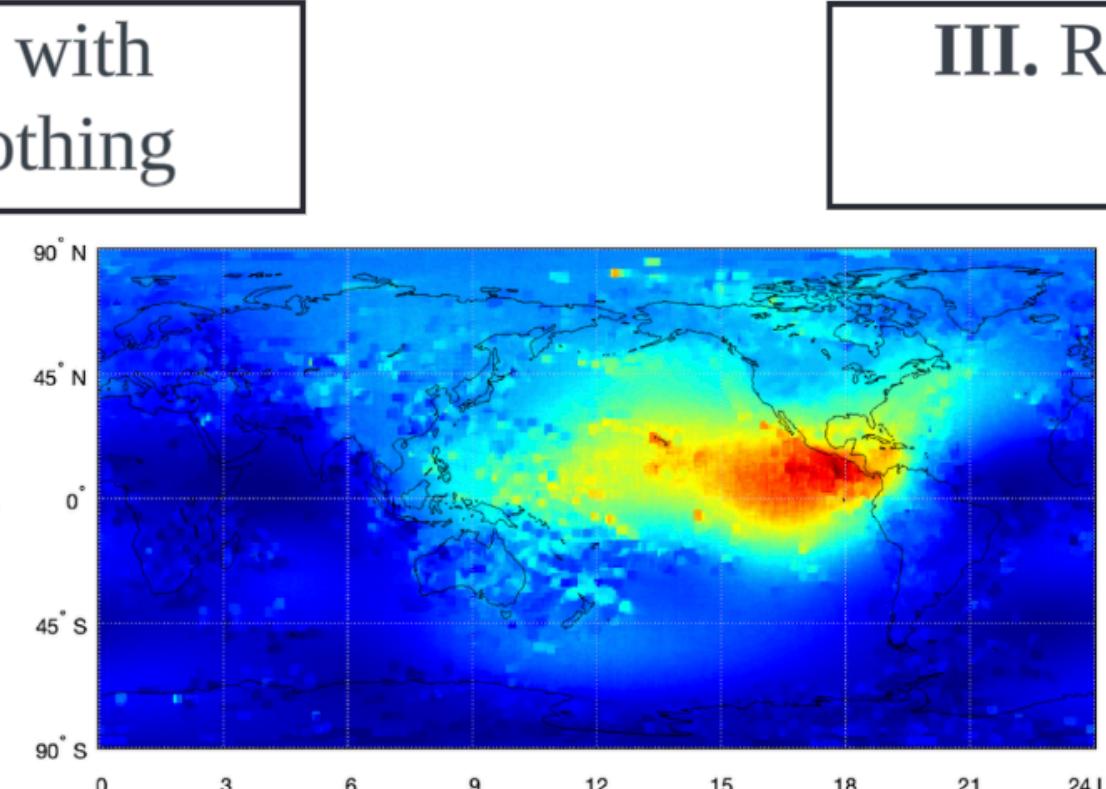
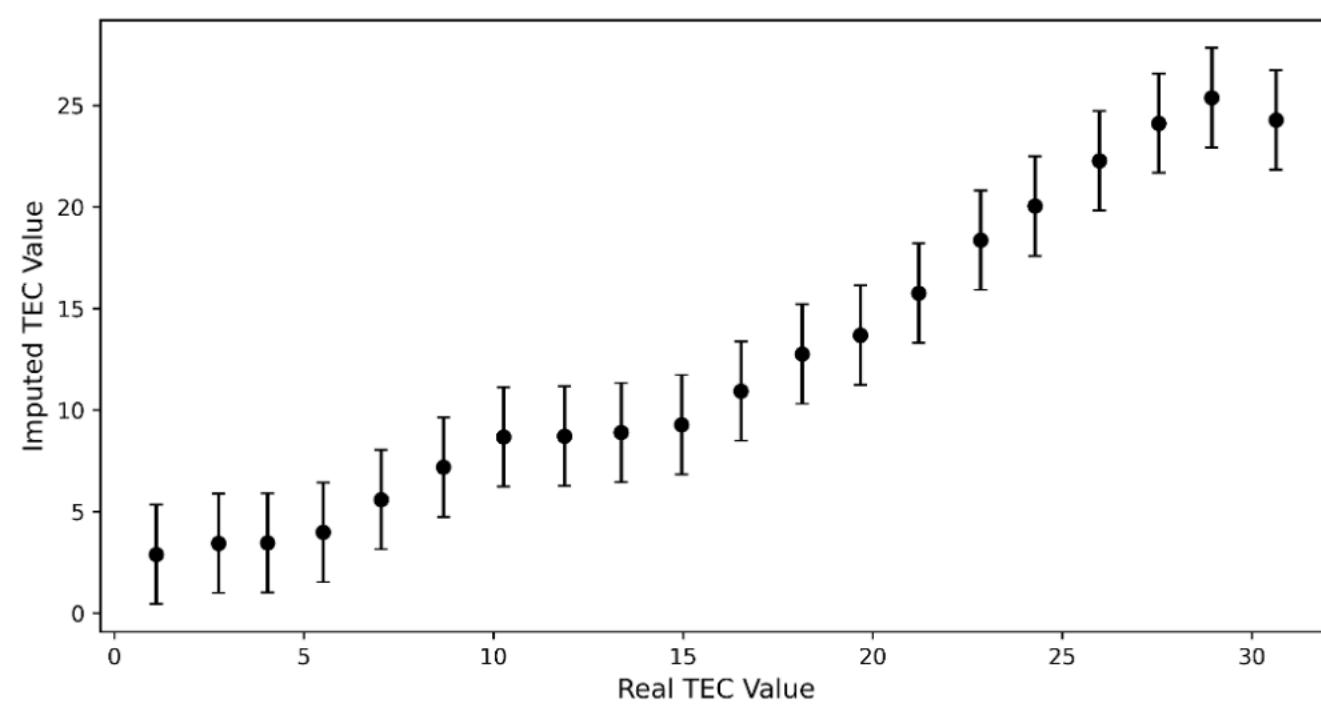
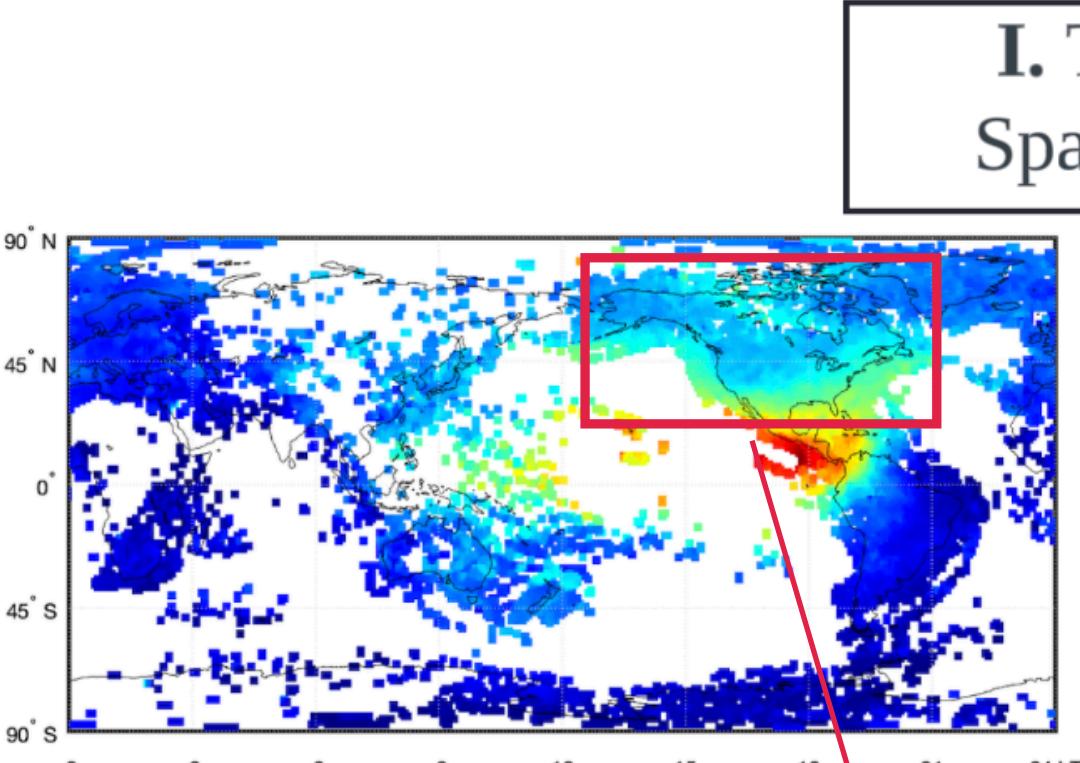


## III. RKHS Method for Multi-modal Autoregression Model



## II. Uncertainty Quantification for Tensor Completion

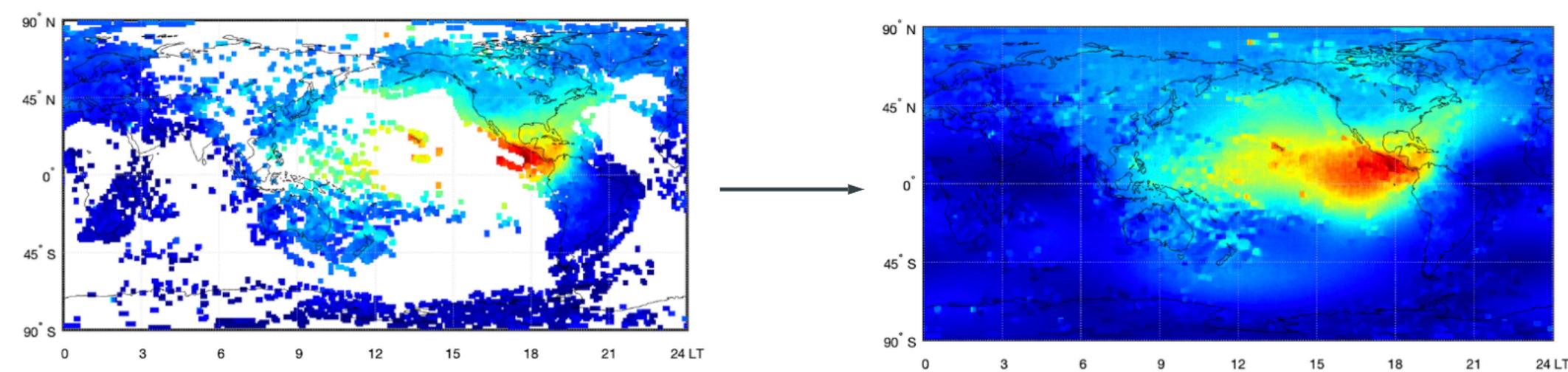
# Dissertation Summary



Solar Wind Parameters

**IV. Solar Flare Intensity Forecast with Scalar-on-Tensor Regression**

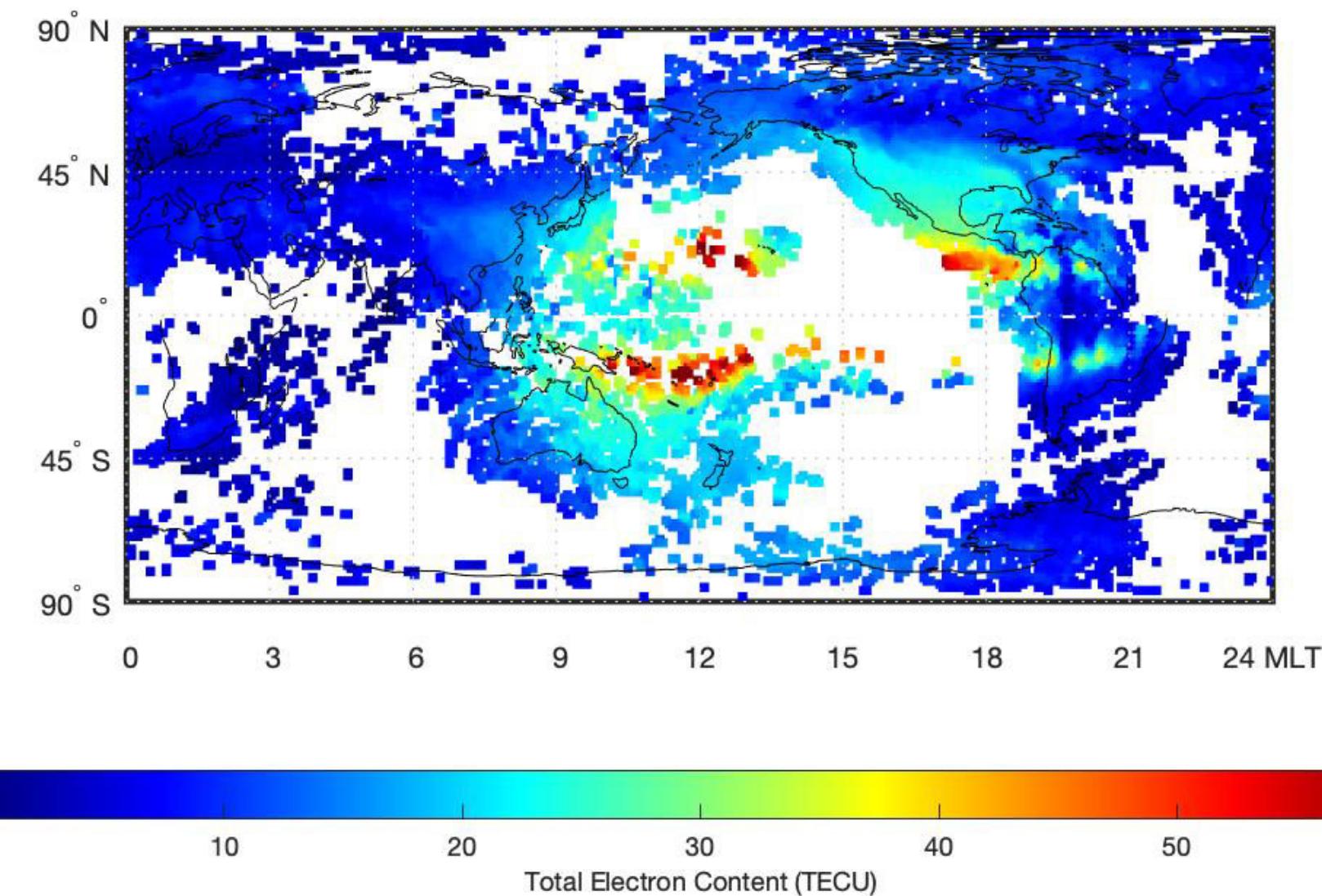
# I. Tensor Completion with Spatio-Temporal Smoothing



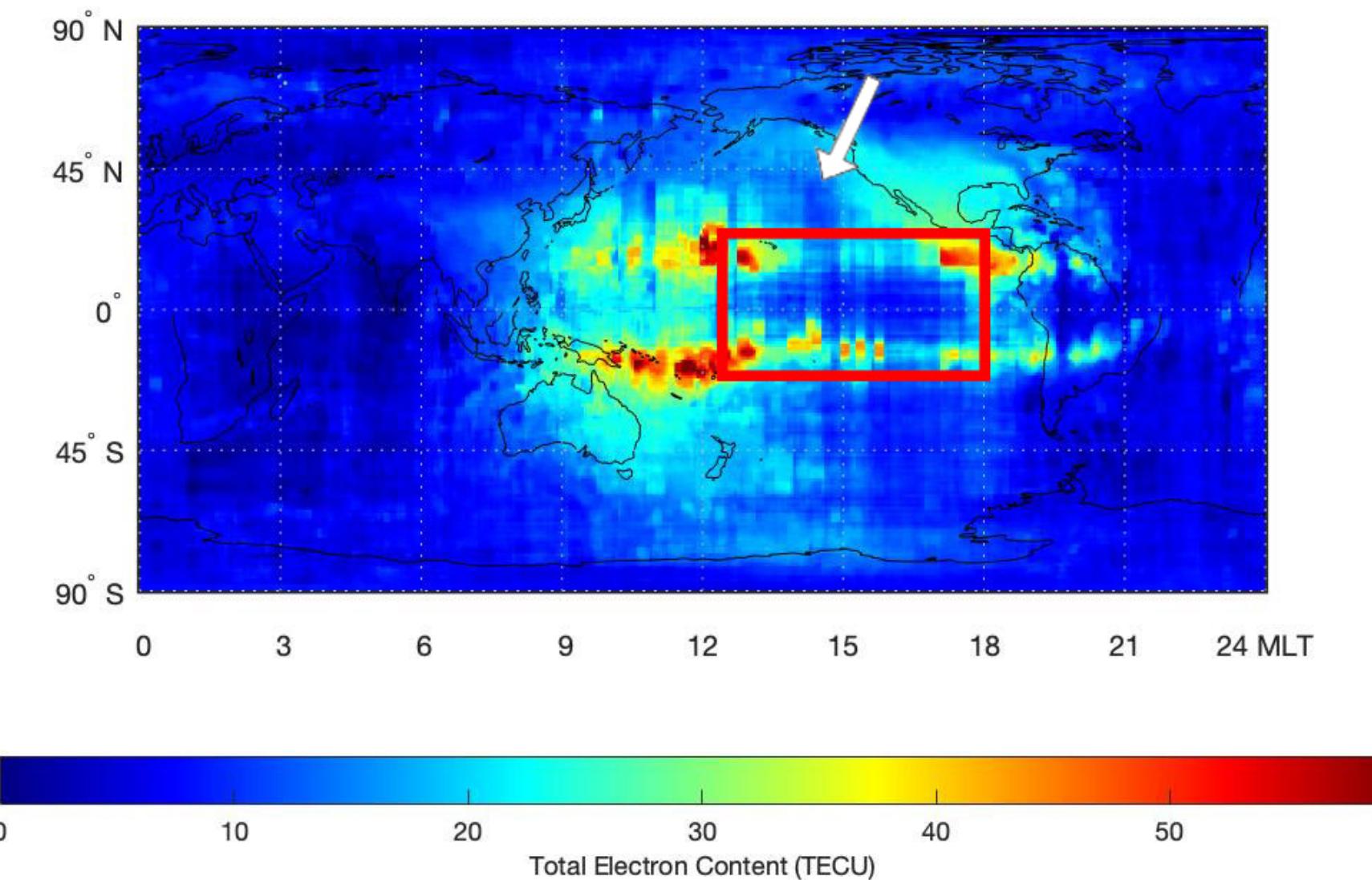
**How to generate accurate imputations when  
data are missing in spatio-temporal patches?**

# Low-rank Completion

(A) Original Map

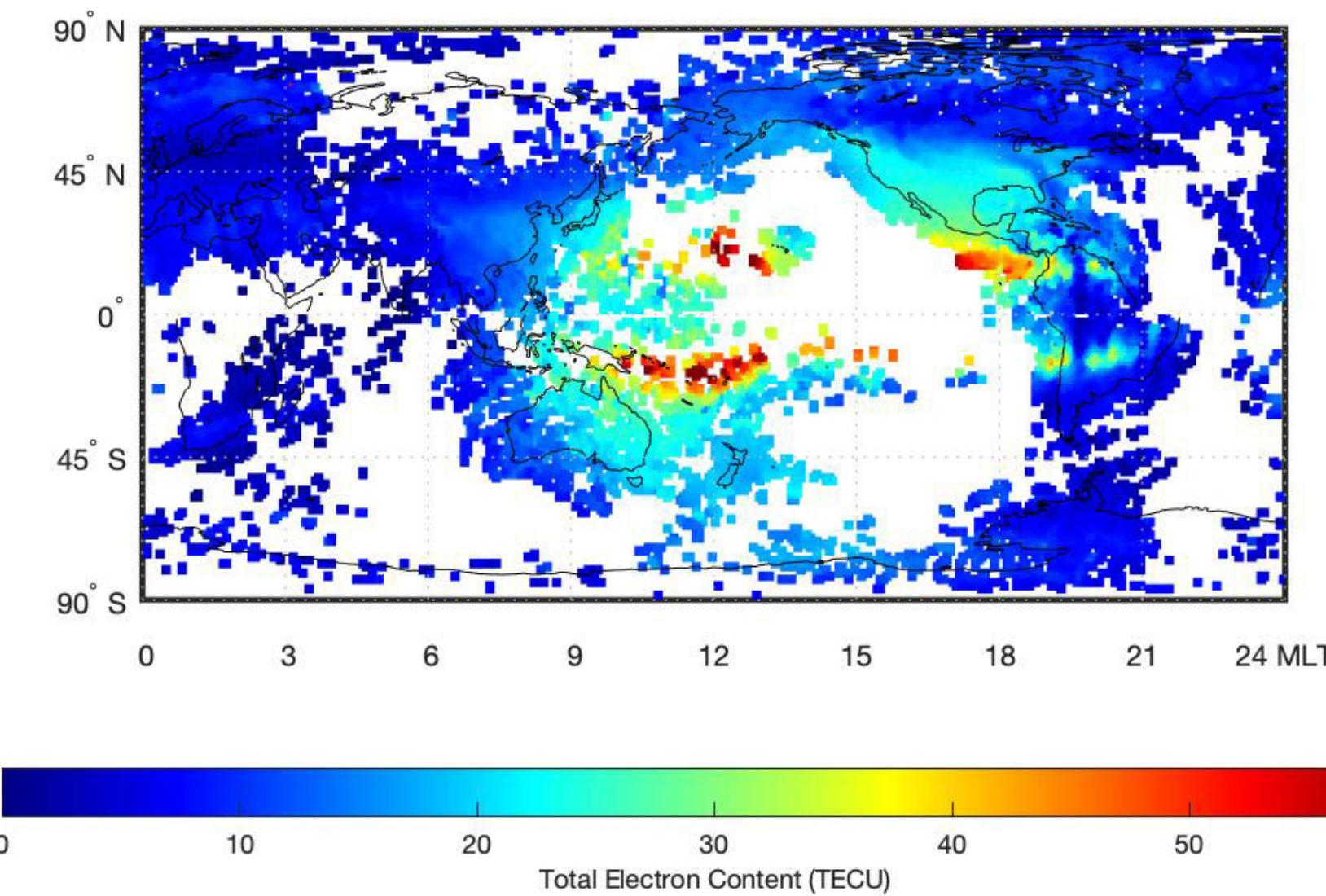


(B) SoftImpute

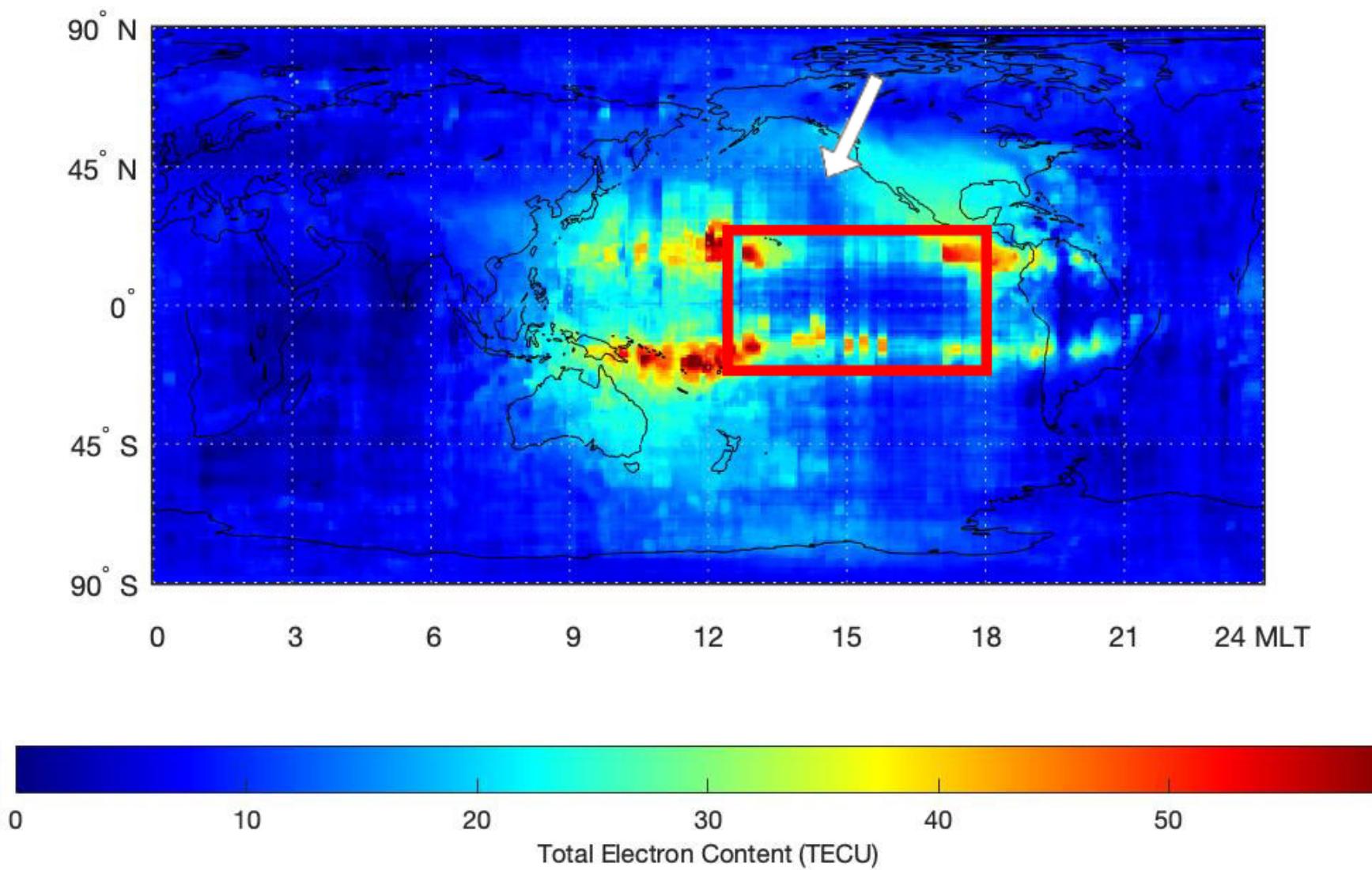


# Low-rank Completion

(A) Original Map



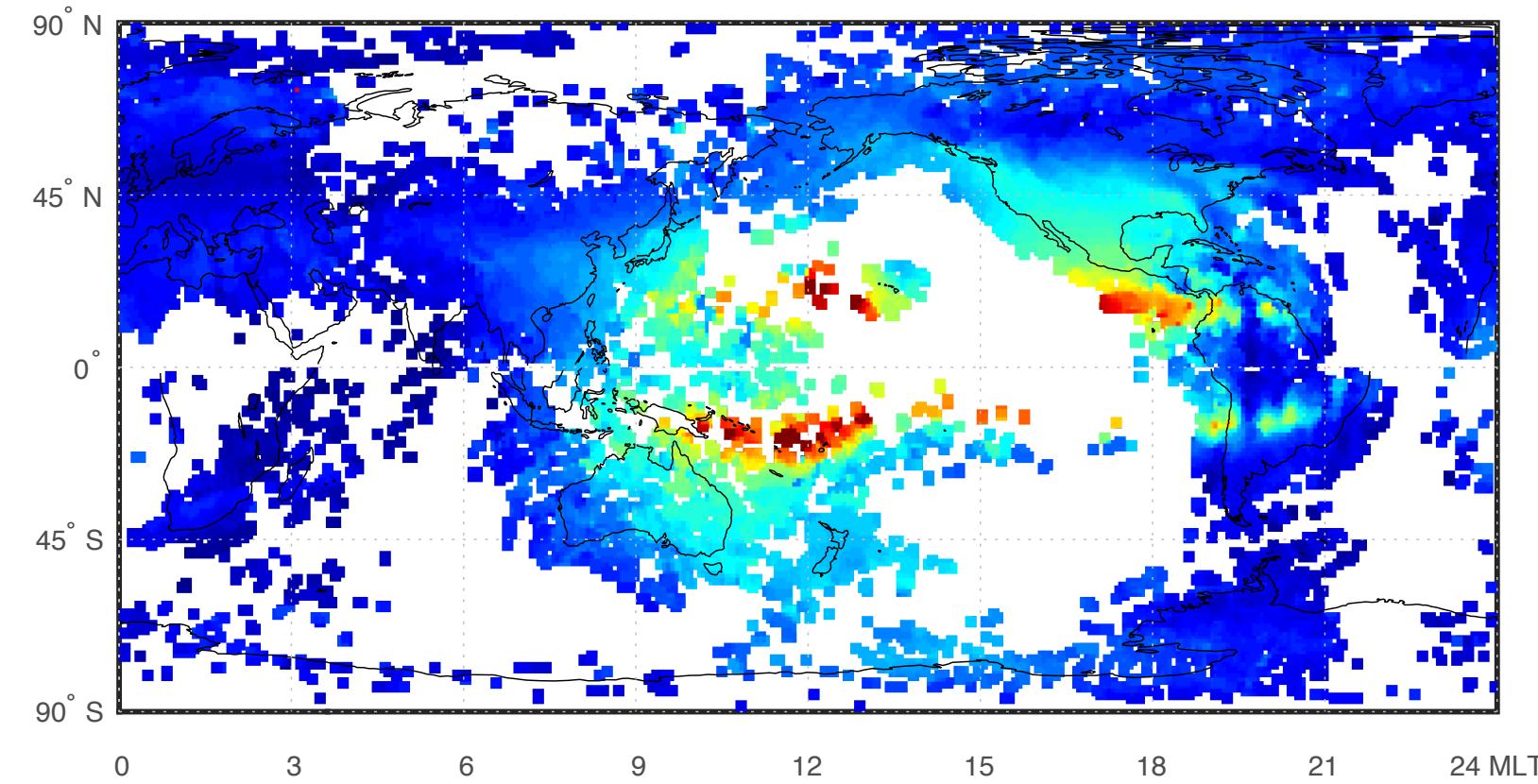
(B) SoftImpute



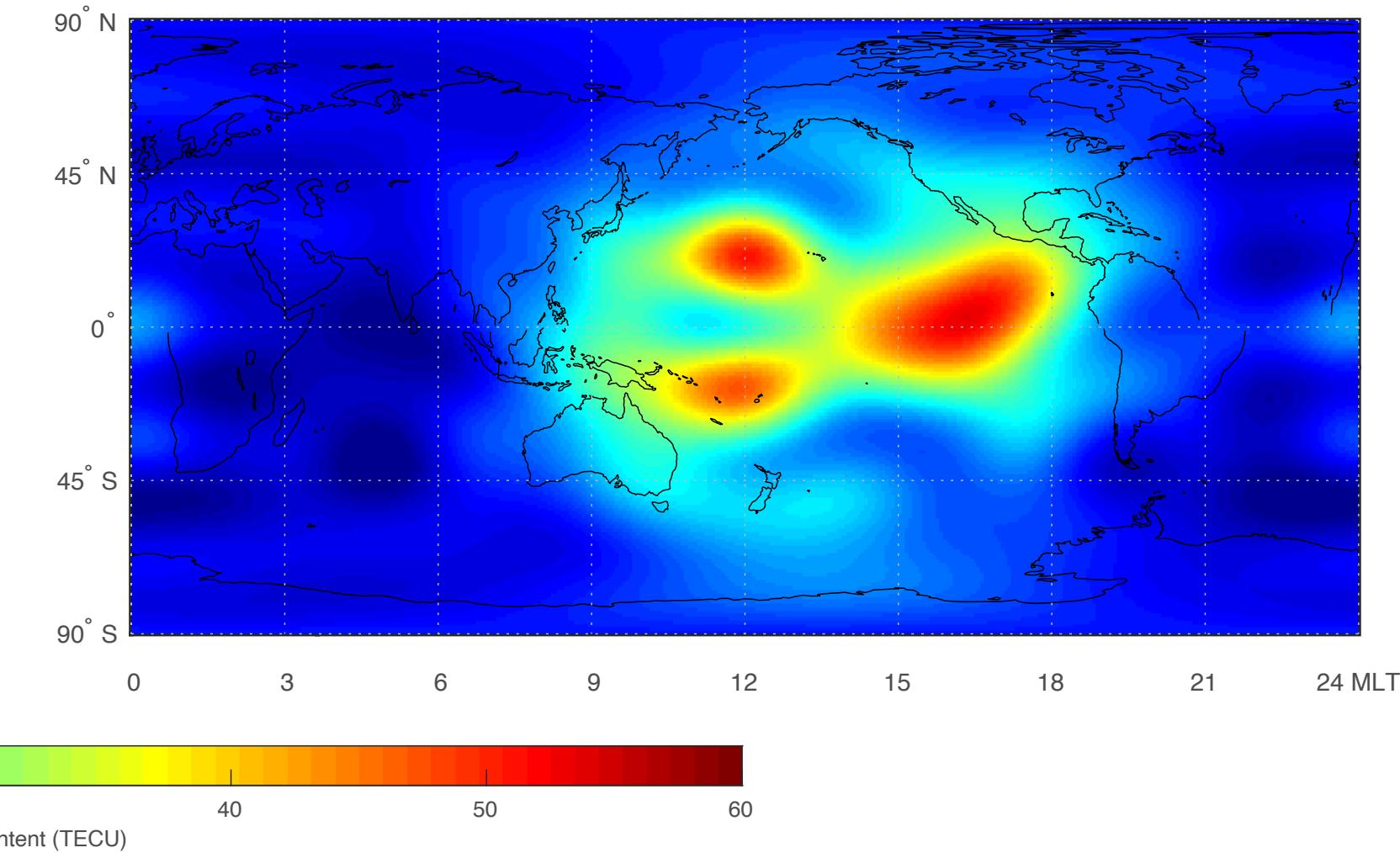
- Traditional matrix/tensor completion method [Hastie et al. (2015)] cannot impute regions with almost all entries missing (called the missing patches).

# Spatio-Temporal Kernel Smoothing

(A) Madrigal TEC map

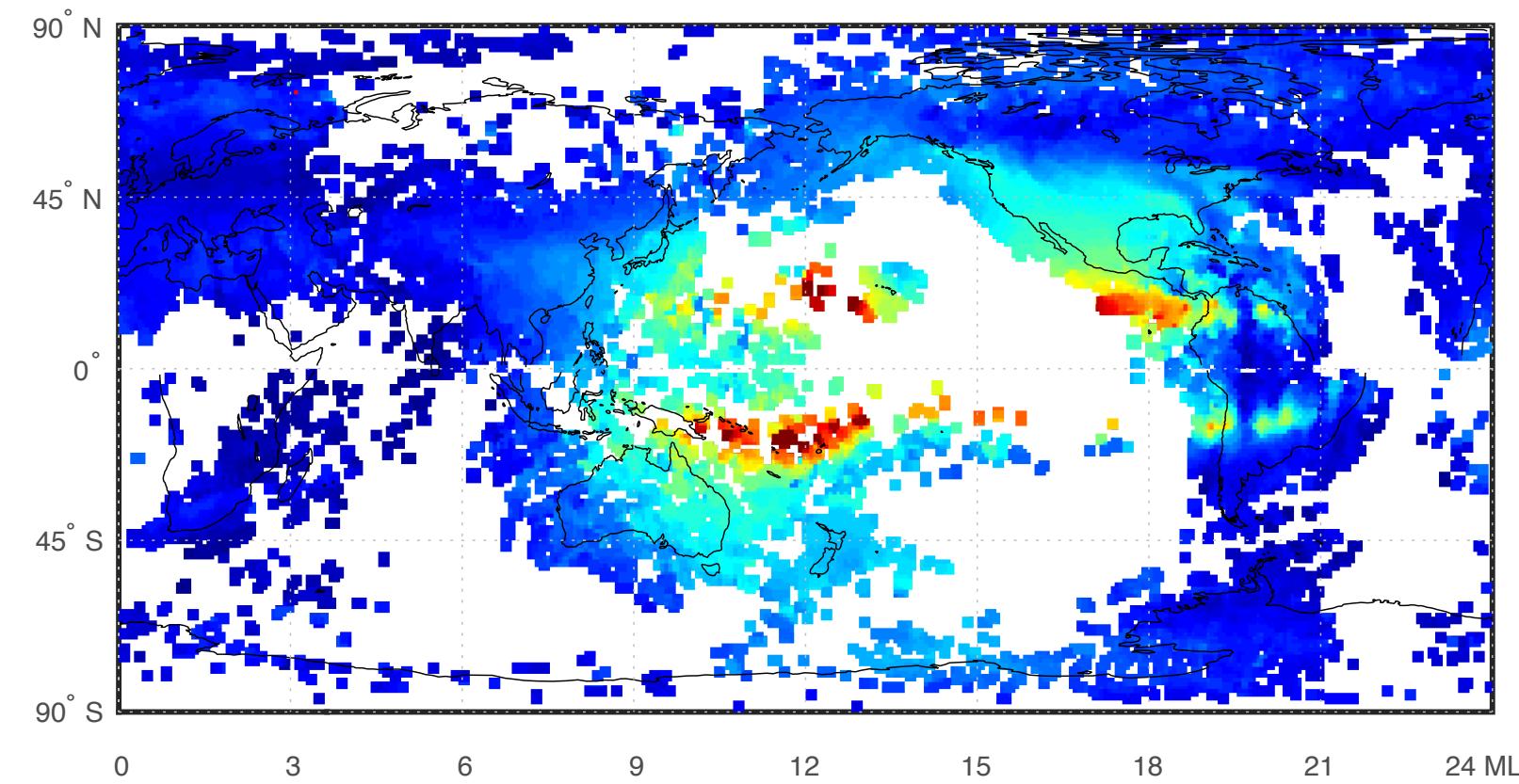


(B) Spherical Harmonics Fitting

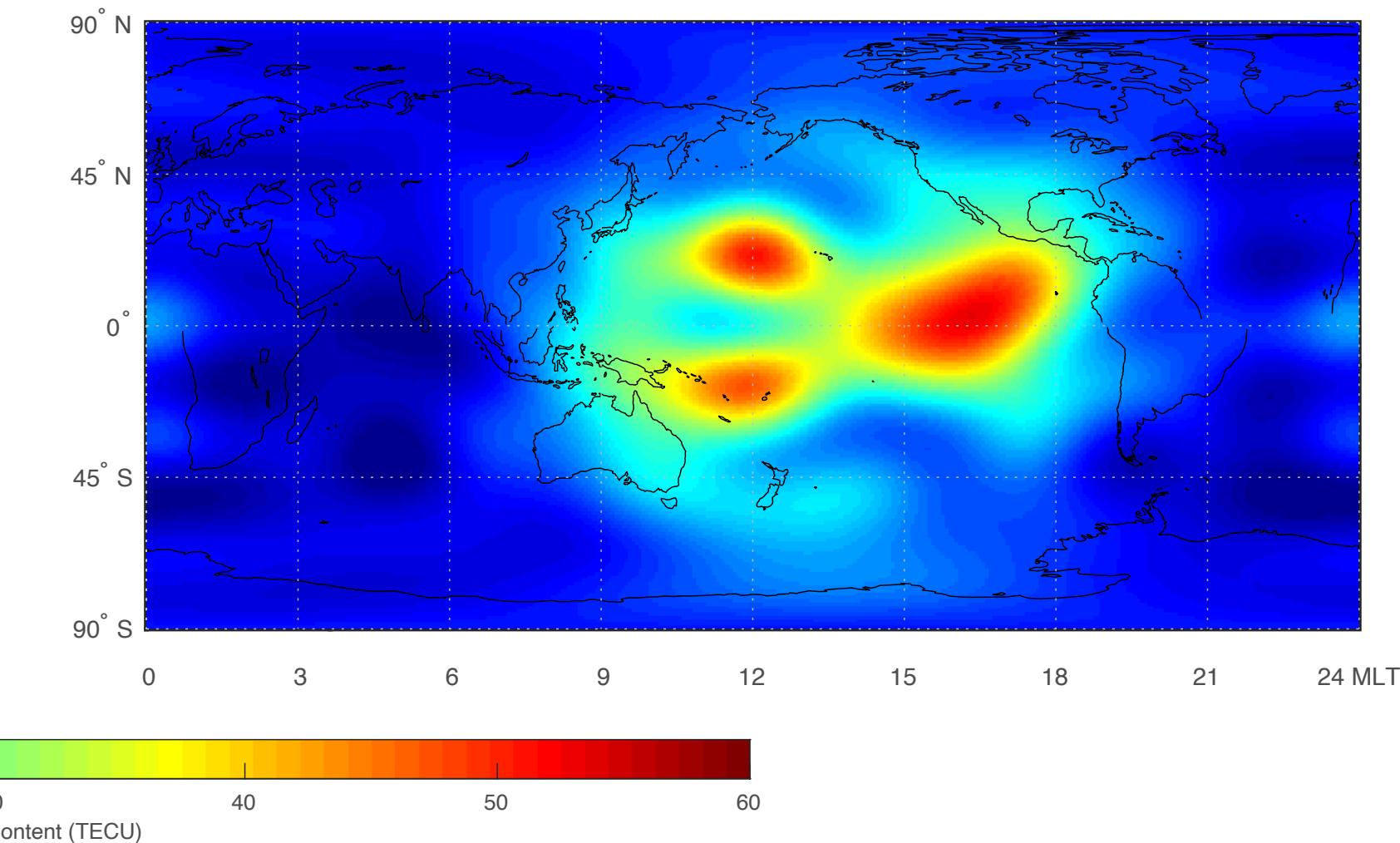


# Spatio-Temporal Kernel Smoothing

(A) Madrigal TEC map



(B) Spherical Harmonics Fitting



- Spatio-temporal kernel smoothing can impute the missing patches with more reasonable values, but tends to be overly-smooth.

# Our Method

# Our Method

- 3-D spatio-temporal tensor (as a matrix time series):  $\mathbf{X}_1, \dots, \mathbf{X}_T \in \mathbb{R}^{M \times N}$

# Our Method

- 3-D spatio-temporal tensor (as a matrix time series):  $\mathbf{X}_1, \dots, \mathbf{X}_T \in \mathbb{R}^{M \times N}$
- We aim at imputing  $\mathbf{X}_1, \dots, \mathbf{X}_T$  via solving the following non-convex optimization problem:

# Our Method

- 3-D spatio-temporal tensor (as a matrix time series):  $\mathbf{X}_1, \dots, \mathbf{X}_T \in \mathbb{R}^{M \times N}$
- We aim at imputing  $\mathbf{X}_1, \dots, \mathbf{X}_T$  via solving the following non-convex optimization problem:

$$\begin{aligned} \min_{\mathbf{A}_t, \mathbf{B}_t, t \in [T]} & \frac{1}{2} \sum_t \left\| \mathbf{P}_{\Omega_t} (\mathbf{X}_t - \mathbf{A}_t \mathbf{B}_t^\top) \right\|_F^2 + \lambda_1 \sum_t \left( \|\mathbf{A}_t\|_F^2 + \|\mathbf{B}_t\|_F^2 \right) \\ & + \lambda_2 \sum_{t \geq 2} \left\| \mathbf{A}_t \mathbf{B}_t^\top - \mathbf{A}_{t-1} \mathbf{B}_{t-1}^\top \right\|_F^2 + \lambda_3 \sum_t \left\| \mathbf{A}_t \mathbf{B}_t^\top - \mathbf{Y}_t \right\|_F^2 \end{aligned}$$

# Our Method

- 3-D spatio-temporal tensor (as a matrix time series):  $\mathbf{X}_1, \dots, \mathbf{X}_T \in \mathbb{R}^{M \times N}$
- We aim at imputing  $\mathbf{X}_1, \dots, \mathbf{X}_T$  via solving the following non-convex optimization problem:

$$\begin{aligned} & \min_{\mathbf{A}_t, \mathbf{B}_t, t \in [T]} \frac{1}{2} \sum_t \left\| \mathbf{P}_{\Omega_t} (\mathbf{X}_t - \mathbf{A}_t \mathbf{B}_t^\top) \right\|_F^2 + \lambda_1 \sum_t \left( \|\mathbf{A}_t\|_F^2 + \|\mathbf{B}_t\|_F^2 \right) \\ & + \lambda_2 \sum_{t \geq 2} \left\| \mathbf{A}_t \mathbf{B}_t^\top - \mathbf{A}_{t-1} \mathbf{B}_{t-1}^\top \right\|_F^2 + \lambda_3 \sum_t \left\| \mathbf{A}_t \mathbf{B}_t^\top - \mathbf{Y}_t \right\|_F^2 \end{aligned}$$

↑  
reconstruction loss

# Our Method

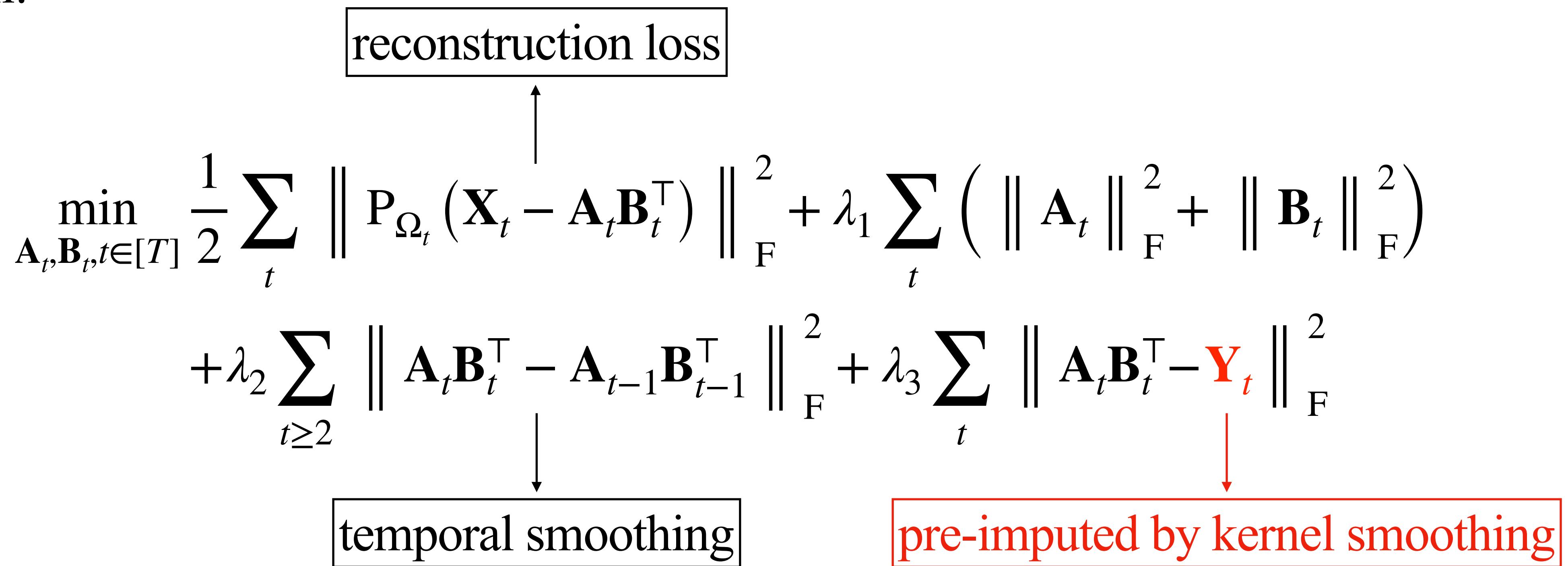
- 3-D spatio-temporal tensor (as a matrix time series):  $\mathbf{X}_1, \dots, \mathbf{X}_T \in \mathbb{R}^{M \times N}$
- We aim at imputing  $\mathbf{X}_1, \dots, \mathbf{X}_T$  via solving the following non-convex optimization problem:

$$\begin{aligned} & \min_{\mathbf{A}_t, \mathbf{B}_t, t \in [T]} \frac{1}{2} \sum_t \left\| \mathbf{P}_{\Omega_t} (\mathbf{X}_t - \mathbf{A}_t \mathbf{B}_t^\top) \right\|_F^2 + \lambda_1 \sum_t \left( \|\mathbf{A}_t\|_F^2 + \|\mathbf{B}_t\|_F^2 \right) \\ & + \lambda_2 \sum_{t \geq 2} \left\| \mathbf{A}_t \mathbf{B}_t^\top - \mathbf{A}_{t-1} \mathbf{B}_{t-1}^\top \right\|_F^2 + \lambda_3 \sum_t \left\| \mathbf{A}_t \mathbf{B}_t^\top - \mathbf{Y}_t \right\|_F^2 \end{aligned}$$

↑  
reconstruction loss  
↓  
temporal smoothing

# Our Method

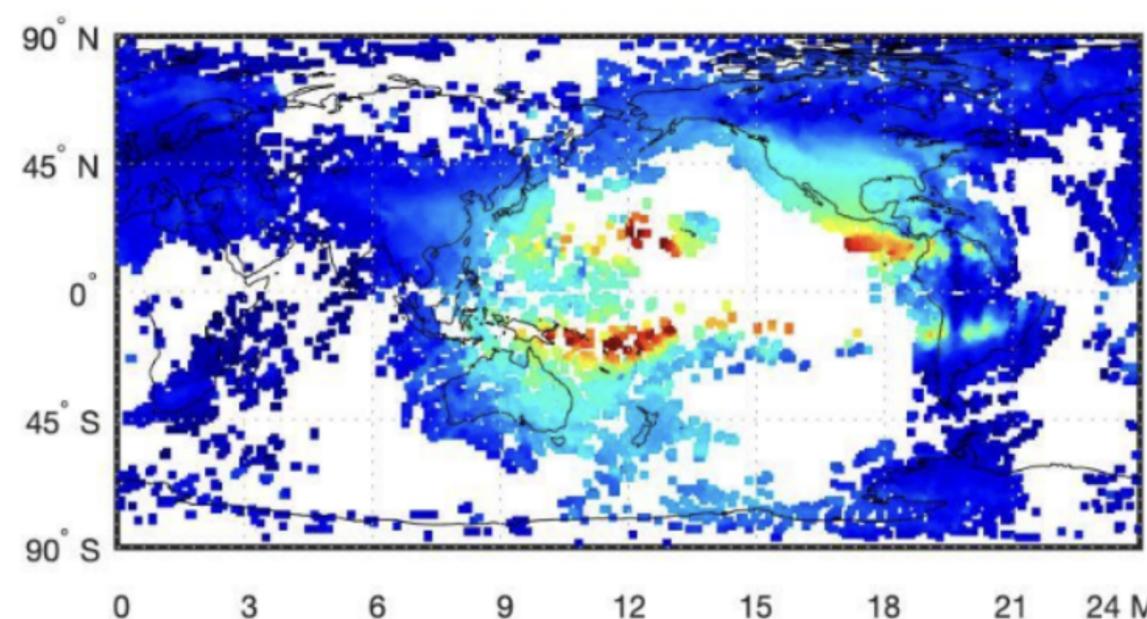
- 3-D spatio-temporal tensor (as a matrix time series):  $\mathbf{X}_1, \dots, \mathbf{X}_T \in \mathbb{R}^{M \times N}$
- We aim at imputing  $\mathbf{X}_1, \dots, \mathbf{X}_T$  via solving the following non-convex optimization problem:



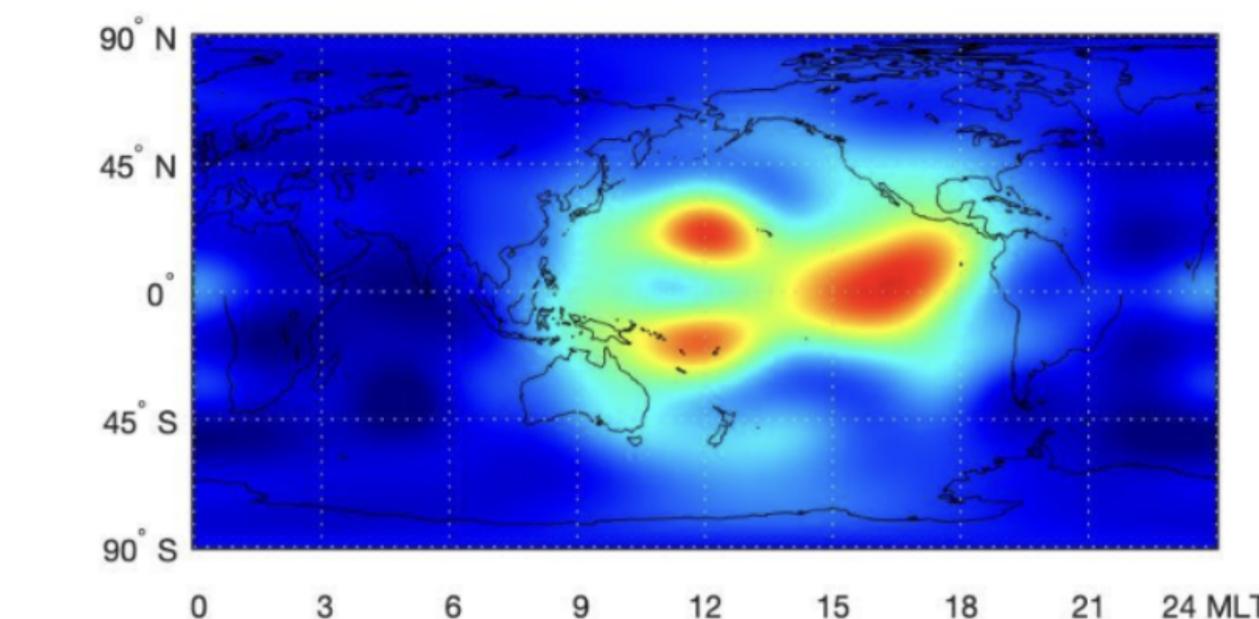
# Results

## Global Total Electron Content (TEC) Reconstruction

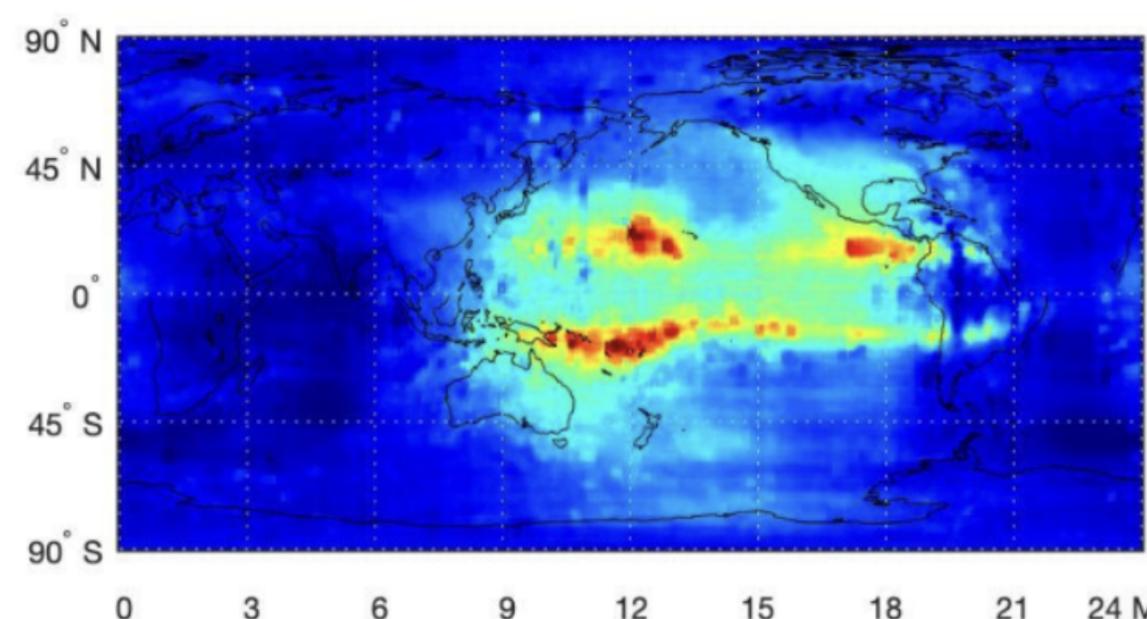
(A) Original Data



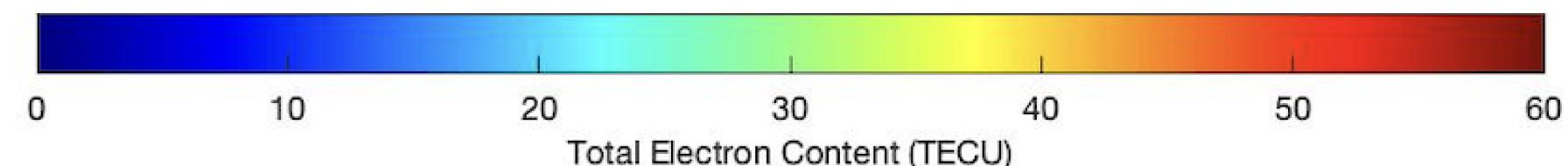
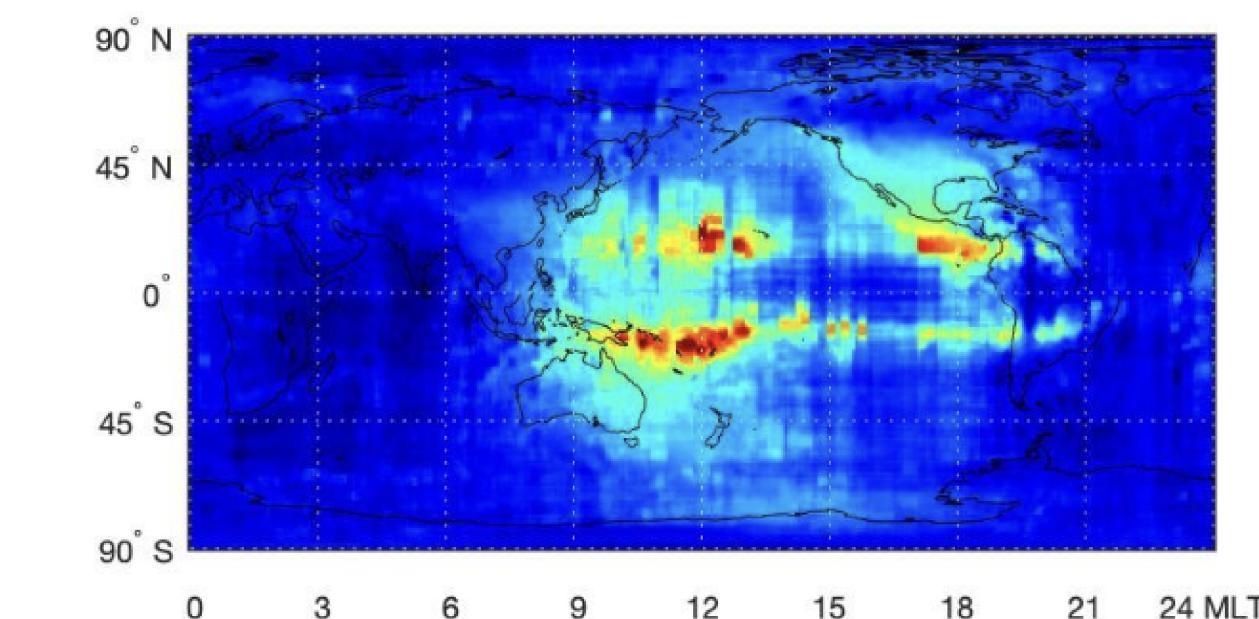
(B) Kernel Smoothing



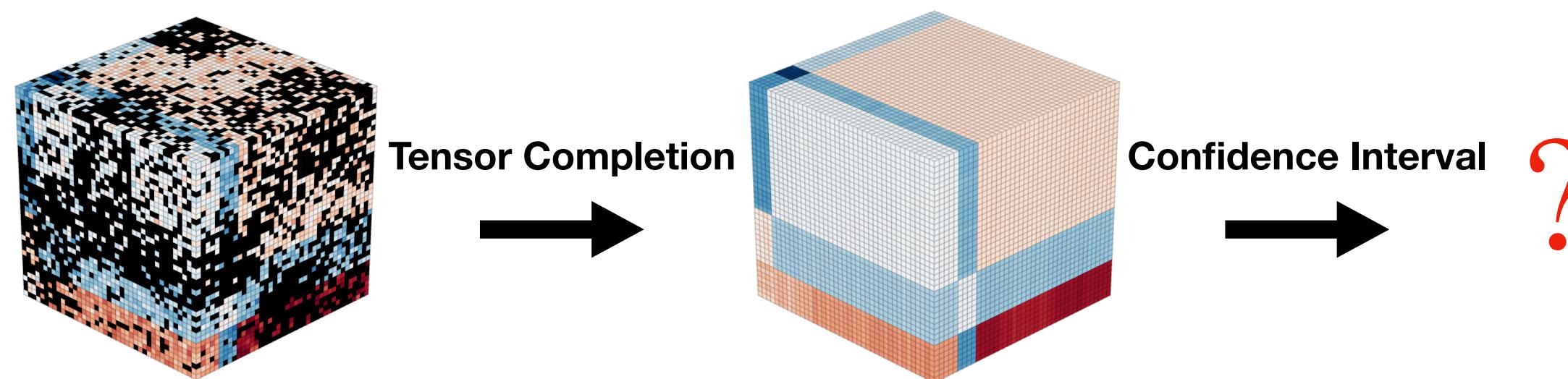
(C) Our Method



(D) Matrix Completion



## II. Conformalized Tensor Completion with Riemannian Optimization

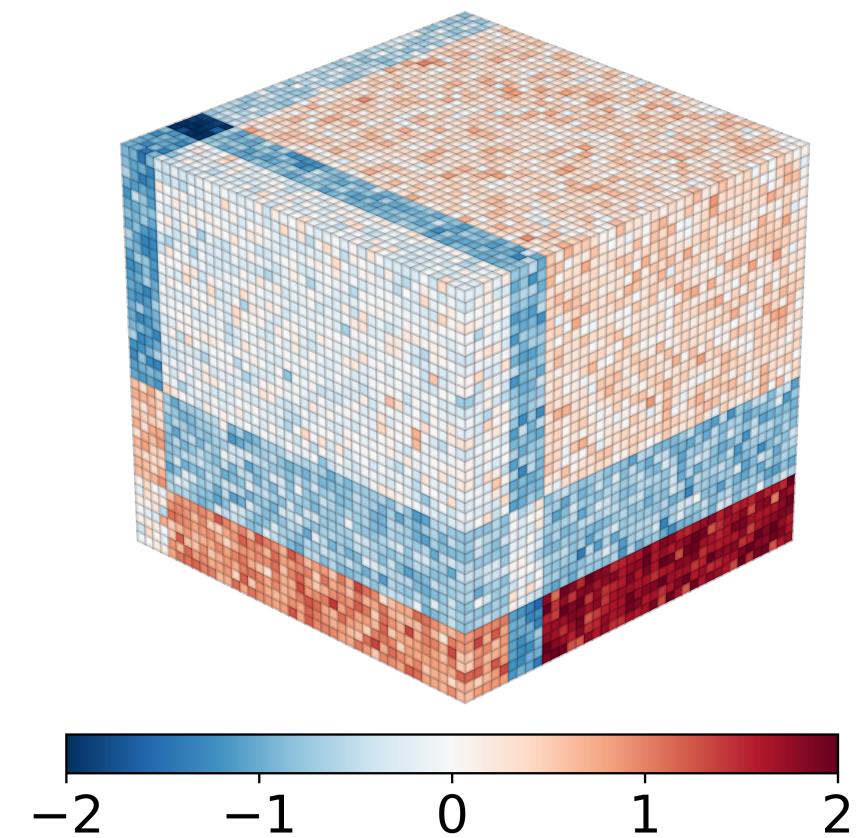


# Research Question

# Research Question

Noisy Data Tensor

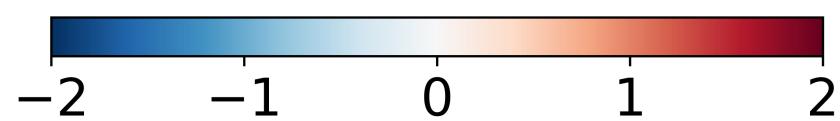
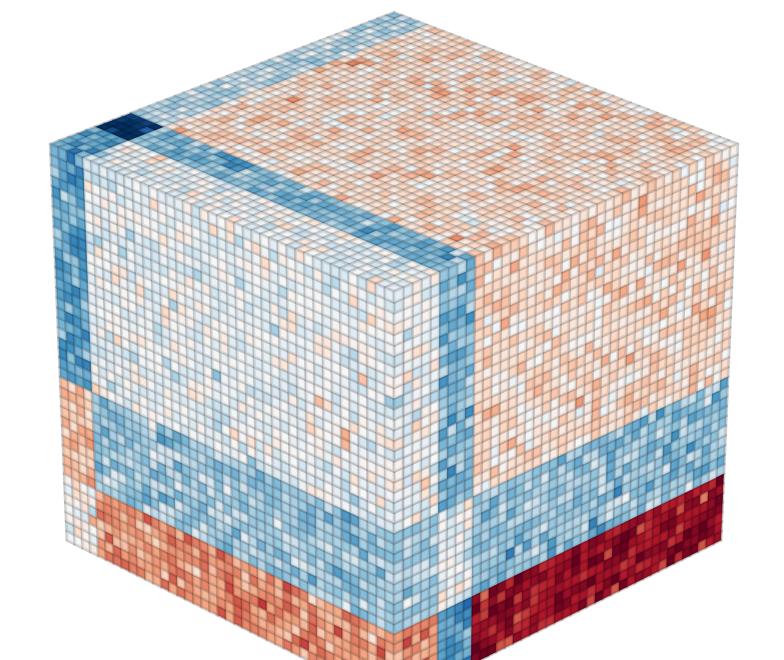
$$\mathcal{X} \in \mathbb{R}^{d_1 \times \dots \times d_K}$$



# Research Question

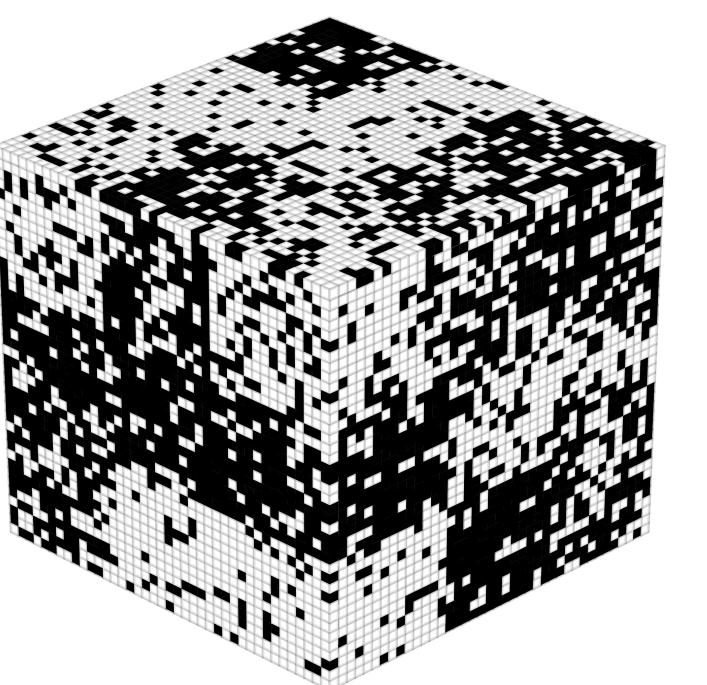
Noisy Data Tensor

$$\mathcal{X} \in \mathbb{R}^{d_1 \times \dots \times d_K}$$

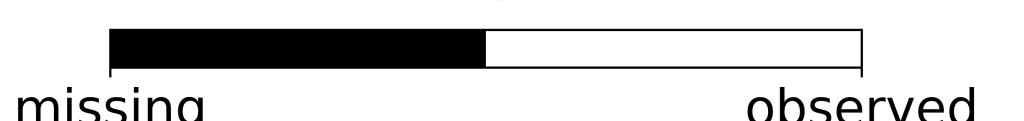


Missing Data

$$\mathcal{W} \in \{-1, 1\}^{d_1 \times \dots \times d_K}$$

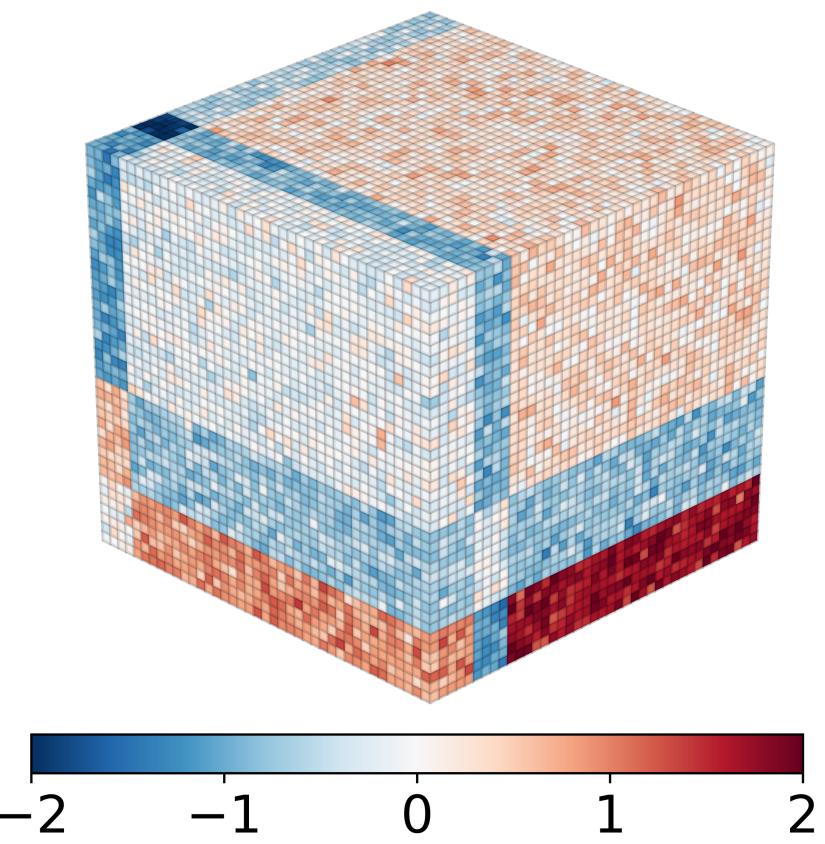


+

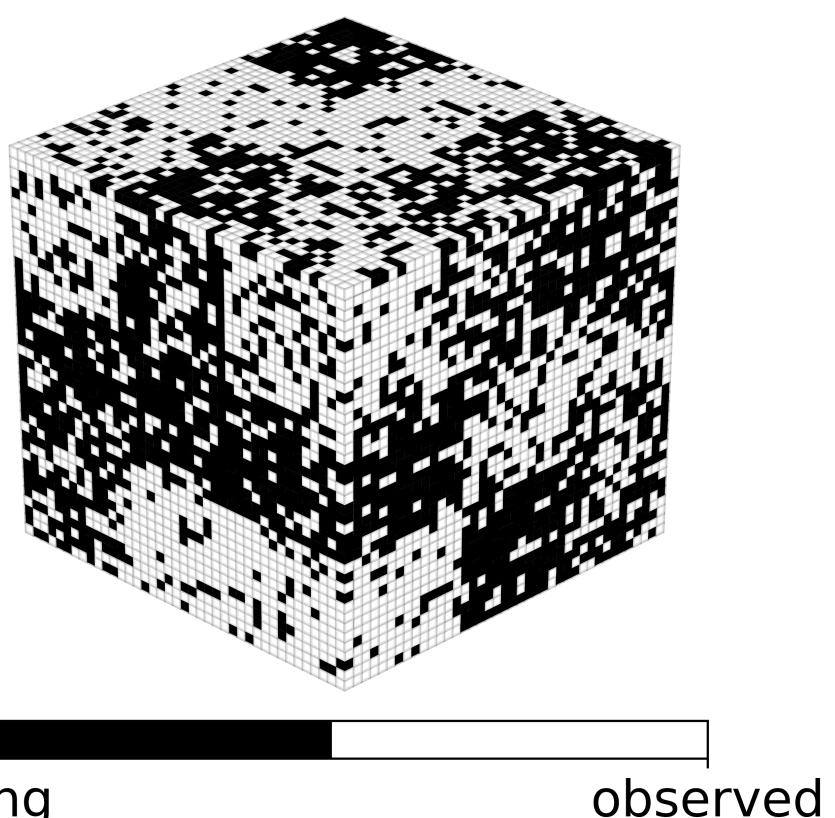


# Research Question

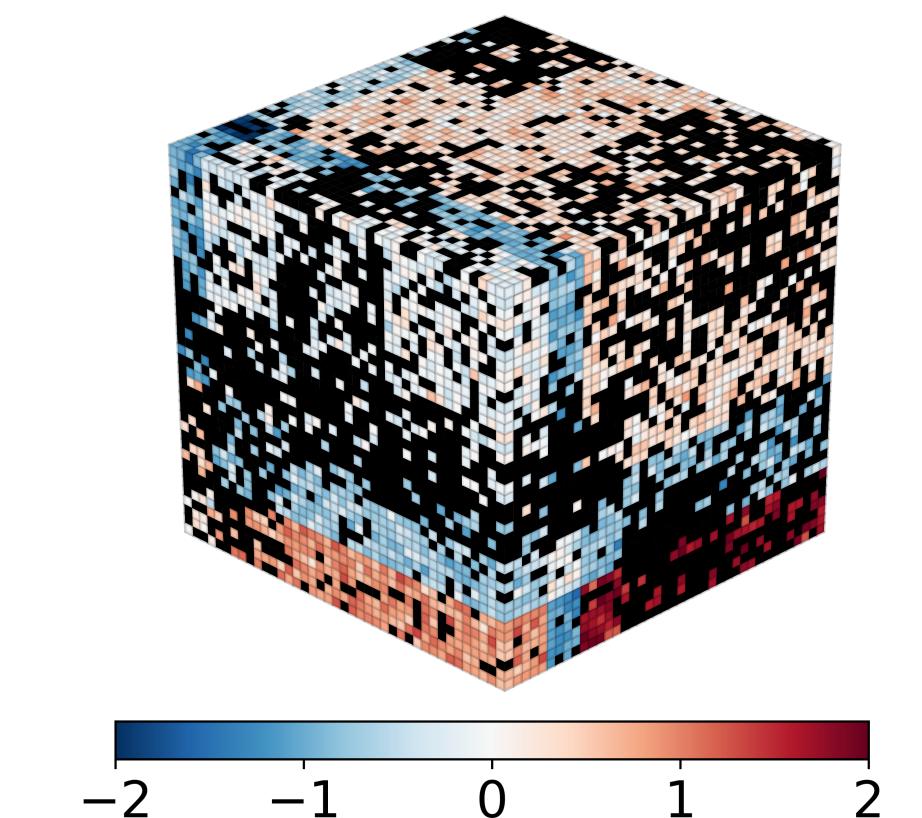
Noisy Data Tensor  
 $\mathcal{X} \in \mathbb{R}^{d_1 \times \dots \times d_K}$



Missing Data  
 $\mathcal{W} \in \{-1,1\}^{d_1 \times \dots \times d_K}$

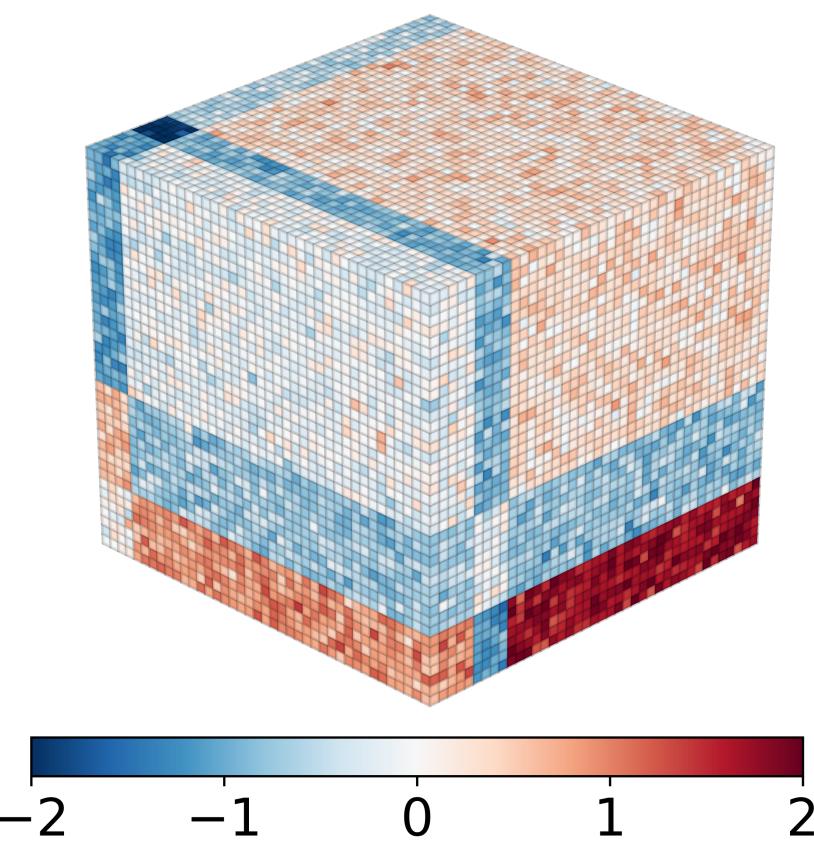


Partially-Observed Data Tensor  $\mathcal{X}_\Omega$   
 $\Omega = \{s \in [d_1] \times \dots \times [d_K] \mid [\mathcal{W}]_s = 1\}$

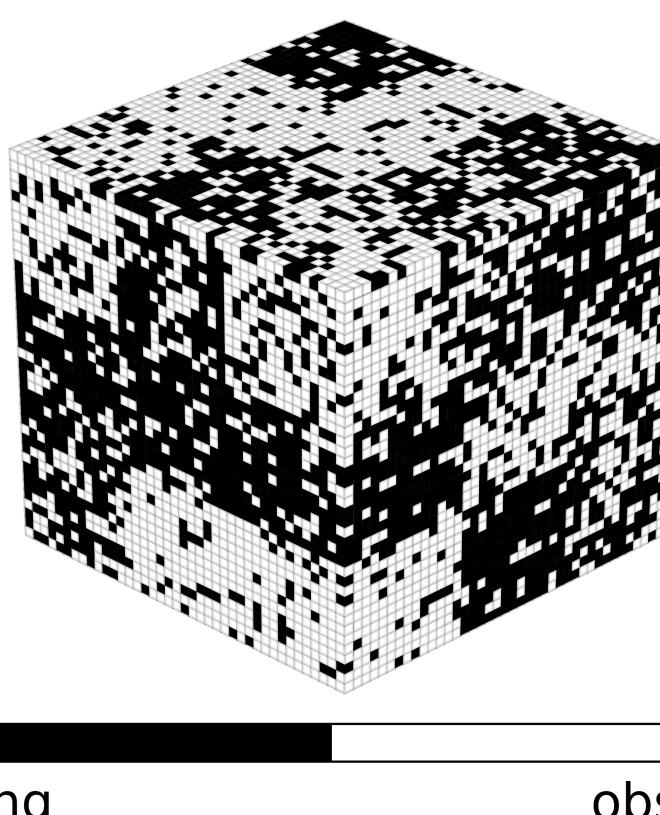


# Research Question

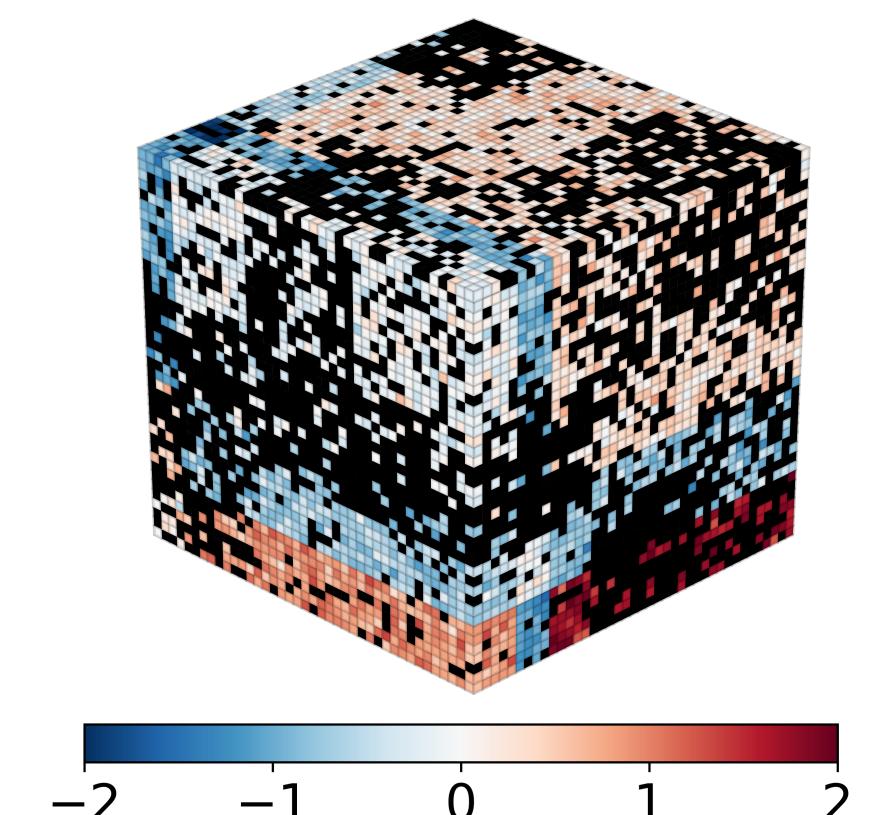
Noisy Data Tensor  
 $\mathcal{X} \in \mathbb{R}^{d_1 \times \dots \times d_K}$



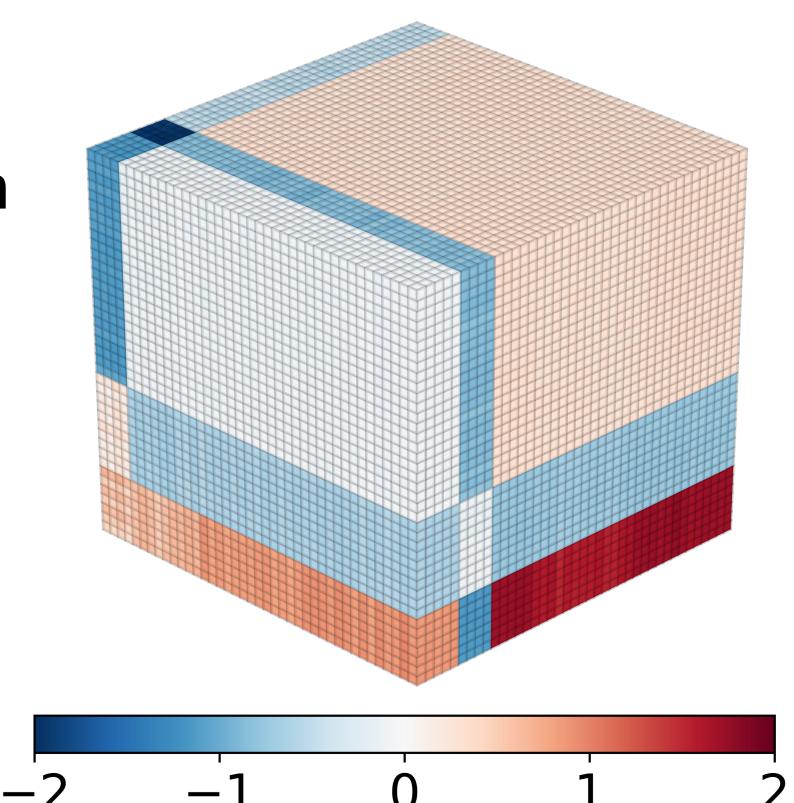
Missing Data  
 $\mathcal{W} \in \{-1, 1\}^{d_1 \times \dots \times d_K}$



Partially-Observed Data Tensor  $\mathcal{X}_\Omega$   
 $\Omega = \{s \in [d_1] \times \dots \times [d_K] \mid [\mathcal{W}]_s = 1\}$



Completed Tensor  $\widehat{\mathcal{X}}$

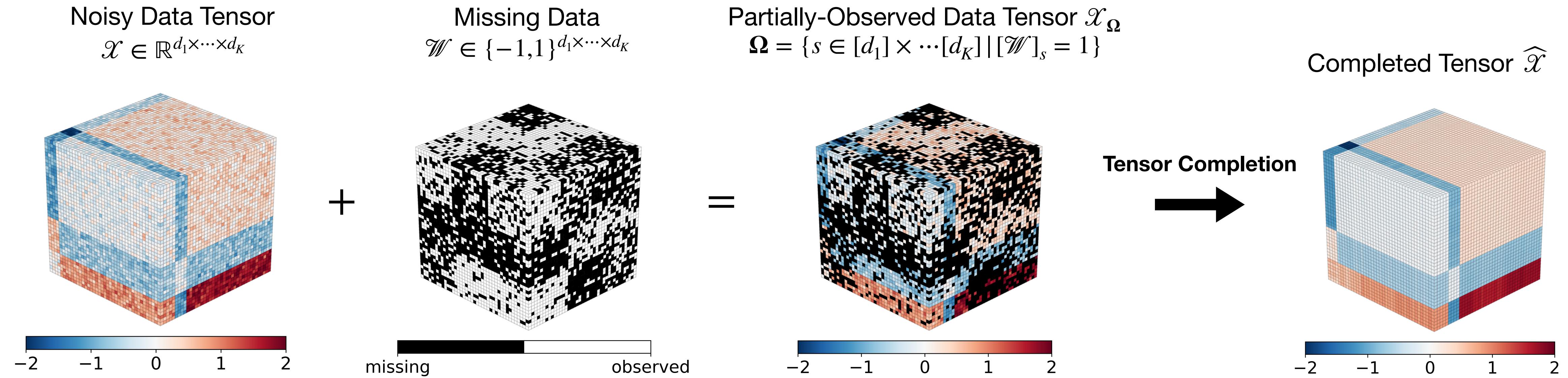


+

=

Tensor Completion

# Research Question



For any  $s^* \in \Omega^c = \{s \in [d_1] \times \dots \times [d_K] | [\mathcal{W}]_s = -1\}$ , how to construct a confidence region  $\widehat{C}_{s^*}$  such that:

$$P\left([\mathcal{X}]_{s^*} \in \widehat{C}_{s^*}\right) \geq 1 - \alpha,$$

with  $1 - \alpha$  being a pre-specified level of confidence?

# **Literature Review**

## **Two Types of Uncertainty Quantification**

# Literature Review

## Two Types of Uncertainty Quantification

Suppose the data tensor  $\mathcal{X}$  is generated by:

# Literature Review

## Two Types of Uncertainty Quantification

Suppose the data tensor  $\mathcal{X}$  is generated by:

$$\mathcal{X} = \mathcal{X}^* + \mathcal{E},$$

# Literature Review

## Two Types of Uncertainty Quantification

Suppose the data tensor  $\mathcal{X}$  is generated by:

$$\mathcal{X} = \mathcal{X}^* + \mathcal{E},$$

existing works construct confidence sets of  $\widehat{\mathcal{X}}$  with two distinct goals:

# Literature Review

## Two Types of Uncertainty Quantification

Suppose the data tensor  $\mathcal{X}$  is generated by:

$$\mathcal{X} = \mathcal{X}^* + \mathcal{E},$$

existing works construct confidence sets of  $\widehat{\mathcal{X}}$  with two distinct goals:

- inference on true tensor:  $P\left([\mathcal{X}^*]_{s^*} \in \widehat{C}_{s^*}\right) \geq 1 - \alpha$

# Literature Review

## Two Types of Uncertainty Quantification

Suppose the data tensor  $\mathcal{X}$  is generated by:

$$\mathcal{X} = \mathcal{X}^* + \mathcal{E},$$

existing works construct confidence sets of  $\widehat{\mathcal{X}}$  with two distinct goals:

- inference on true tensor:  $P\left([\mathcal{X}^*]_{s^*} \in \widehat{C}_{s^*}\right) \geq 1 - \alpha$
- predictive UQ:  $P\left([\mathcal{X}]_{s^*} \in \widehat{C}_{s^*}\right) \geq 1 - \alpha$

# **Literature Review**

## **Inference on True Tensor/Matrix**

# Literature Review

## Inference on True Tensor/Matrix

- Convert bound on  $\|\widehat{\mathcal{X}} - \mathcal{X}^*\|_F$  to confidence sets [Carpentier et al. (2018, 2019)]

# Literature Review

## Inference on True Tensor/Matrix

- Convert bound on  $\|\widehat{\mathcal{X}} - \mathcal{X}^*\|_F$  to confidence sets [Carpentier et al. (2018, 2019)]
  - ▶ relies on specific algorithm and can be extremely loose

# Literature Review

## Inference on True Tensor/Matrix

- Convert bound on  $\|\widehat{\mathcal{X}} - \mathcal{X}^*\|_F$  to confidence sets [Carpentier et al. (2018, 2019)]
  - ▶ relies on specific algorithm and can be extremely loose
- Distributional analysis of  $[\widehat{\mathcal{X}} - \mathcal{X}^*]_{s^*}$  [Chen et al. (2019), Xia and Yuan (2021), Cai et al. (2022), Farias et al. (2022)]

# Literature Review

## Inference on True Tensor/Matrix

- Convert bound on  $\|\widehat{\mathcal{X}} - \mathcal{X}^*\|_F$  to confidence sets [Carpentier et al. (2018, 2019)]
  - ▶ relies on specific algorithm and can be extremely loose
- Distributional analysis of  $[\widehat{\mathcal{X}} - \mathcal{X}^*]_{s^*}$  [Chen et al. (2019), Xia and Yuan (2021), Cai et al. (2022), Farias et al. (2022)]
  - ▶ relies on specific algorithm, noise distribution, and exact low-rank true tensor form

# Literature Review

## Inference on True Tensor/Matrix

- Convert bound on  $\|\widehat{\mathcal{X}} - \mathcal{X}^*\|_F$  to confidence sets [Carpentier et al. (2018, 2019)]
  - ▶ relies on specific algorithm and can be extremely loose
- Distributional analysis of  $[\widehat{\mathcal{X}} - \mathcal{X}^*]_{s^*}$  [Chen et al. (2019), Xia and Yuan (2021), Cai et al. (2022), Farias et al. (2022)]
  - ▶ relies on specific algorithm, noise distribution, and exact low-rank true tensor form
- Bayesian matrix completion [Mai and Alquier (2015), Yuchi et al. (2023)]

# Literature Review

## Inference on True Tensor/Matrix

- Convert bound on  $\|\widehat{\mathcal{X}} - \mathcal{X}^*\|_F$  to confidence sets [Carpentier et al. (2018, 2019)]
  - ▶ relies on specific algorithm and can be extremely loose
- Distributional analysis of  $[\widehat{\mathcal{X}} - \mathcal{X}^*]_{s^*}$  [Chen et al. (2019), Xia and Yuan (2021), Cai et al. (2022), Farias et al. (2022)]
  - ▶ relies on specific algorithm, noise distribution, and exact low-rank true tensor form
- Bayesian matrix completion [Mai and Alquier (2015), Yuchi et al. (2023)]
- Assume  $[\mathcal{W}]_s \stackrel{i.i.d.}{\sim} \text{Bern}(p)$ , i.e. data missing uniformly at random.

# **Literature Review**

## **Predictive Inference of Matrix Completion**

# Literature Review

## Predictive Inference of Matrix Completion

- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]

# Literature Review

## Predictive Inference of Matrix Completion

- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]
  - ▶ do not have coverage guarantees

# Literature Review

## Predictive Inference of Matrix Completion

- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]
  - ▶ do not have coverage guarantees
- Conformal prediction based matrix completion

# Literature Review

## Predictive Inference of Matrix Completion

- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]
  - ▶ do not have coverage guarantees
- Conformal prediction based matrix completion
  - ▶ data is missing independently, i.e.  $[\mathcal{W}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(p_s)$  [Gui et al. (2023)]

# Literature Review

## Predictive Inference of Matrix Completion

- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]
  - ▶ do not have coverage guarantees
- Conformal prediction based matrix completion
  - ▶ data is missing independently, i.e.  $[\mathcal{W}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(p_s)$  [Gui et al. (2023)]
  - ▶ considers uncertainty of newly-arrived rows [Shao and Zhang (2023)]

# Literature Review

## Predictive Inference of Matrix Completion

- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]
  - ▶ do not have coverage guarantees
- Conformal prediction based matrix completion
  - ▶ data is missing independently, i.e.  $[\mathcal{W}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(p_s)$  [Gui et al. (2023)]
  - ▶ considers uncertainty of newly-arrived rows [Shao and Zhang (2023)]
- Not directly transferrable from matrix to tensor setting.

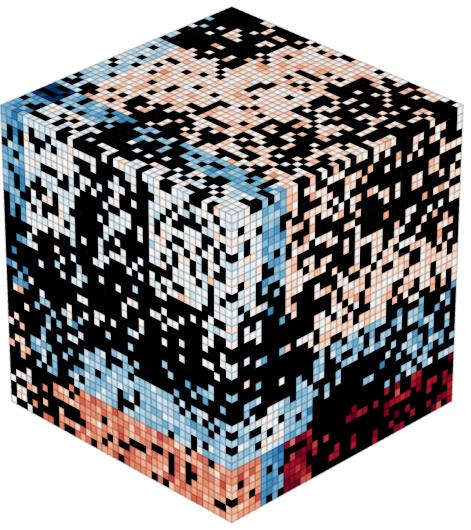
# Literature Review

## Predictive Inference of Matrix Completion

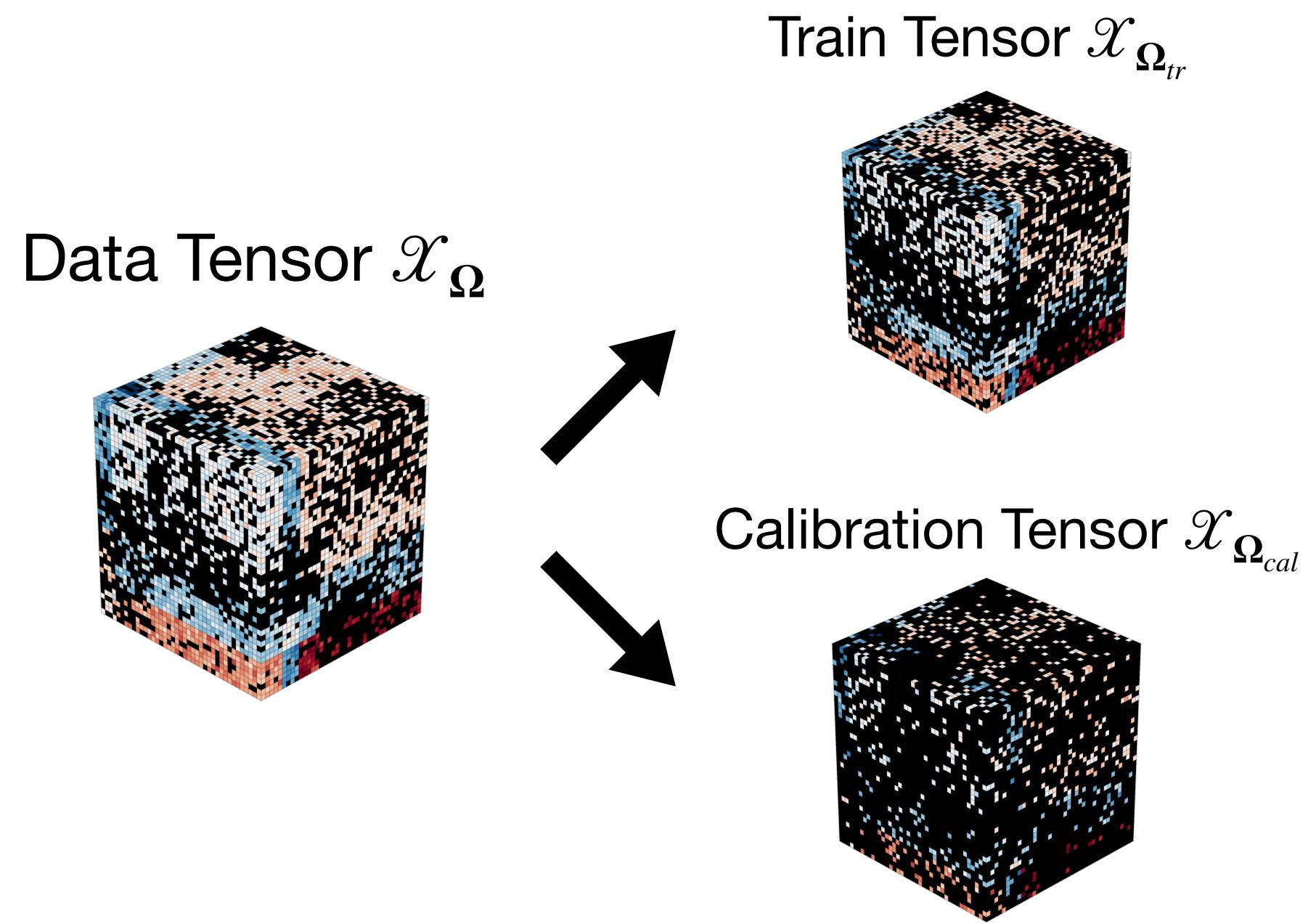
- Deep learning based matrix completion [Zeldes et al. (2017), Kasalicky et al. (2023)]
  - ▶ do not have coverage guarantees
- Conformal prediction based matrix completion
  - ▶ data is missing independently, i.e.  $[\mathcal{W}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(p_s)$  [Gui et al. (2023)]
  - ▶ considers uncertainty of newly-arrived rows [Shao and Zhang (2023)]
- Not directly transferrable from matrix to tensor setting.
- No previous work considers a missing mechanism under the spatio-temporal context.

# Classical Conformal Prediction (CP)

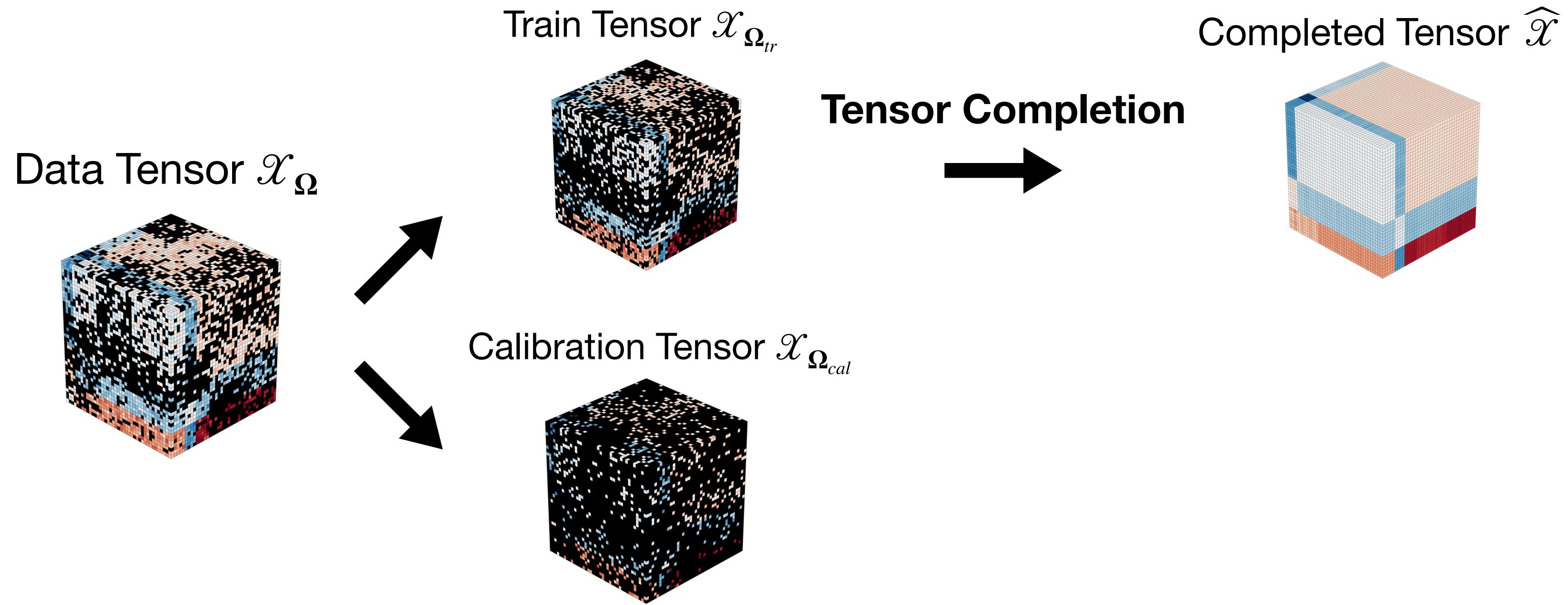
Data Tensor  $\mathcal{X}_\Omega$



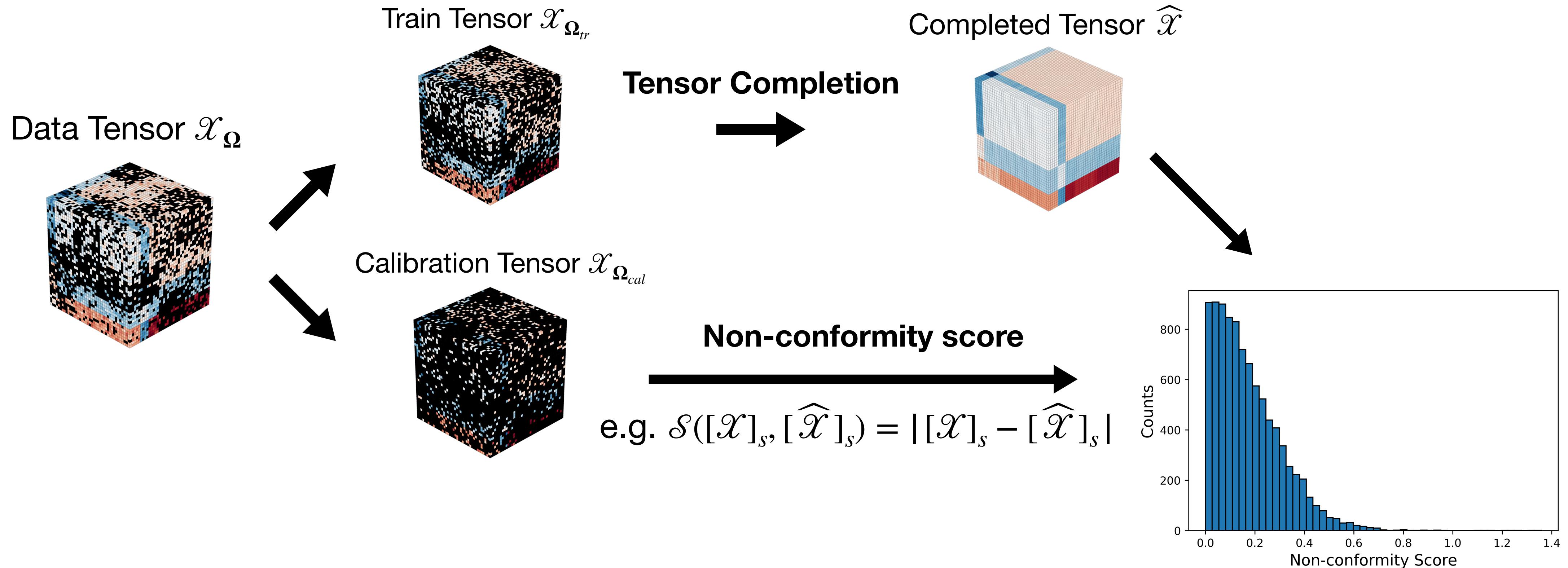
# Classical Conformal Prediction (CP)



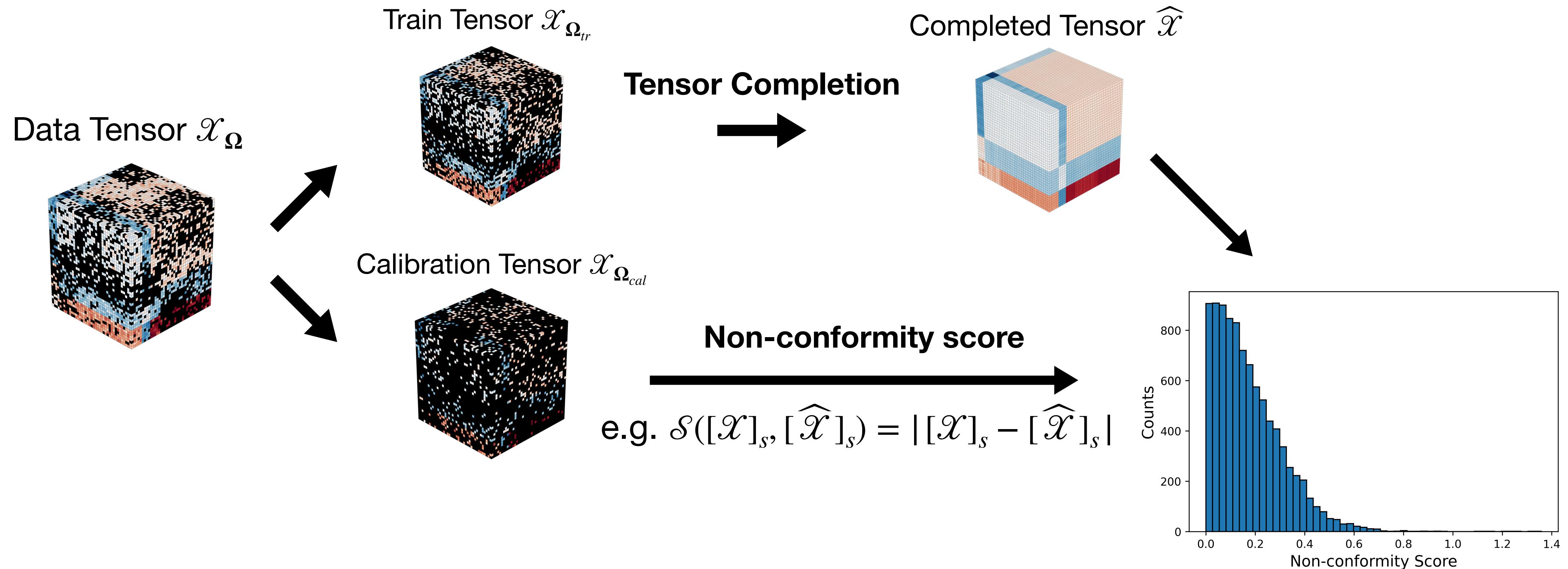
# Classical Conformal Prediction (CP)



# Classical Conformal Prediction (CP)

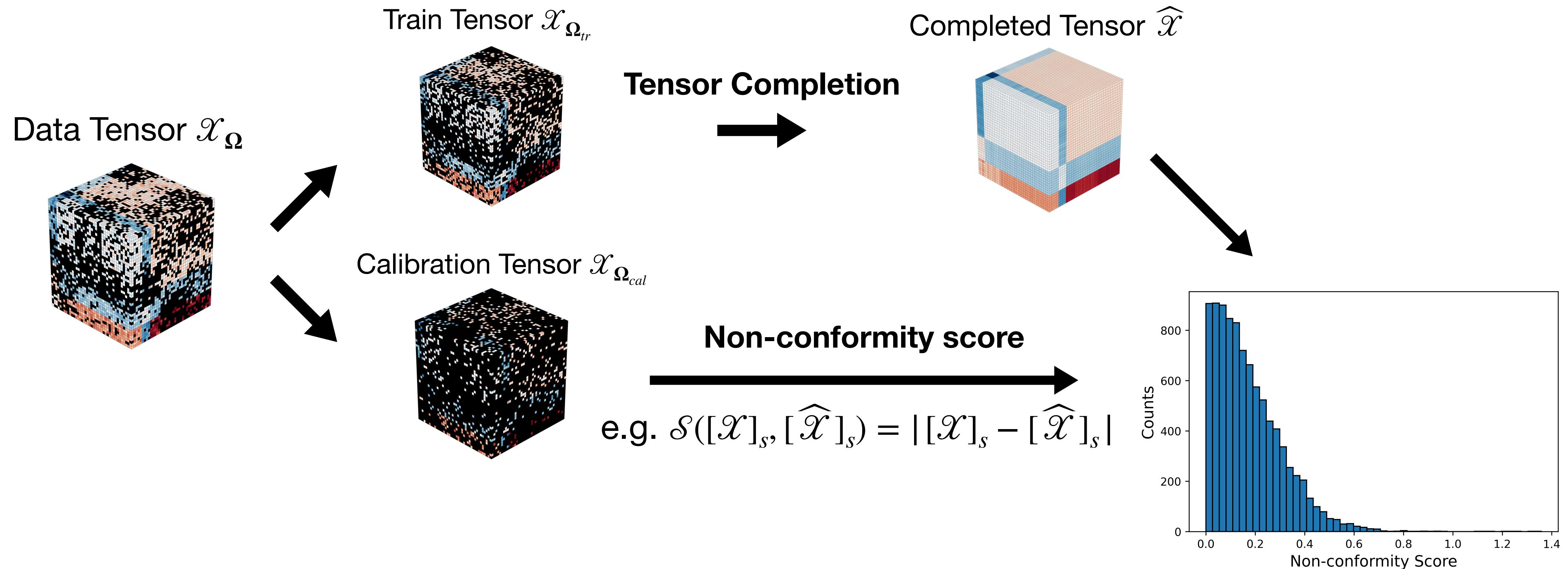


# Classical Conformal Prediction (CP)



$$\widehat{\mathcal{C}}_{1-\alpha,s^*} = \left\{ x \in \mathbb{R} \mid \mathcal{S}(x, [\widehat{\mathcal{X}}]_s) \leq \widehat{q}_{1-\alpha} \left( \sum_{s \in \Omega_{cal}} \frac{1}{|\Omega_{cal}| + 1} \delta_{\mathcal{S}([\mathcal{X}]_s, [\widehat{\mathcal{X}}]_s)} + \frac{1}{|\Omega_{cal}| + 1} \delta_{+\infty} \right) \right\}$$

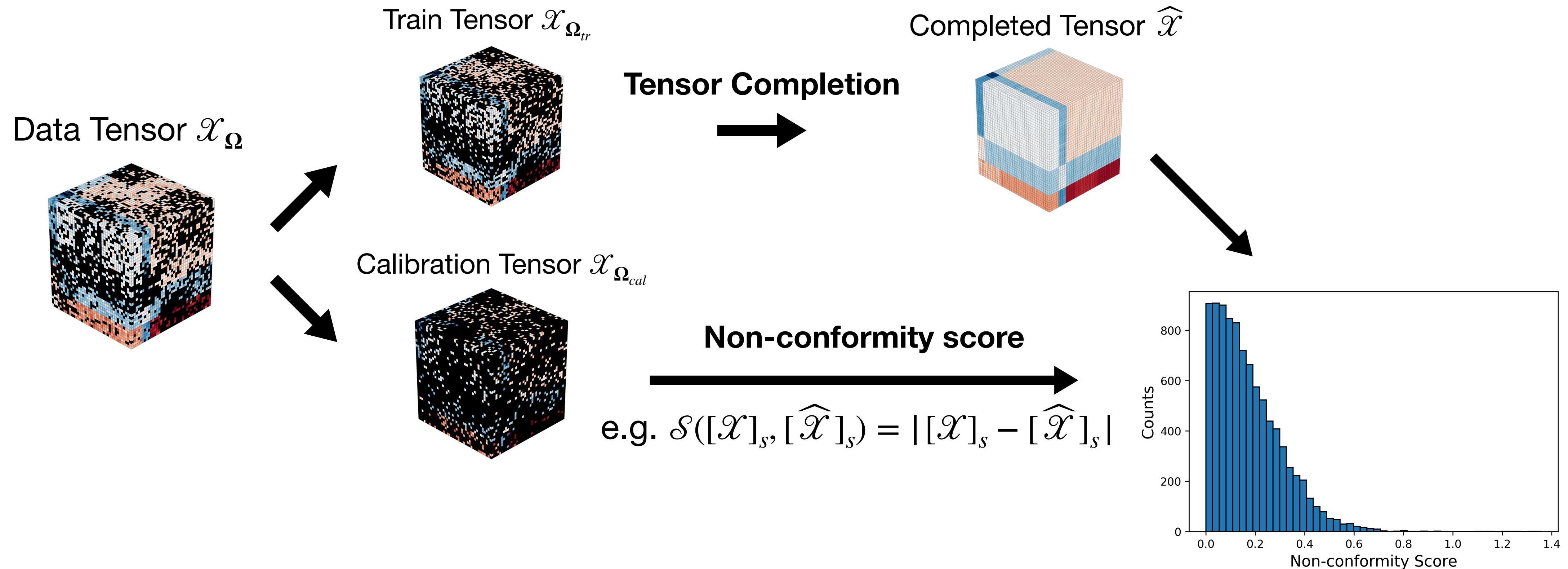
# Classical Conformal Prediction (CP)



$$\widehat{\mathcal{C}}_{1-\alpha,s^*} = \left\{ x \in \mathbb{R} \mid \mathcal{S}(x, [\widehat{\mathcal{X}}]_s) \leq \widehat{q}_{1-\alpha} \left( \sum_{s \in \Omega_{cal}} \frac{1}{|\Omega_{cal}| + 1} \delta_{\mathcal{S}([\mathcal{X}]_s, [\widehat{\mathcal{X}}]_s)} + \frac{1}{|\Omega_{cal}| + 1} \delta_{+\infty} \right) \right\}$$

$\downarrow$   
 $(1 - \alpha)$  level conformal interval at  $s^*$

# Classical Conformal Prediction (CP)

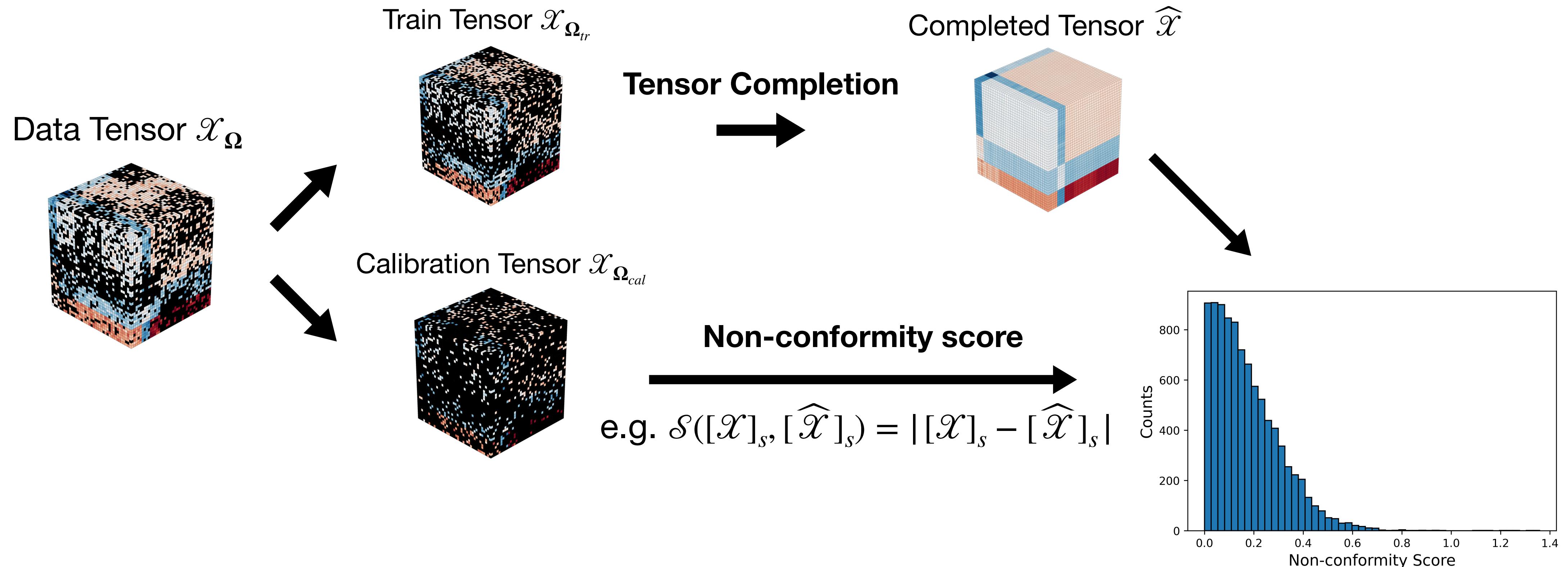


$$\widehat{\mathcal{C}}_{1-\alpha,s^*} = \left\{ x \in \mathbb{R} \mid \mathcal{S}(x, [\widehat{\mathcal{X}}]_s) \leq \widehat{q}_{1-\alpha} \left( \sum_{s \in \Omega_{cal}} \frac{1}{|\Omega_{cal}| + 1} \delta_{\mathcal{S}([\mathcal{X}]_s, [\widehat{\mathcal{X}}]_s)} + \frac{1}{|\Omega_{cal}| + 1} \delta_{+\infty} \right) \right\}$$

$(1 - \alpha)$  level conformal interval at  $s^*$

$(1 - \alpha)$  quantile

# Classical Conformal Prediction (CP)



$$\widehat{\mathcal{C}}_{1-\alpha,s^*} = \left\{ x \in \mathbb{R} \mid \mathcal{S}(x, [\widehat{\mathcal{X}}]_s) \leq \widehat{q}_{1-\alpha} \left( \sum_{s \in \Omega_{cal}} \frac{1}{|\Omega_{cal}| + 1} \delta_{\mathcal{S}([\mathcal{X}]_s, [\widehat{\mathcal{X}}]_s)} + \frac{1}{|\Omega_{cal}| + 1} \delta_{+\infty} \right) \right\}$$

$(1 - \alpha)$  level conformal interval at  $s^*$

$(1 - \alpha)$  quantile

eCDF of the non-conformity score

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

- $[\mathcal{W}]_s \stackrel{i.i.d.}{\sim} \text{Bern}(p)$

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

- $[\mathcal{W}]_s \stackrel{i.i.d.}{\sim} \text{Bern}(p)$
- $s_1, \dots, s_n$ : indices of the  $n$  calibration entries

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

- $[\mathcal{W}]_s \stackrel{i.i.d.}{\sim} \text{Bern}(p)$
- $s_1, \dots, s_n$ : indices of the  $n$  calibration entries
- $s_{n+1}$ : index of a test set entry

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

- $[\mathcal{W}]_s \stackrel{i.i.d.}{\sim} \text{Bern}(p)$
- $s_1, \dots, s_n$ : indices of the  $n$  calibration entries
- $s_{n+1}$ : index of a test set entry
- $v_1, \dots, v_n, v_{n+1}$ : non-conformity scores

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

- $[\mathcal{W}]_s \stackrel{i.i.d.}{\sim} \text{Bern}(p)$
- $s_1, \dots, s_n$ : indices of the  $n$  calibration entries
- $s_{n+1}$ : index of a test set entry
- $v_1, \dots, v_n, v_{n+1}$ : non-conformity scores
- $V$ : the random variable of the non-conformity score at an unseen entry

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

$$(s_1, s_2, \dots, s_n, s_{n+1}) \quad V = v_1$$

$$(s_1, s_2, \dots, s_n, s_{n+1}) \quad V = v_2$$

.....

.....

$$(s_1, s_2, \dots, s_n, s_{n+1}) \quad V = v_{n+1}$$

# Validity of Classic Conformal Prediction

Tibshirani et al. (2019)

$$(\textcolor{red}{s_1}, s_2, \dots, s_n, s_{n+1}) \quad V = v_1$$

$$(s_1, \textcolor{red}{s_2}, \dots, s_n, s_{n+1}) \quad V = v_2$$

.....

.....

$$(s_1, s_2, \dots, s_n, \textcolor{red}{s_{n+1}}) \quad V = v_{n+1}$$

$$P(V = v_k) = \frac{P([\mathcal{W}]_s = 1, \forall s \in \Omega_{tr} \cup \{s_1, \dots, s_{n+1}\} \setminus \{\textcolor{red}{s_k}\}, [\mathcal{W}]_s = -1, \text{o.w.})}{\sum_{l=1}^{n+1} P([\mathcal{W}]_s = 1, \forall s \in \Omega_{tr} \cup \{s_1, \dots, s_{n+1}\} \setminus \{\textcolor{red}{s_l}\}, [\mathcal{W}]_s = -1, \text{o.w.})} = \frac{1}{n+1}$$

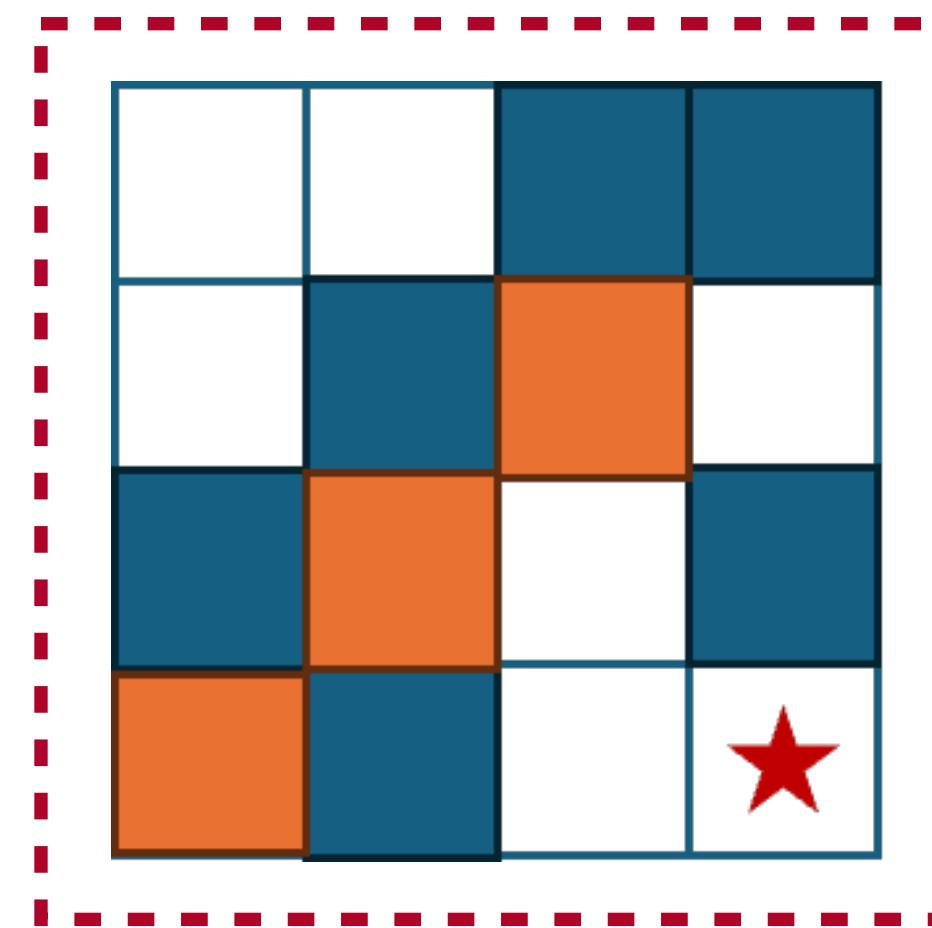
# **Weighted Conformal Prediction**

## **Leave-one-out Likelihood Weighting**

# Weighted Conformal Prediction

## Leave-one-out Likelihood Weighting

■ : training   ■ : calibration   □ : missing   ★ : testing

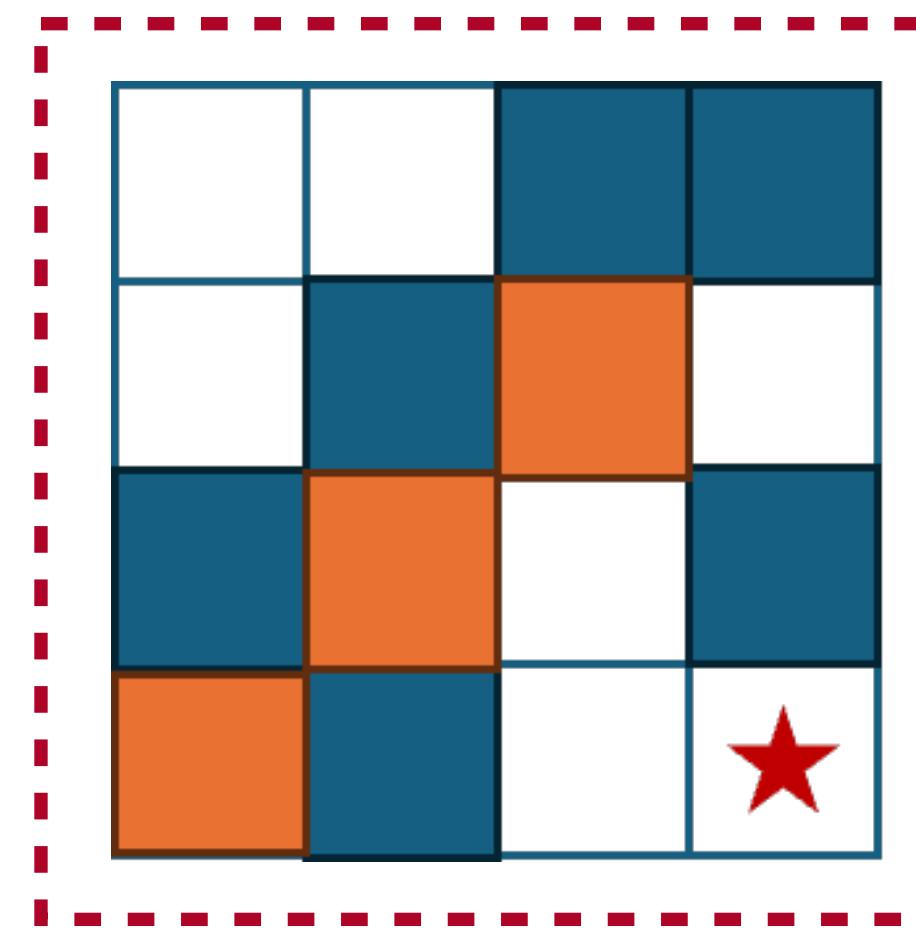


Observation

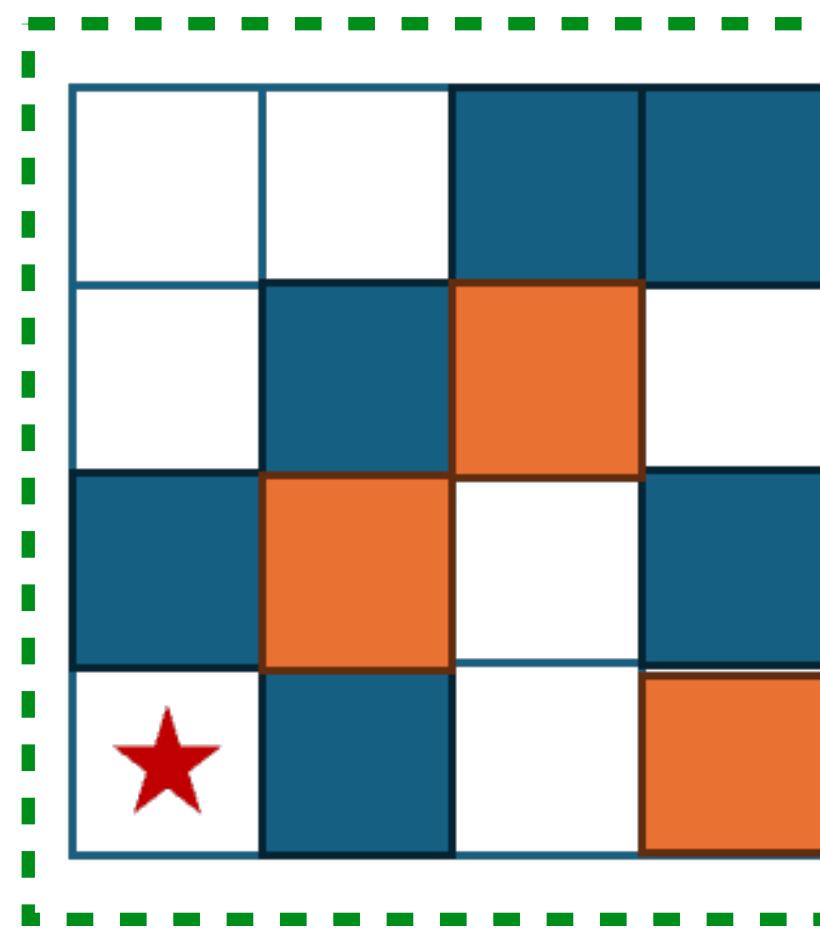
# Weighted Conformal Prediction

## Leave-one-out Likelihood Weighting

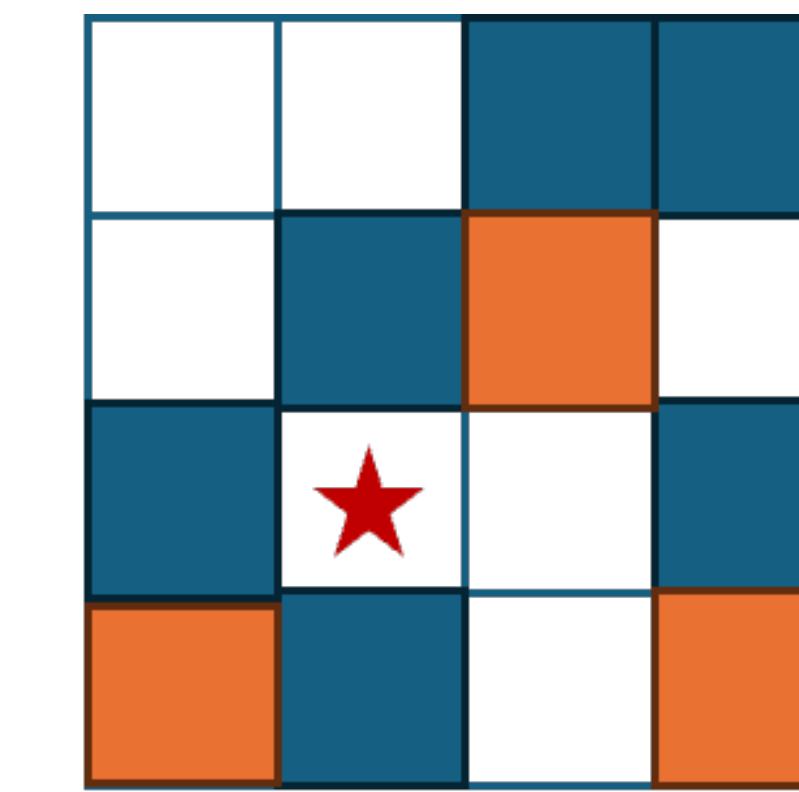
■ : training   ■ : calibration   □ : missing   ★ : testing



Observation



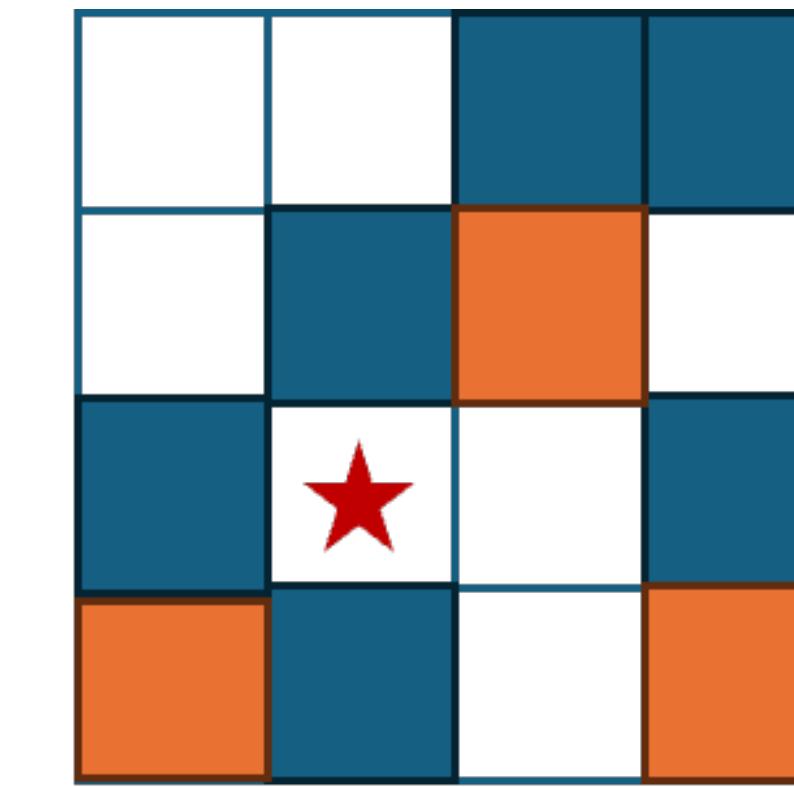
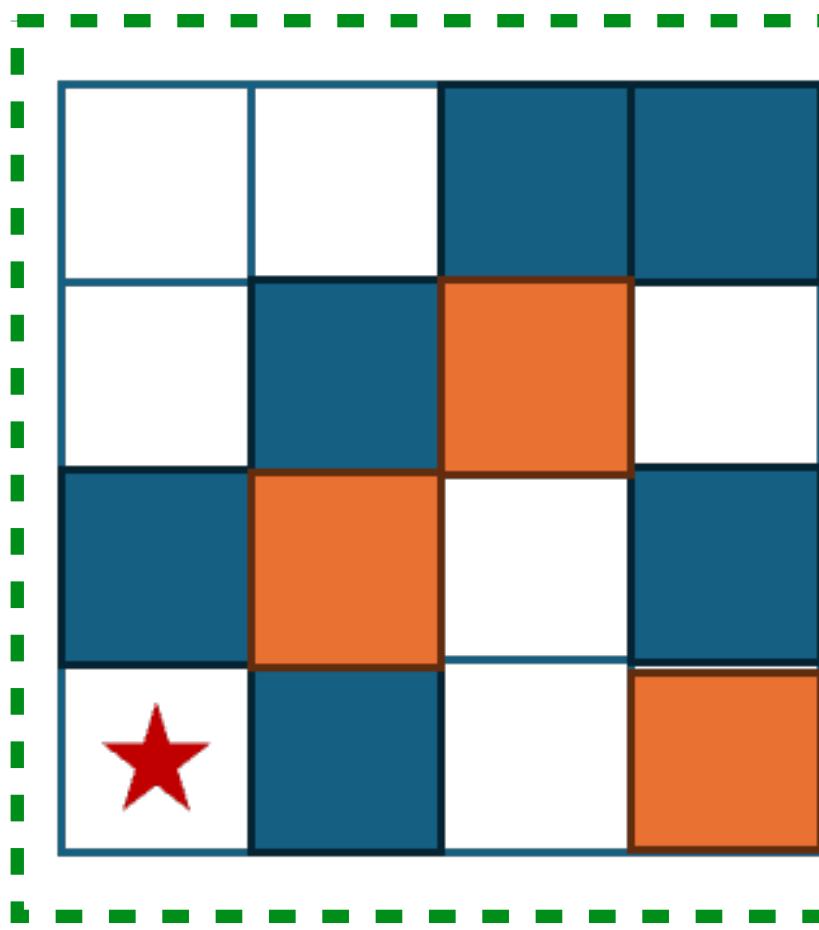
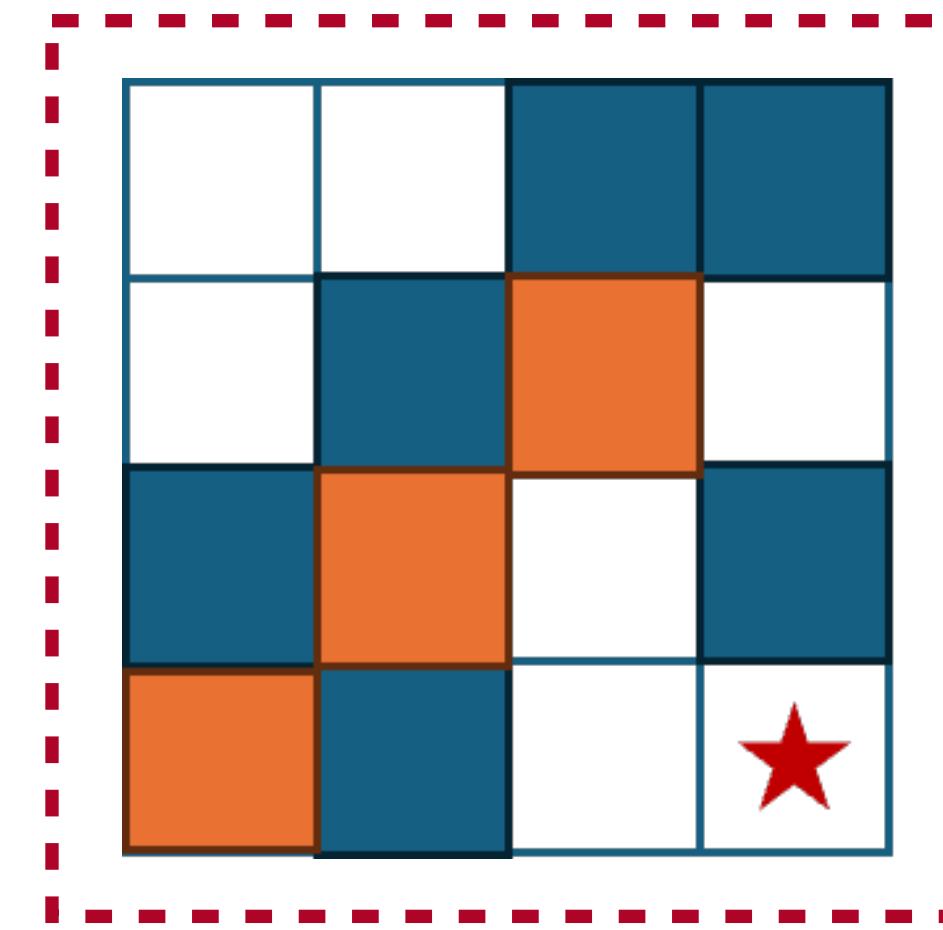
Hypothetical scenarios



# Weighted Conformal Prediction

## Leave-one-out Likelihood Weighting

■ : training   ■ : calibration   □ : missing   ★ : testing



probability

$p_1$

$p_2$

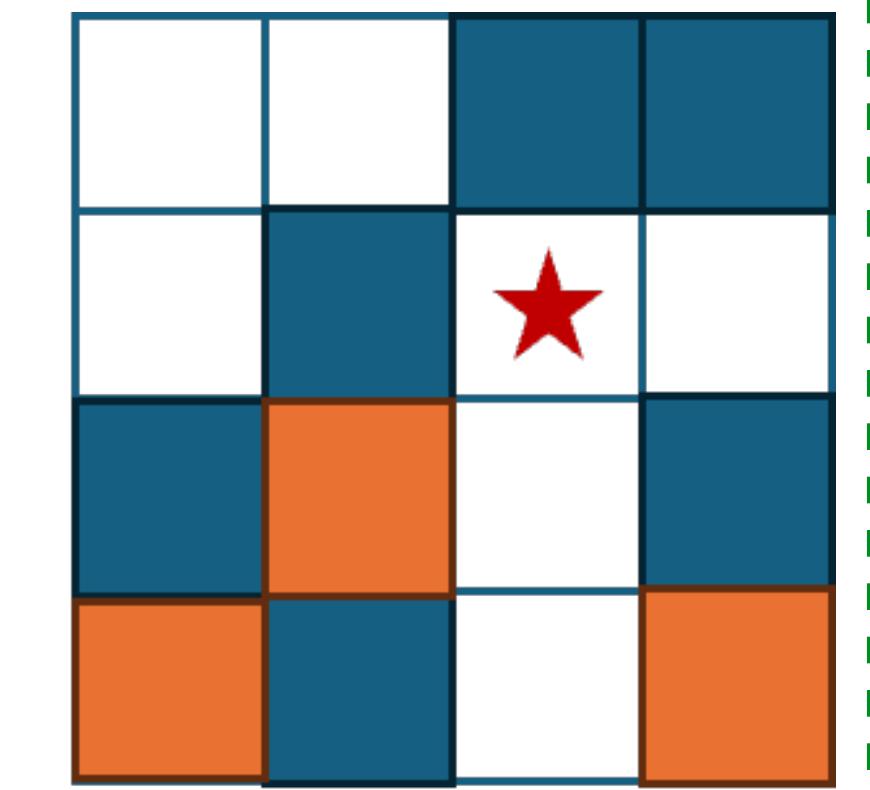
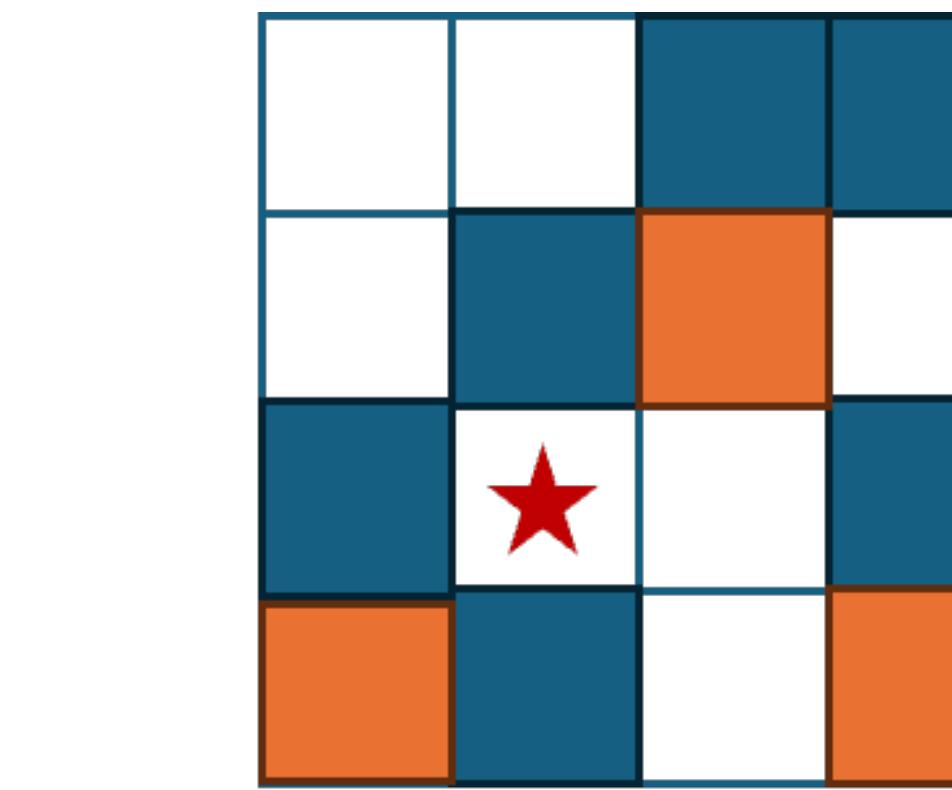
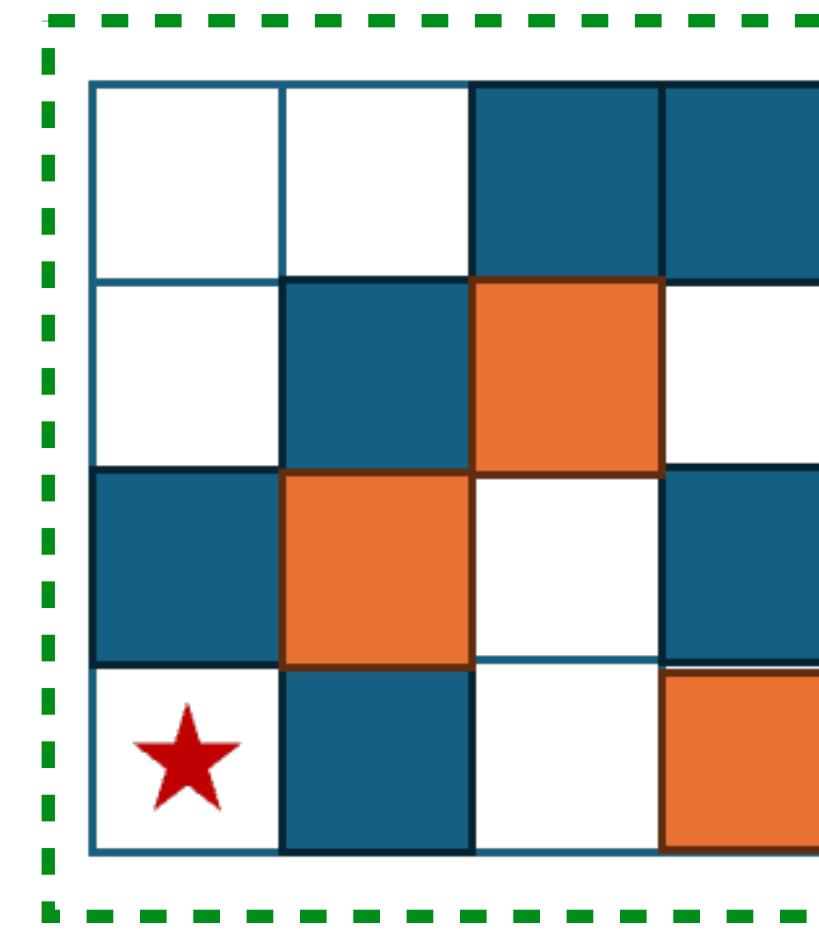
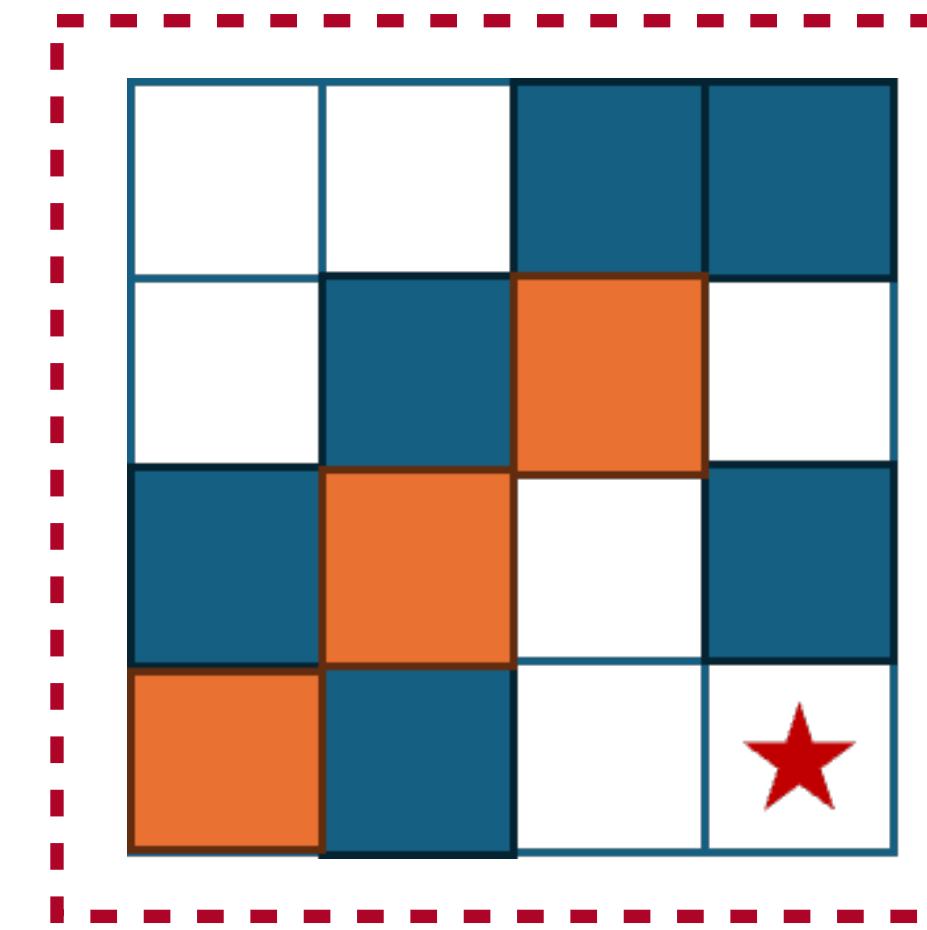
$p_3$

$p_4$

# Weighted Conformal Prediction

## Leave-one-out Likelihood Weighting

■ : training   ■ : calibration   □ : missing   ★ : testing



Observation

Hypothetical scenarios

probability

$p_1$

$p_2$

$p_3$

$p_4$

$$\text{weight of } \star \frac{p_1}{p_1 + p_2 + p_3 + p_4}$$

$$\frac{p_2}{p_1 + p_2 + p_3 + p_4}$$

$$\frac{p_3}{p_1 + p_2 + p_3 + p_4}$$

$$\frac{p_4}{p_1 + p_2 + p_3 + p_4}$$

# Weighted Conformal Prediction

## Leave-one-out Likelihood Weighting

### Theorem 1

Let  $\omega_k \propto P([\mathcal{W}]_s = 1, \forall s \in \Omega_{tr} \cup \{s_1, \dots, s_{n_{cal}+1}\} \setminus \{s_k\}, [\mathcal{W}]_s = -1, \text{o.w.})$ , then

$$\widehat{C}_{1-\alpha, s^*} = \left\{ x \in \mathbb{R} \mid \mathcal{S}(x, [\widehat{\mathcal{X}}]_{s^*}) \leq \widehat{q}_{1-\alpha} \left( \sum_{s \in \Omega_{cal}} \omega_s \cdot \delta_{\mathcal{S}([\mathcal{X}]_s, [\widehat{\mathcal{X}}]_s)} + \omega_{s^*} \cdot \delta_{+\infty} \right) \right\}$$

is still valid, i.e.  $P([\mathcal{X}]_{s^*} \in \widehat{C}_{1-\alpha, s^*}) \geq 1 - \alpha$ .

# Weighted Conformal Prediction

## Leave-one-out Likelihood Weighting

### Theorem 1

Let  $\omega_k \propto P([\mathcal{W}]_s = 1, \forall s \in \Omega_{tr} \cup \{s_1, \dots, s_{n_{cal}+1}\} \setminus \{s_k\}, [\mathcal{W}]_s = -1, \text{o.w.})$ , then

$$\widehat{C}_{1-\alpha, s^*} = \left\{ x \in \mathbb{R} \mid \mathcal{S}(x, [\widehat{\mathcal{X}}]_{s^*}) \leq \widehat{q}_{1-\alpha} \left( \sum_{s \in \Omega_{cal}} \omega_s \cdot \delta_{\mathcal{S}([\mathcal{X}]_s, [\widehat{\mathcal{X}}]_s)} + \omega_{s^*} \cdot \delta_{+\infty} \right) \right\}$$

is still valid, i.e.  $P([\mathcal{X}]_{s^*} \in \widehat{C}_{1-\alpha, s^*}) \geq 1 - \alpha$ .

Now we need to estimate the likelihood  $p(\mathcal{W})$  with a single sample.

# Missing Propensity Model

## Motivation

# Missing Propensity Model

## Motivation

- We want to build a probabilistic model for a binary tensor  $\mathcal{W}$  that can:

# Missing Propensity Model

## Motivation

- We want to build a probabilistic model for a binary tensor  $\mathcal{W}$  that can:
  - model independent missingness with heterogeneous probabilities

# Missing Propensity Model

## Motivation

- We want to build a probabilistic model for a binary tensor  $\mathcal{W}$  that can:
  - model independent missingness with heterogeneous probabilities

$$[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}([\mathcal{P}]_s), \quad \mathcal{P} \in [0,1]^{d_1 \times \dots \times d_K}$$

# Missing Propensity Model

## Motivation

- We want to build a probabilistic model for a binary tensor  $\mathcal{W}$  that can:
  - model independent missingness with heterogeneous probabilities

$$[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}([\mathcal{P}]_s), \quad \mathcal{P} \in [0,1]^{d_1 \times \dots \times d_K}$$

- model locally-dependent missingness

# Missing Propensity Model

## Motivation

- We want to build a probabilistic model for a binary tensor  $\mathcal{W}$  that can:
  - model independent missingness with heterogeneous probabilities

$$[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}([\mathcal{P}]_s), \quad \mathcal{P} \in [0,1]^{d_1 \times \dots \times d_K}$$

- model locally-dependent missingness

$$[\mathcal{W}]_{s_1} \perp [\mathcal{W}]_{s_2} \mid [\mathcal{W}]_j, j \in N(s_1)$$

# Missing Propensity Model

## Motivation

- We want to build a probabilistic model for a binary tensor  $\mathcal{W}$  that can:
  - model independent missingness with heterogeneous probabilities

$$[\mathcal{W}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}([\mathcal{P}]_s), \quad \mathcal{P} \in [0,1]^{d_1 \times \dots \times d_K}$$

- model locally-dependent missingness

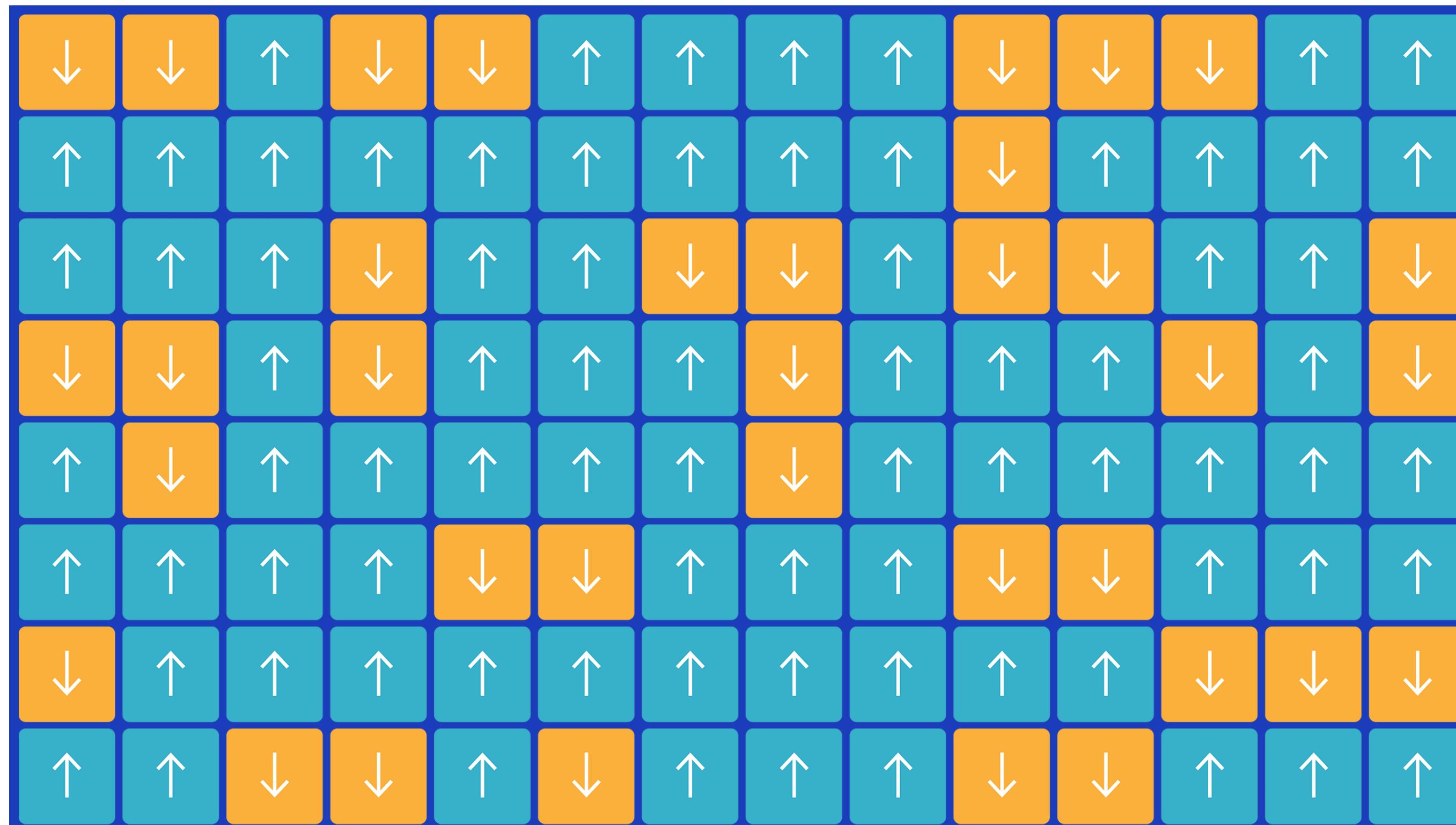
$$[\mathcal{W}]_{s_1} \perp [\mathcal{W}]_{s_2} \mid [\mathcal{W}]_j, j \in N(s_1)$$

where  $N(s_1) = \{s \mid s \sim s_1\}$ , and  $s \sim s_1$  means  $s, s_1$  are neighboring tensor entries.

# Missing Propensity Model

## Tensor Ising Model

- We adopt the idea from the Ising model [Cipra (1987)], which models the probability distribution of atomic spins on a lattice grid.

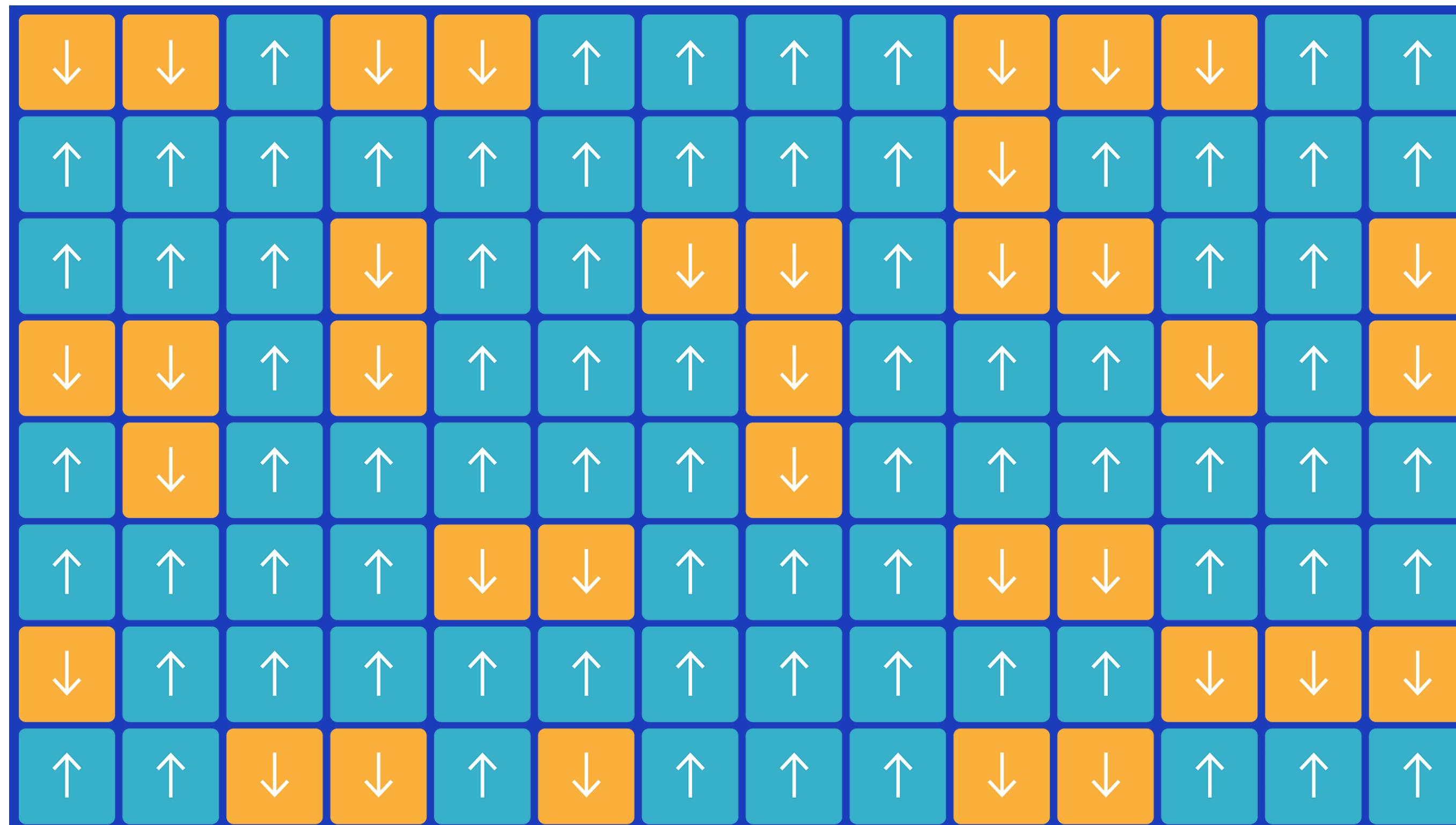


[source]

# Missing Propensity Model

## Tensor Ising Model

- We adopt the idea from the Ising model [Cipra (1987)], which models the probability distribution of atomic spins on a lattice grid.



[source]

# Missing Propensity Model

## Tensor Ising Model

- We model  $p(\mathcal{W})$  via a Boltzmann distribution:

$$p(\mathcal{W}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g_{s_1 s_2} \cdot \overbrace{[\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2}}^{\text{"co-missingness"}} + \sum_{s_1} h_{s_1} \cdot [\mathcal{W}]_{s_1} \right]$$

# Missing Propensity Model

## Tensor Ising Model

- We model  $p(\mathcal{W})$  via a Boltzmann distribution:

$$p(\mathcal{W}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g_{s_1 s_2} \cdot \overbrace{[\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2}}^{\text{"co-missingness"}} + \sum_{s_1} h_{s_1} \cdot [\mathcal{W}]_{s_1} \right]$$

↓  
interaction strength

# Missing Propensity Model

## Tensor Ising Model

- We model  $p(\mathcal{W})$  via a Boltzmann distribution:

$$p(\mathcal{W}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g_{s_1 s_2} \cdot \overbrace{[\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2}}^{\text{"co-missingness"}} + \sum_{s_1} h_{s_1} \cdot [\mathcal{W}]_{s_1} \right]$$

$\downarrow \qquad \qquad \qquad \downarrow$

interaction strength      individual missing propensity

# Missing Propensity Model

## Parameter Estimation Problem of Tensor Ising Model

# Missing Propensity Model

## Parameter Estimation Problem of Tensor Ising Model

- Given  $\mathcal{W}$ , we want to estimate  $g_{s_1, s_2}$  for arbitrary  $s_1 \sim s_2$  and  $h_{s_1}$  for any  $s_1$ .

# Missing Propensity Model

## Parameter Estimation Problem of Tensor Ising Model

- Given  $\mathcal{W}$ , we want to estimate  $g_{s_1, s_2}$  for arbitrary  $s_1 \sim s_2$  and  $h_{s_1}$  for any  $s_1$ .
- Existing works [Ravikumar, et al. (2010), Barber and Drton (2015), Liu et al. (2024)] on learning the graph structure of Ising model deal with multiple samples of 1-D or 2-D Ising model.

# Missing Propensity Model

## Parameter Estimation Problem of Tensor Ising Model

- Given  $\mathcal{W}$ , we want to estimate  $g_{s_1, s_2}$  for arbitrary  $s_1 \sim s_2$  and  $h_{s_1}$  for any  $s_1$ .
- Existing works [Ravikumar, et al. (2010), Barber and Drton (2015), Liu et al. (2024)] on learning the graph structure of Ising model deal with multiple samples of 1-D or 2-D Ising model.
- In our work, we have a single sample but an order- $K$  Ising model.

# **Missing Propensity Model**

## **Latent Space Tensor Ising Model**

# Missing Propensity Model

## Latent Space Tensor Ising Model

- We assume each tensor entry  $s$  has a latent feature  $[\mathcal{B}]_s \in \mathbb{R}$ , and:

# Missing Propensity Model

## Latent Space Tensor Ising Model

- We assume each tensor entry  $s$  has a latent feature  $[\mathcal{B}]_s \in \mathbb{R}$ , and:

$$p(\mathcal{W} | \mathcal{B}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g([\mathcal{B}]_{s_1}, [\mathcal{B}]_{s_2}) \cdot [\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2} + \sum_{s_1} h([\mathcal{B}]_{s_1}) \cdot [\mathcal{W}]_{s_1} \right],$$

# Missing Propensity Model

## Latent Space Tensor Ising Model

- We assume each tensor entry  $s$  has a latent feature  $[\mathcal{B}]_s \in \mathbb{R}$ , and:

$$p(\mathcal{W} | \mathcal{B}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g([\mathcal{B}]_{s_1}, [\mathcal{B}]_{s_2}) \cdot [\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2} + \sum_{s_1} h([\mathcal{B}]_{s_1}) \cdot [\mathcal{W}]_{s_1} \right],$$

where:

# Missing Propensity Model

## Latent Space Tensor Ising Model

- We assume each tensor entry  $s$  has a latent feature  $[\mathcal{B}]_s \in \mathbb{R}$ , and:

$$p(\mathcal{W} | \mathcal{B}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g([\mathcal{B}]_{s_1}, [\mathcal{B}]_{s_2}) \cdot [\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2} + \sum_{s_1} h([\mathcal{B}]_{s_1}) \cdot [\mathcal{W}]_{s_1} \right],$$

where:

- $g(\cdot, \cdot)$  is symmetric;

# Missing Propensity Model

## Latent Space Tensor Ising Model

- We assume each tensor entry  $s$  has a latent feature  $[\mathcal{B}]_s \in \mathbb{R}$ , and:

$$p(\mathcal{W} | \mathcal{B}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g([\mathcal{B}]_{s_1}, [\mathcal{B}]_{s_2}) \cdot [\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2} + \sum_{s_1} h([\mathcal{B}]_{s_1}) \cdot [\mathcal{W}]_{s_1} \right],$$

where:

- $g(\cdot, \cdot)$  is symmetric;
- $h(\cdot)$  is twice-continuously differentiable.

# Example I: Bernoulli Model

# Example I: Bernoulli Model

- If  $g(\cdot, \cdot) = 0$ , then:

# Example I: Bernoulli Model

- If  $g(\cdot, \cdot) = 0$ , then:

$$[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}([\mathcal{P}]_s), \quad [\mathcal{P}]_s = \frac{\exp[2h([\mathcal{B}]_s)]}{1 + \exp[2h([\mathcal{B}]_s)]}$$

# Example I: Bernoulli Model

- If  $g(\cdot, \cdot) = 0$ , then:

$$[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}([\mathcal{P}]_s), \quad [\mathcal{P}]_s = \frac{\exp[2h([\mathcal{B}]_s)]}{1 + \exp[2h([\mathcal{B}]_s)]}$$

- The conformal weight satisfies:

# Example I: Bernoulli Model

- If  $g(\cdot, \cdot) = 0$ , then:

$$[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}([\mathcal{P}]_s), \quad [\mathcal{P}]_s = \frac{\exp[2h([\mathcal{B}]_s)]}{1 + \exp[2h([\mathcal{B}]_s)]}$$

- The conformal weight satisfies:

$$\omega_s \propto \frac{1 - [\mathcal{P}]_s}{[\mathcal{P}]_s} = \exp[-2h([\mathcal{B}]_s)]$$

# Example II: Interaction-only Ising Model

# Example II: Interaction-only Ising Model

- If  $h(\cdot) = 0$ ,

# Example II: Interaction-only Ising Model

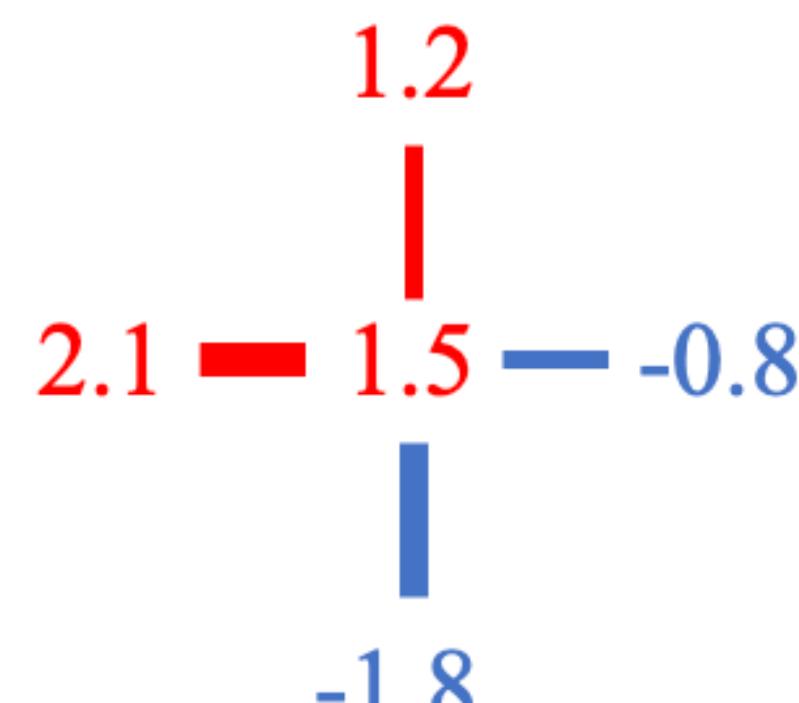
- If  $h(\cdot) = 0$ ,

$$\tilde{p}_s = P([\mathcal{W}]_s = 1 | [\mathcal{W}]_{-s}, \mathcal{B}) = \frac{\exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}{1 + \exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}$$

# Example II: Interaction-only Ising Model

- If  $h(\cdot) = 0$ ,

$$\tilde{p}_s = P([\mathcal{W}]_s = 1 | [\mathcal{W}]_{-s}, \mathcal{B}) = \frac{\exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}{1 + \exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}$$

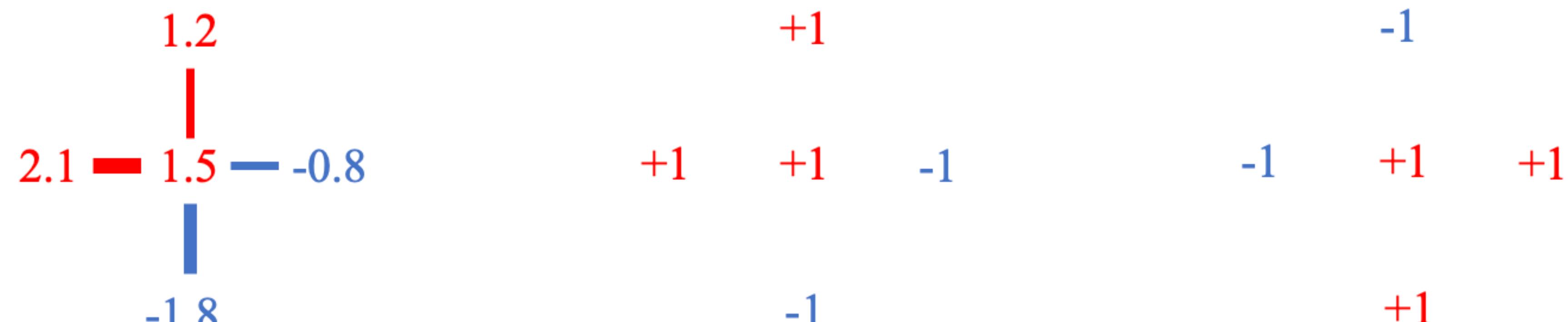


Example  $\mathcal{B}$

# Example II: Interaction-only Ising Model

- If  $h(\cdot) = 0$ ,

$$\tilde{p}_s = P([\mathcal{W}]_s = 1 | [\mathcal{W}]_{-s}, \mathcal{B}) = \frac{\exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}{1 + \exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}$$

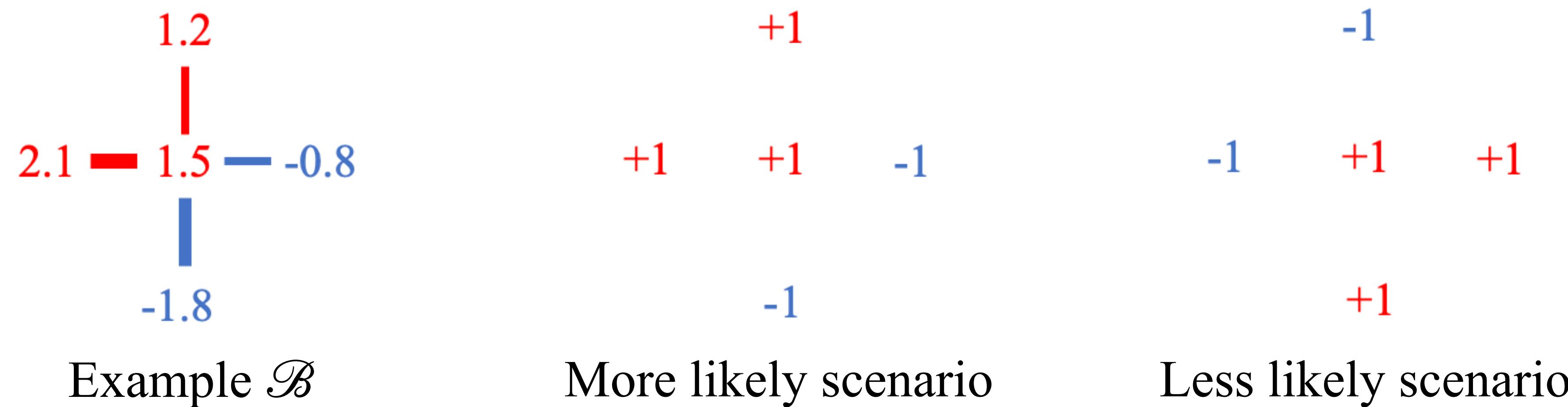


Example  $\mathcal{B}$

# Example II: Interaction-only Ising Model

- If  $h(\cdot) = 0$ ,

$$\tilde{p}_s = P([\mathcal{W}]_s = 1 | [\mathcal{W}]_{-s}, \mathcal{B}) = \frac{\exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}{1 + \exp \left[ 2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]}$$



# Example II: Interaction-only Ising Model

- The conformal weight (almost always) satisfies:

$$\omega_s \propto \frac{1 - \tilde{p}_s}{\tilde{p}_s} = \exp \left[ -2 \sum_{j \in N(s)} g([\mathcal{B}]_s, [\mathcal{B}]_j) \cdot [\mathcal{W}]_j \right]$$

Exceptions are samples in  $\Omega_{cal} \cap N(s_{n+1})$ , where  $[\mathcal{W}]_{s_{n+1}} = -1$ , but the leave-one-out procedure requires  $[\mathcal{W}]_{s_{n+1}} = 1$ .

# **Estimation Method**

## **Maximum Pseudo-Likelihood Estimation (MPLE)**

# Estimation Method

## Maximum Pseudo-Likelihood Estimation (MPLE)

- For theoretical reason [Barber et al. (2023)], we can only fit the Ising model with  $\mathcal{W}_{\Omega_{tr}}$ , defined as:

# Estimation Method

## Maximum Pseudo-Likelihood Estimation (MPLE)

- For theoretical reason [Barber et al. (2023)], we can only fit the Ising model with  $\mathcal{W}_{\Omega_{tr}}$ , defined as:

$$[\mathcal{W}_{\Omega_{tr}}]_s = 1, \quad \text{if and only if } s \in \Omega_{tr}$$

# Estimation Method

## Maximum Pseudo-Likelihood Estimation (MPLE)

- For theoretical reason [Barber et al. (2023)], we can only fit the Ising model with  $\mathcal{W}_{\Omega_{tr}}$ , defined as:

$$[\mathcal{W}_{\Omega_{tr}}]_s = 1, \quad \text{if and only if } s \in \Omega_{tr}$$

- The pseudo-likelihood of  $\mathcal{W}_{\Omega_{tr}}$  posits that  $[\mathcal{W}_{\Omega_{tr}}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(q \cdot \tilde{p}_s)$ , thus the MPLE is:

# Estimation Method

## Maximum Pseudo-Likelihood Estimation (MPLE)

- For theoretical reason [Barber et al. (2023)], we can only fit the Ising model with  $\mathcal{W}_{\Omega_{tr}}$ , defined as:

$$[\mathcal{W}_{\Omega_{tr}}]_s = 1, \quad \text{if and only if } s \in \Omega_{tr}$$

- The pseudo-likelihood of  $\mathcal{W}_{\Omega_{tr}}$  posits that  $[\mathcal{W}_{\Omega_{tr}}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(q \cdot \tilde{p}_s)$ , thus the MPLE is:

$$\widehat{\mathcal{B}} = \operatorname{argmin}_{\mathcal{B} \in \mathbb{F}} - \sum_{s: [\mathcal{W}_{\Omega_{tr}}]_s = 1} \log(q\tilde{p}_s) - \sum_{s: [\mathcal{W}_{\Omega_{tr}}]_s = -1} \log(1 - q\tilde{p}_s)$$

# Estimation Method

## Maximum Pseudo-Likelihood Estimation (MPLE)

- For theoretical reason [Barber et al. (2023)], we can only fit the Ising model with  $\mathcal{W}_{\Omega_{tr}}$ , defined as:

$$[\mathcal{W}_{\Omega_{tr}}]_s = 1, \quad \text{if and only if } s \in \Omega_{tr}$$

- The pseudo-likelihood of  $\mathcal{W}_{\Omega_{tr}}$  posits that  $[\mathcal{W}_{\Omega_{tr}}]_s \stackrel{\text{ind.}}{\sim} \text{Bern}(q \cdot \tilde{p}_s)$ , thus the MPLE is:

$$\widehat{\mathcal{B}} = \operatorname{argmin}_{\mathcal{B} \in \mathbb{F}} - \sum_{s: [\mathcal{W}_{\Omega_{tr}}]_s = 1} \log(q\tilde{p}_s) - \sum_{s: [\mathcal{W}_{\Omega_{tr}}]_s = -1} \log(1 - q\tilde{p}_s)$$

where  $q \in (0,1)$  is the probability of assigning data to the training set, and  $\mathbb{F}$  is a feasible set of  $\mathcal{B}$ .

# **Estimation Method**

## **Low-rank MPLE**

# Estimation Method

## Low-rank MPLE

- Given a single sample  $\mathcal{W}_{\Omega_{tr}}$ , we need to impose a parsimonious structure over  $\mathcal{B}$  for estimation.

# Estimation Method

## Low-rank MPLE

- Given a single sample  $\mathcal{W}_{\Omega_{tr}}$ , we need to impose a parsimonious structure over  $\mathcal{B}$  for estimation.
- We set  $\mathbb{F}$  to be the set of low tensor-train (TT) rank tensors [Oseledets (2011)]:

# Estimation Method

## Low-rank MPLE

- Given a single sample  $\mathcal{W}_{\Omega_{tr}}$ , we need to impose a parsimonious structure over  $\mathcal{B}$  for estimation.
- We set  $\mathbb{F}$  to be the set of low tensor-train (TT) rank tensors [Oseledets (2011)]:

$$\mathbb{F} = \{\mathcal{B} \in \mathbb{R}^{d_1 \times \dots \times d_K} \mid [\mathcal{B}]_{i_1 \dots i_K} = \prod_{k=1}^K [\mathcal{T}^k]_{:i_k:}, \mathcal{T}^k \in \mathbb{R}^{r_{k-1} \times d_k \times r_k}\},$$

# Estimation Method

## Low-rank MPLE

- Given a single sample  $\mathcal{W}_{\Omega_{tr}}$ , we need to impose a parsimonious structure over  $\mathcal{B}$  for estimation.
- We set  $\mathbb{F}$  to be the set of low tensor-train (TT) rank tensors [Oseledets (2011)]:

$$\mathbb{F} = \{\mathcal{B} \in \mathbb{R}^{d_1 \times \dots \times d_K} \mid [\mathcal{B}]_{i_1 \dots i_K} = \prod_{k=1}^K [\mathcal{T}^k]_{:i_k:}, \mathcal{T}^k \in \mathbb{R}^{r_{k-1} \times d_k \times r_k}\},$$

with  $r_0 = r_K = 1$ , and  $\mathbf{r} = (r_1, \dots, r_{K-1})$  being the TT-rank.

# **Estimation Algorithm**

## **Riemannian Gradient Descent (RGrad)**

# Estimation Algorithm

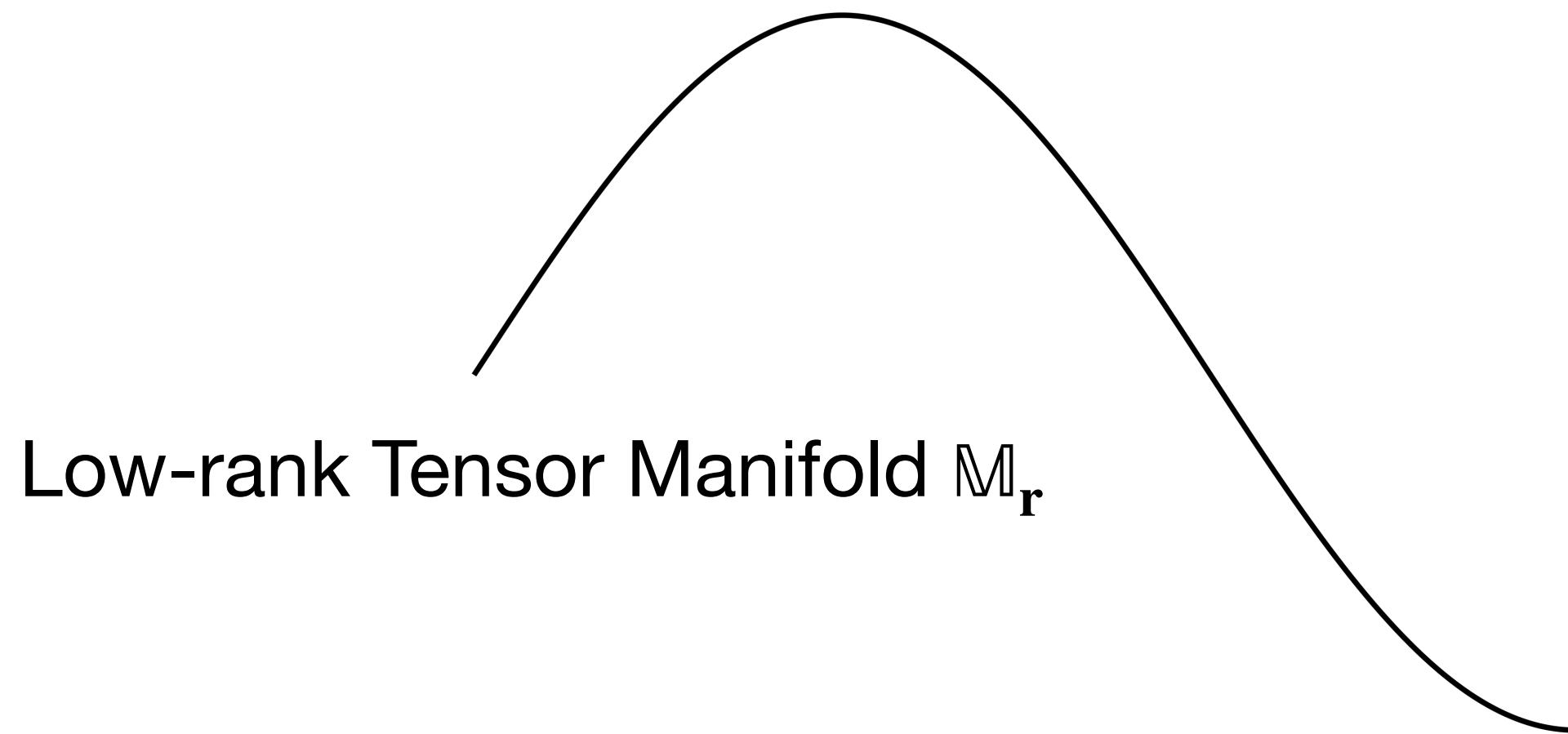
## Riemannian Gradient Descent (RGrad)

- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]

# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

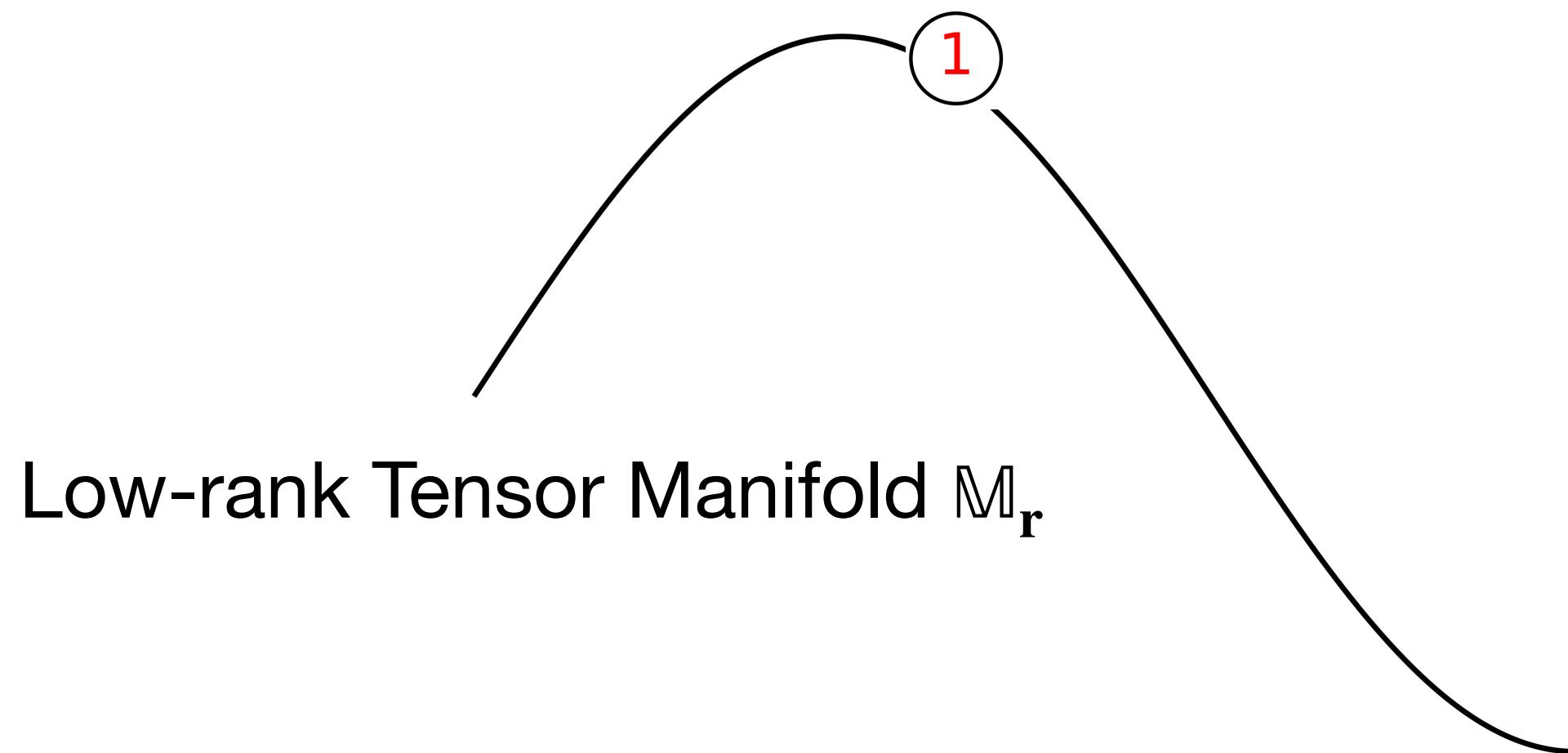
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

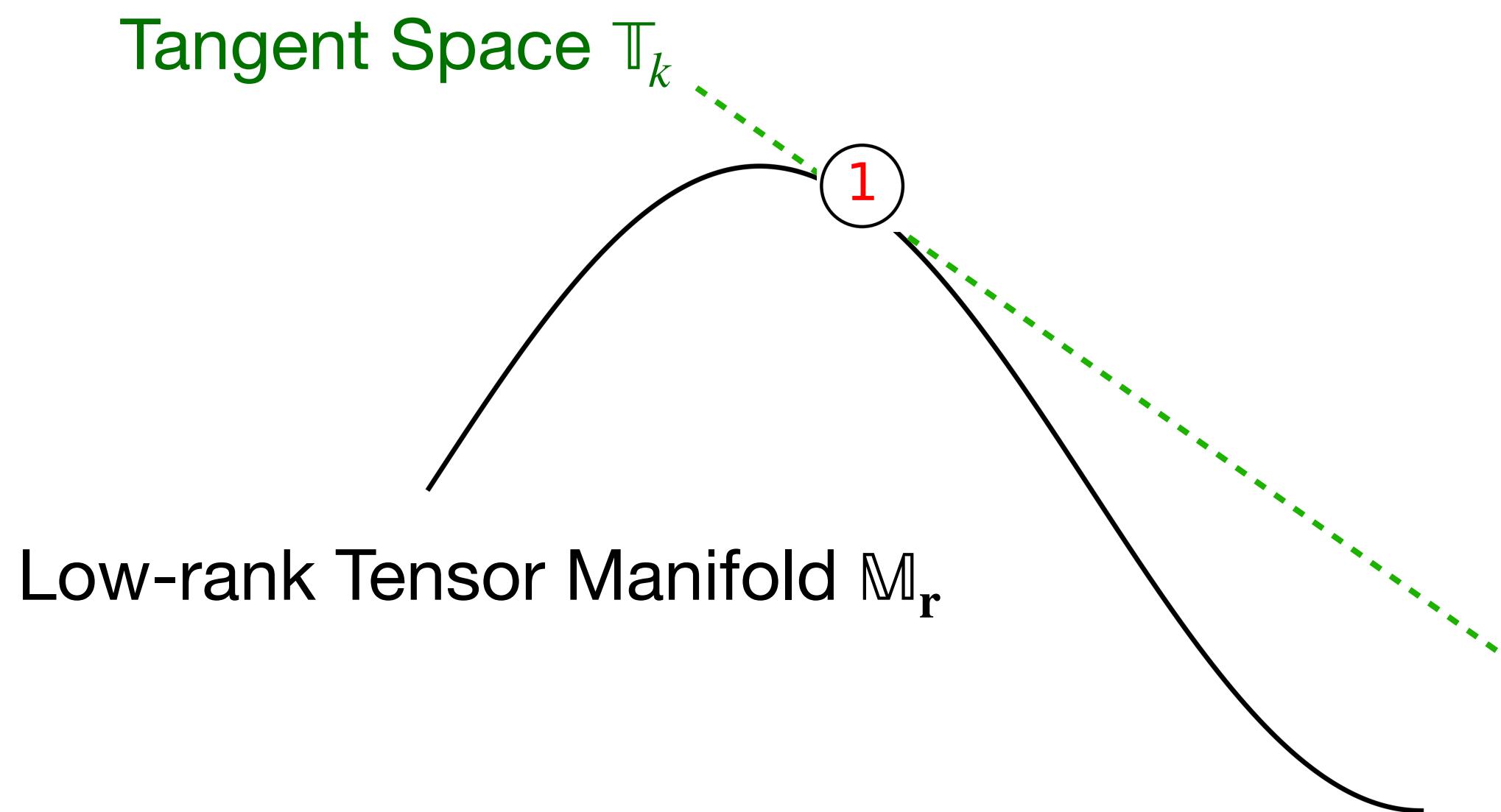
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

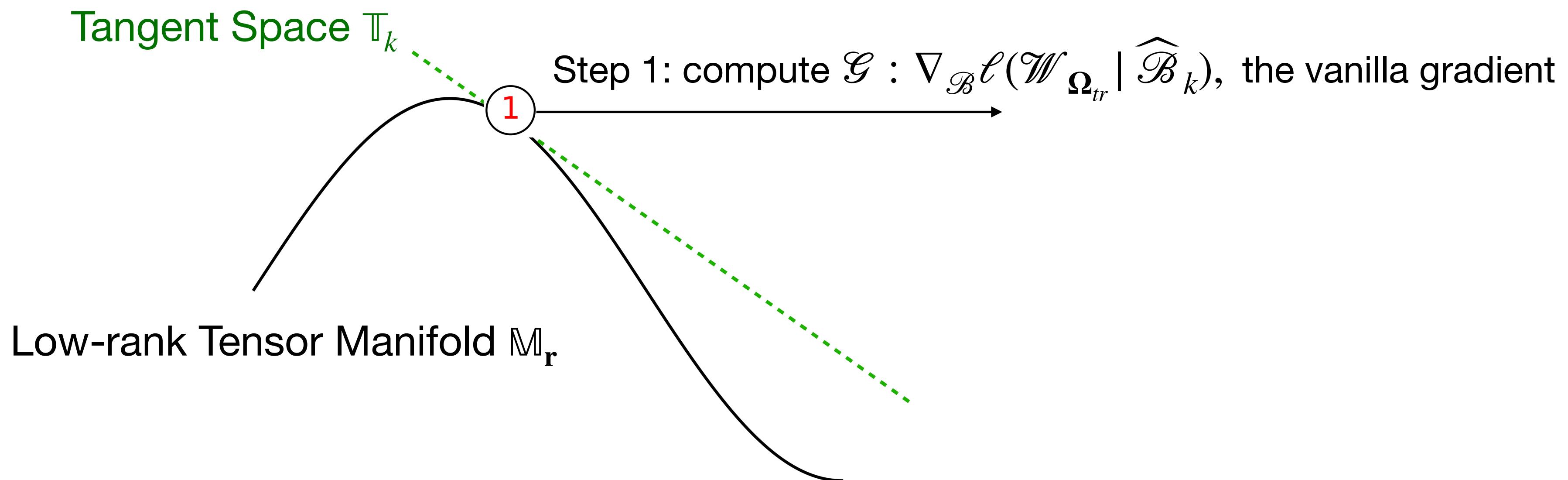
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

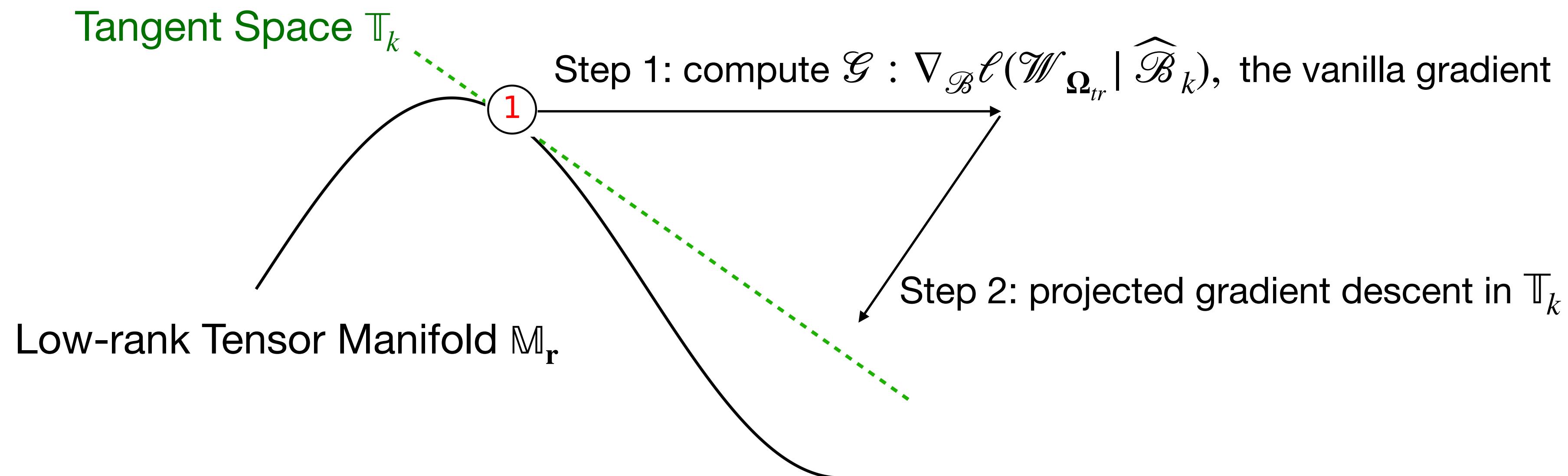
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

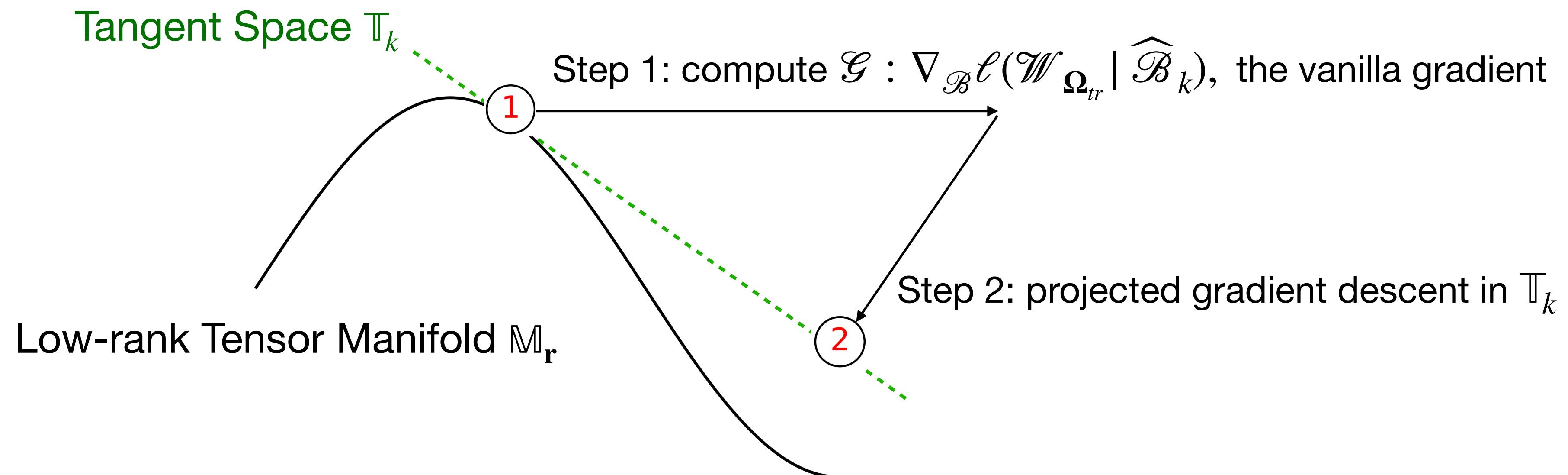
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

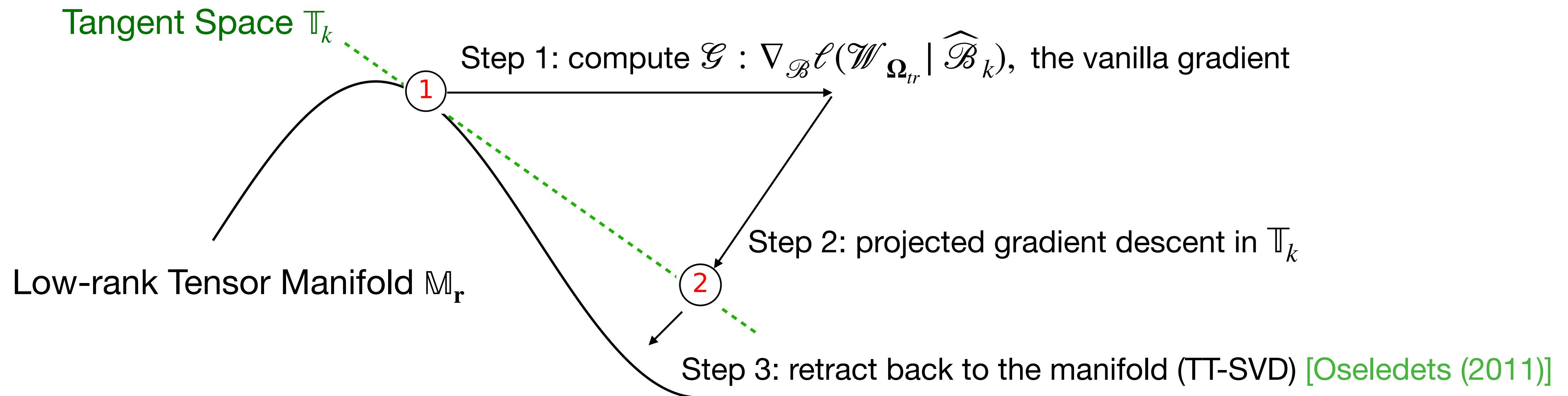
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



# Estimation Algorithm

## Riemannian Gradient Descent (RGrad)

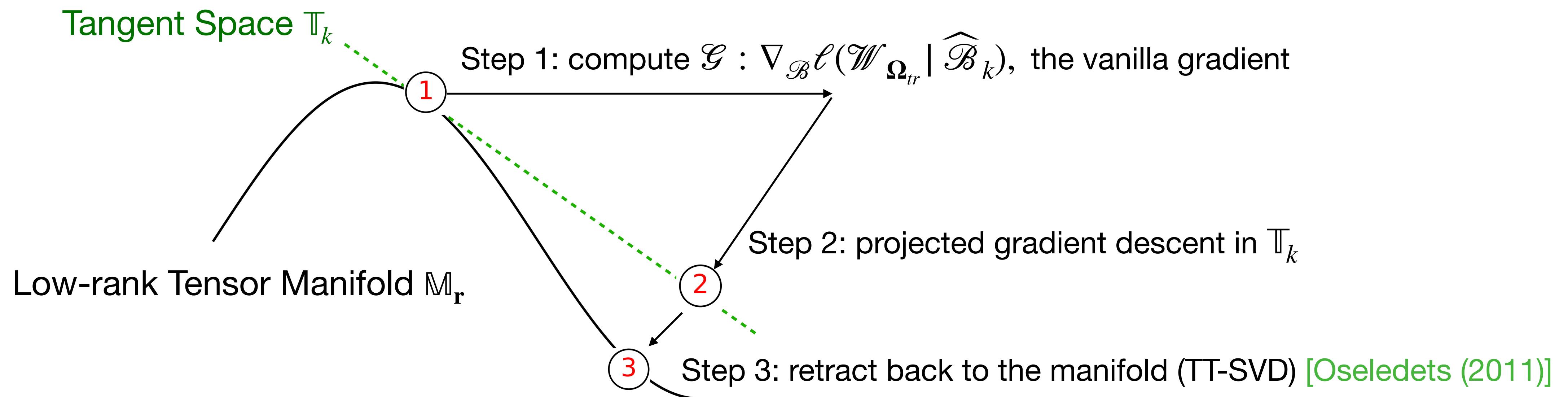
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]



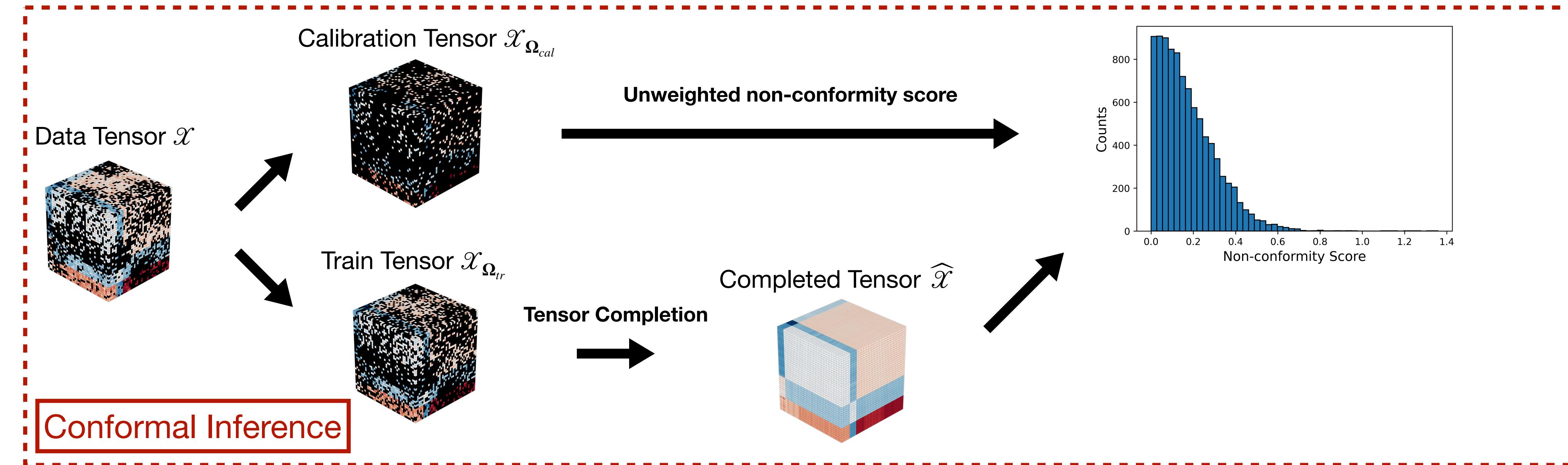
# Estimation Algorithm

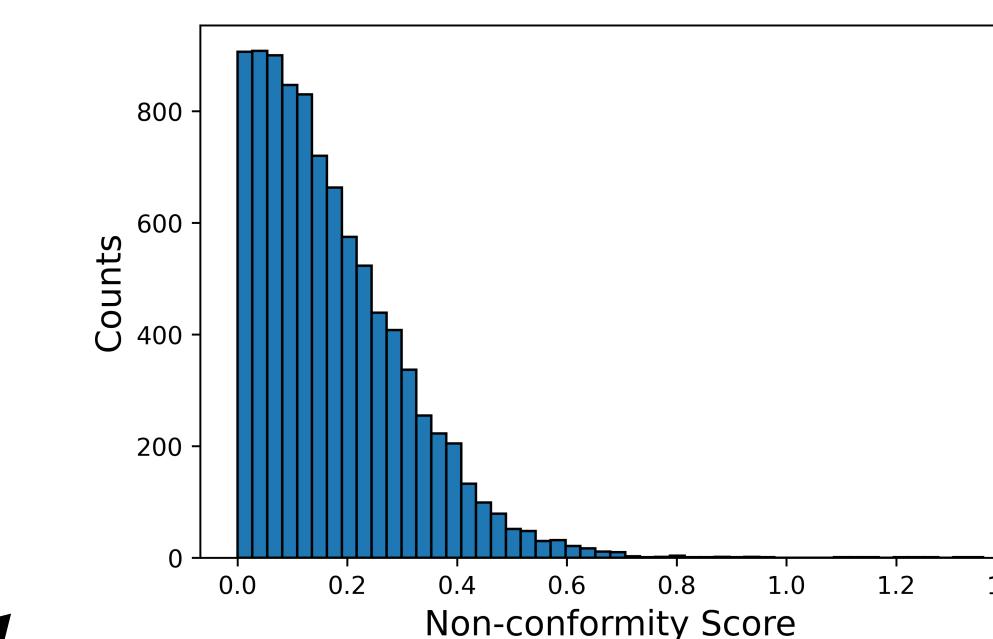
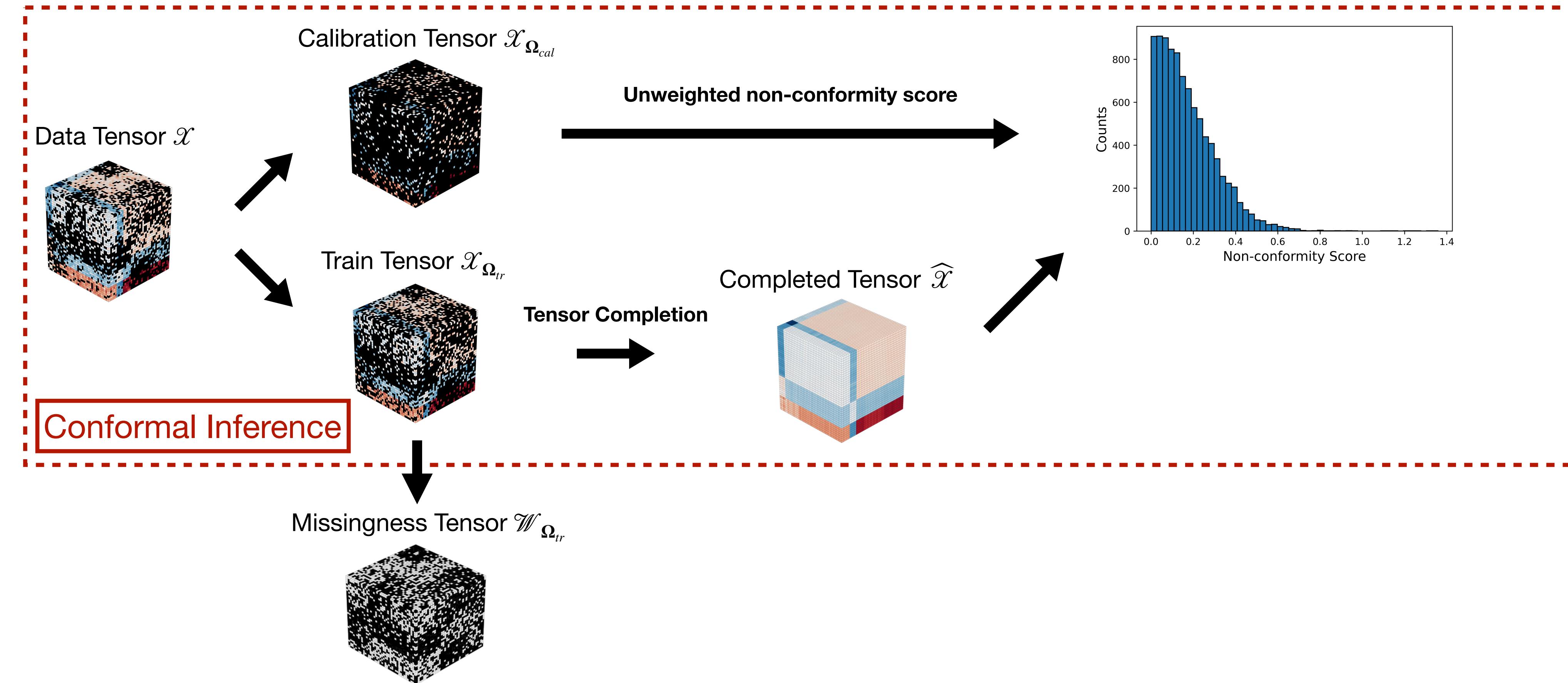
## Riemannian Gradient Descent (RGrad)

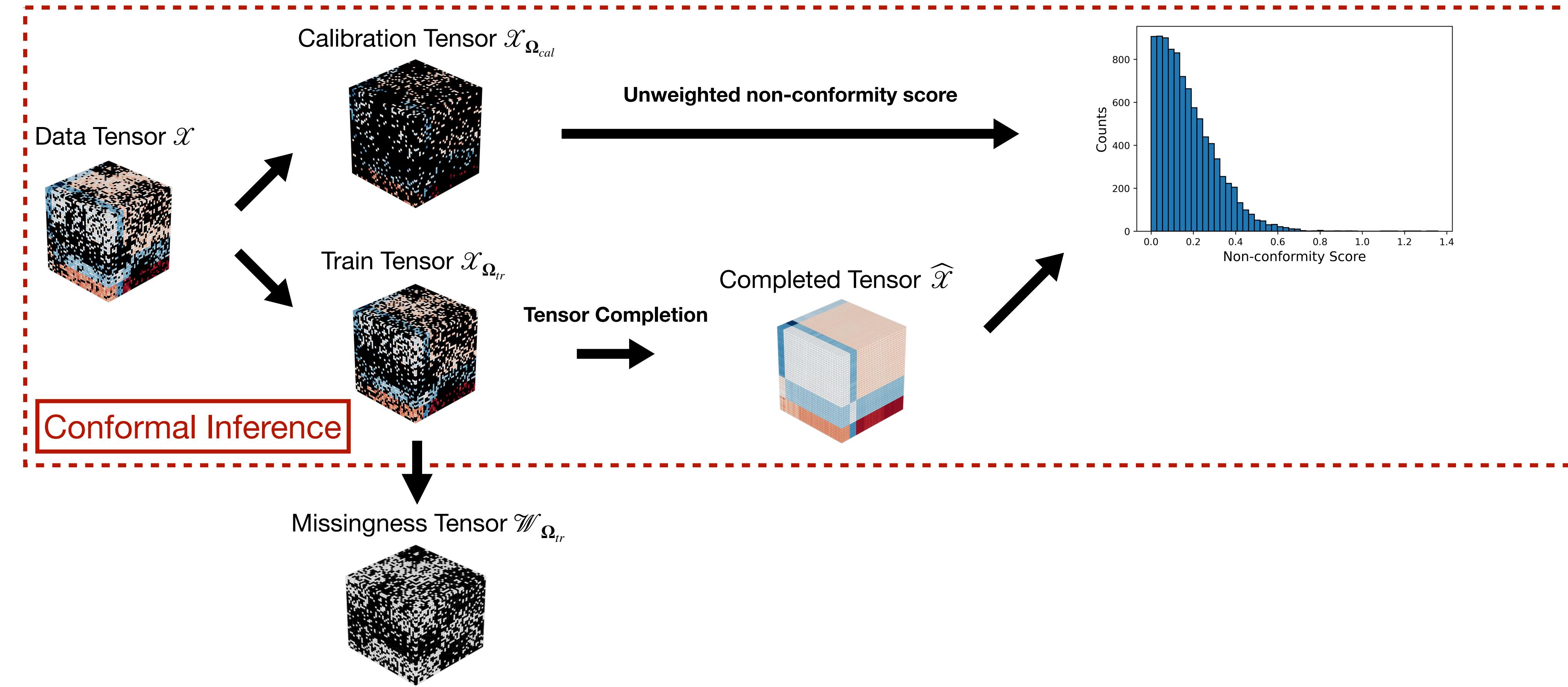
- Low tensor-train rank tensors lie on a smooth manifold  $\mathbb{M}_r$  [Holtz et al. (2012)]







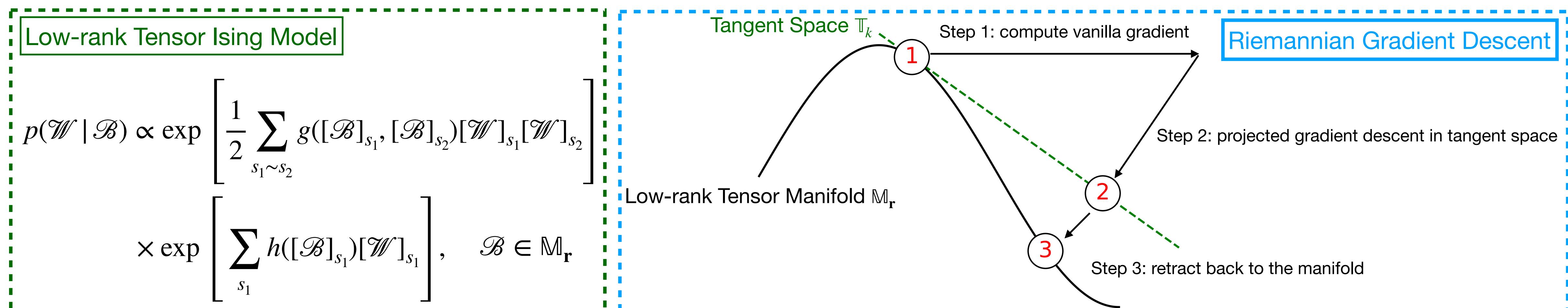
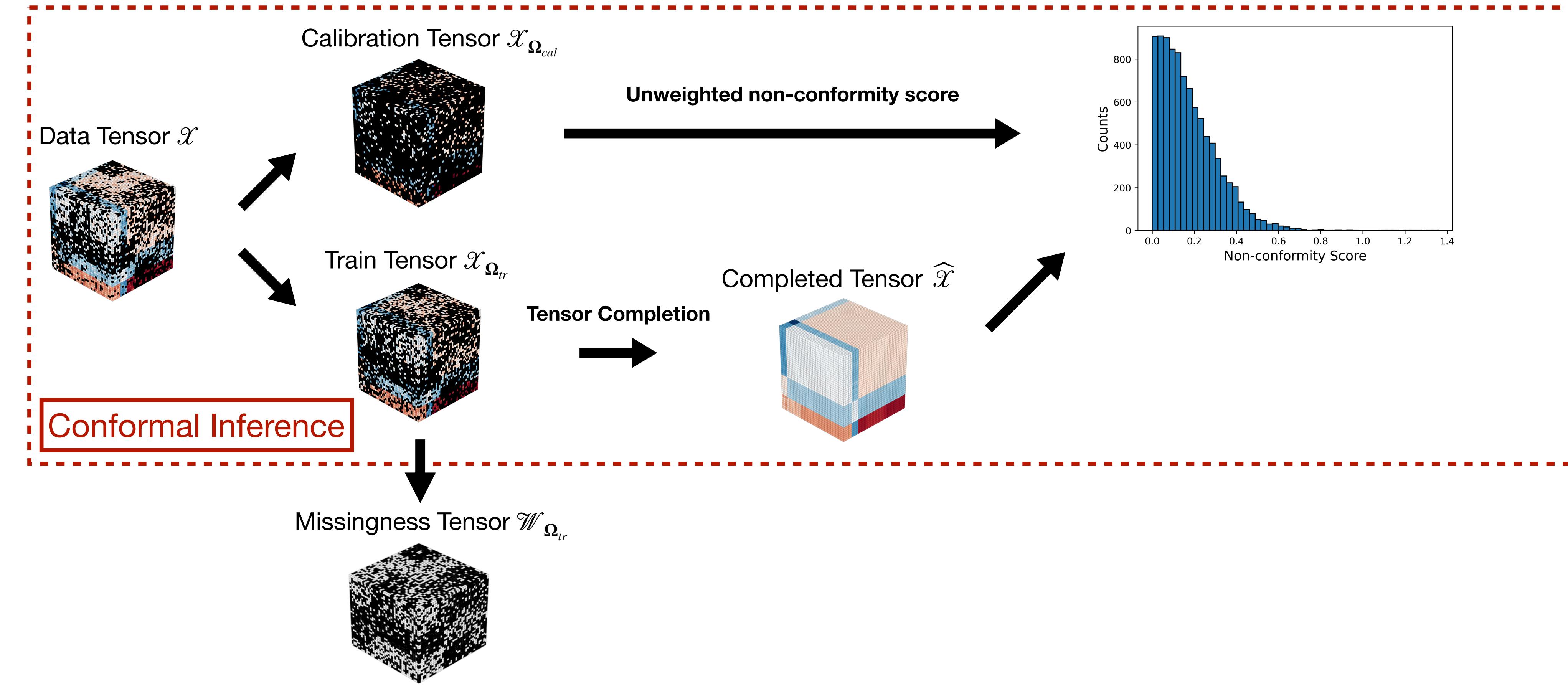


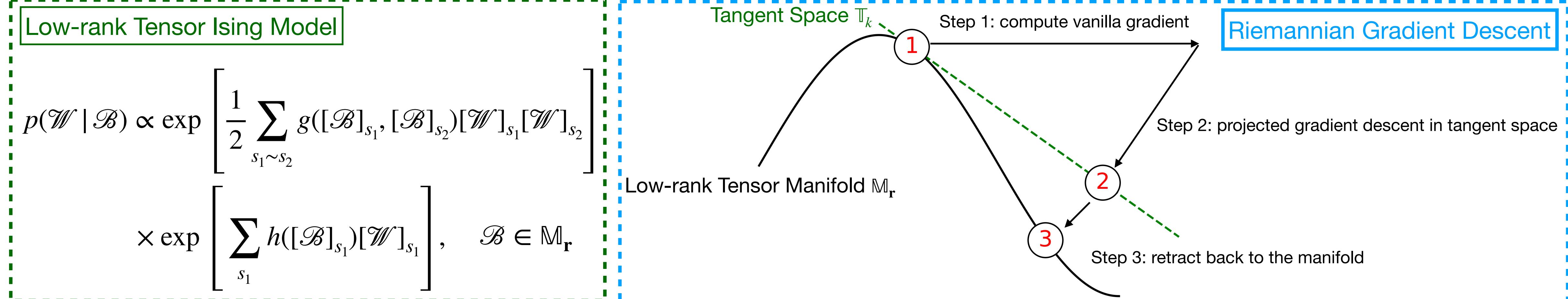
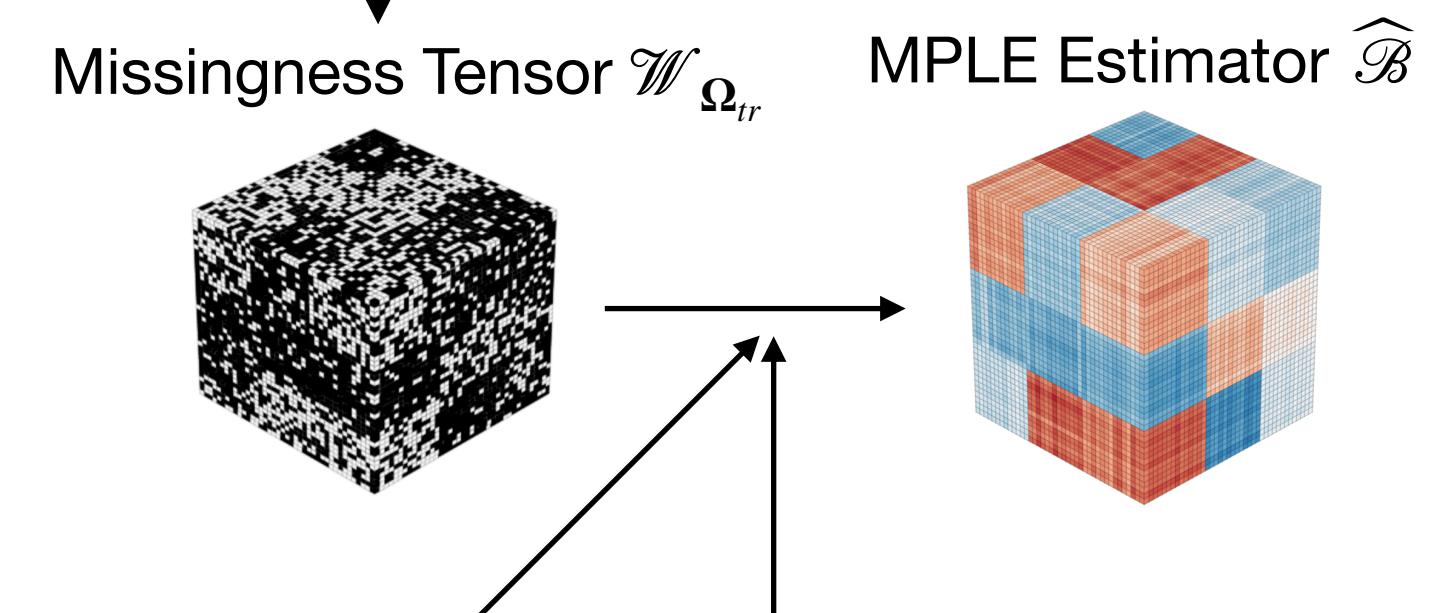
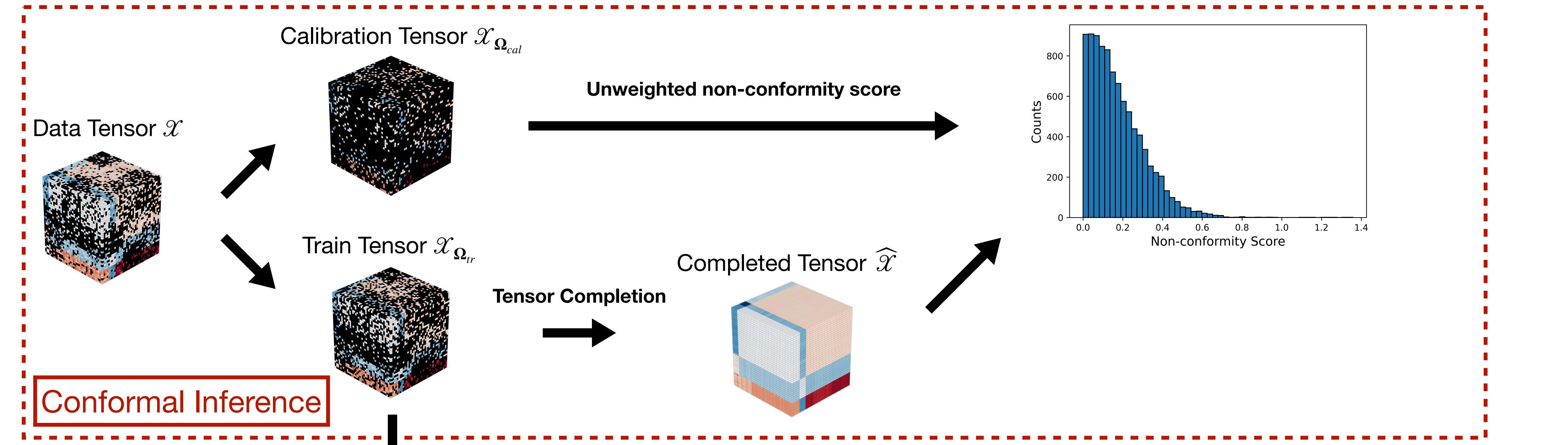


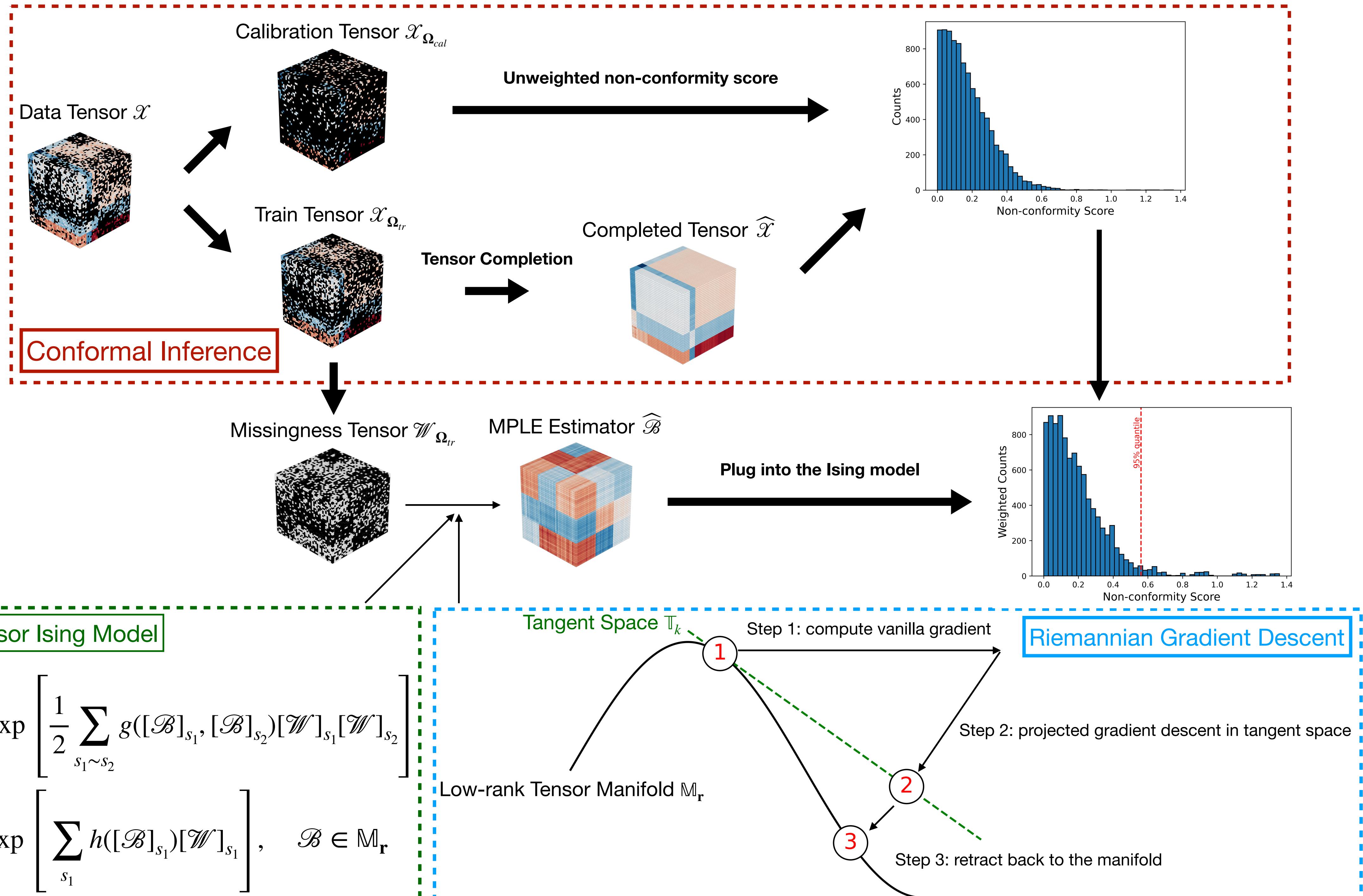
### Low-rank Tensor Ising Model

$$p(\mathcal{W} | \mathcal{B}) \propto \exp \left[ \frac{1}{2} \sum_{s_1 \sim s_2} g([\mathcal{B}]_{s_1}, [\mathcal{B}]_{s_2}) [\mathcal{W}]_{s_1} [\mathcal{W}]_{s_2} \right]$$

$$\times \exp \left[ \sum_{s_1} h([\mathcal{B}]_{s_1}) [\mathcal{W}]_{s_1} \right], \quad \mathcal{B} \in \mathbb{M}_r$$







# Theoretical Results

## MPLE error bound

# Theoretical Results

## MPLE error bound

**Theorem 2** Assume that:

- $g(\cdot, \cdot) = 0$
- $h(\cdot)$  is twice-continuously differentiable and  $h''(\cdot) \geq 0$
- $\|\mathcal{B}^*\|_\infty \leq \xi$

# Theoretical Results

## MPLE error bound

**Theorem 2** Assume that:

- $g(\cdot, \cdot) = 0$
- $h(\cdot)$  is twice-continuously differentiable and  $h''(\cdot) \geq 0$
- $\|\mathcal{B}^*\|_\infty \leq \xi$

$$P\left[\frac{1}{\sqrt{d^*}}\|\widehat{\mathcal{B}} - \mathcal{B}^*\|_F \leq 2C_{K,\xi}\sqrt{\frac{r^*\bar{d}}{d^*}}\right] \geq 1 - \exp(-C_1\bar{d}\log K),$$

where  $r^*, \bar{d}, d^*$  are  $\prod_k r_k, \sum_k d_k, \prod_k d_k$  for  $\mathcal{B}^* \in \mathbb{R}^{d_1 \times \dots \times d_K}$  with tensor-train rank  $\mathbf{r} = (r_1, \dots, r_{K-1})$ .

# Theoretical Results

## MPLE error bound

**Theorem 2** Assume that:

- $g(\cdot, \cdot) = 0$
- $h(\cdot)$  is twice-continuously differentiable and  $h''(\cdot) \geq 0$
- $\|\mathcal{B}^*\|_\infty \leq \xi$

$$P \left[ \frac{1}{\sqrt{d^*}} \|\widehat{\mathcal{B}} - \mathcal{B}^*\|_F \leq 2C_{K,\xi} \sqrt{\frac{r^* \bar{d}}{d^*}} \right] \geq 1 - \exp(-C_1 \bar{d} \log K),$$

↓  
**RMSE of the MPLE**

where  $r^*, \bar{d}, d^*$  are  $\prod_k r_k, \sum_k d_k, \prod_k d_k$  for  $\mathcal{B}^* \in \mathbb{R}^{d_1 \times \dots \times d_K}$  with tensor-train rank  $\mathbf{r} = (r_1, \dots, r_{K-1})$ .

# Theoretical Results

## MPLE error bound

**Theorem 2** Assume that:

- $g(\cdot, \cdot) = 0$
- $h(\cdot)$  is twice-continuously differentiable and  $h''(\cdot) \geq 0$
- $\|\mathcal{B}^*\|_\infty \leq \xi$

positive constants related to  $K$  and  $\xi$

$$P\left[\frac{1}{\sqrt{d^*}}\|\widehat{\mathcal{B}} - \mathcal{B}^*\|_F \leq 2C_{K,\xi}\sqrt{\frac{r^*\bar{d}}{d^*}}\right] \geq 1 - \exp(-C_1\bar{d}\log K),$$

RMSE of the MPLE

where  $r^*, \bar{d}, d^*$  are  $\prod_k r_k$ ,  $\sum_k d_k$ ,  $\prod_k d_k$  for  $\mathcal{B}^* \in \mathbb{R}^{d_1 \times \dots \times d_K}$  with tensor-train rank  $\mathbf{r} = (r_1, \dots, r_{K-1})$ .

# Theoretical Results

## MPLE error bound

**Theorem 2** Assume that:

- $g(\cdot, \cdot) = 0$
- $h(\cdot)$  is twice-continuously differentiable and  $h''(\cdot) \geq 0$
- $\|\mathcal{B}^*\|_\infty \leq \xi$

positive constants related to  $K$  and  $\xi$

$$P \left[ \frac{1}{\sqrt{d^*}} \|\widehat{\mathcal{B}} - \mathcal{B}^*\|_F \leq 2C_{K,\xi} \sqrt{\frac{r^* \bar{d}}{d^*}} \right] \geq 1 - \exp(-C_1 \bar{d} \log K),$$

**RMSE of the MPLE**

$$\asymp (r/d)^{(K-1)/2}, \text{ if } r_k \asymp r, d_k \asymp d, \forall k$$

where  $r^*, \bar{d}, d^*$  are  $\prod_k r_k, \sum_k d_k, \prod_k d_k$  for  $\mathcal{B}^* \in \mathbb{R}^{d_1 \times \dots \times d_K}$  with tensor-train rank  $\mathbf{r} = (r_1, \dots, r_{K-1})$ .

# Theoretical Results

## Coverage Guarantee

**Theorem 3** Assume the same assumptions as Theorem 2, then for any  $0 < c < 1$ :

# Theoretical Results

## Coverage Guarantee

**Theorem 3** Assume the same assumptions as Theorem 2, then for any  $0 < c < 1$ :

$$\begin{aligned} \text{E} \left[ \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{1-\alpha,s}\}} \right] &\geq 1 - \alpha - \frac{2c_{K,\xi}}{(1-c)(1-q)} \sqrt{\frac{r^* \bar{d}}{d^*}} \\ &\quad - \exp[-C_1 \bar{d} \log K] - \exp \left[ -\frac{c^2(1-q)d^* l_\xi}{2} \right] \end{aligned}$$

where  $q \in (0,1)$  is the train-calibration split ratio.

# Theoretical Results

## Coverage Guarantee

**Theorem 3** Assume the same assumptions as Theorem 2, then for any  $0 < c < 1$ :

average coverage on  $\Omega^c$

$$\mathbb{E} \left[ \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{1-\alpha,s}\}} \right] \geq 1 - \alpha - \frac{2c_{K,\xi}}{(1-c)(1-q)} \sqrt{\frac{r^* \bar{d}}{d^*}} - \exp[-C_1 \bar{d} \log K] - \exp \left[ -\frac{c^2(1-q)d^* l_\xi}{2} \right]$$

where  $q \in (0,1)$  is the train-calibration split ratio .

# Theoretical Results

## Coverage Guarantee

**Theorem 3** Assume the same assumptions as Theorem 2, then for any  $0 < c < 1$ :

$$\begin{aligned}
 & \text{average coverage on } \Omega^c \quad \text{target coverage} \\
 & \uparrow \qquad \uparrow \\
 & \mathbb{E} \left[ \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{1-\alpha,s}\}} \right] \geq 1 - \alpha - \frac{2c_{K,\xi}}{(1-c)(1-q)} \sqrt{\frac{r^* \bar{d}}{d^*}} \\
 & \quad - \exp[-C_1 \bar{d} \log K] - \exp \left[ -\frac{c^2(1-q)d^*l_\xi}{2} \right]
 \end{aligned}$$

where  $q \in (0,1)$  is the train-calibration split ratio.

# Theoretical Results

## Coverage Guarantee

**Theorem 3** Assume the same assumptions as Theorem 2, then for any  $0 < c < 1$ :

$$\begin{aligned}
 & \text{average coverage on } \Omega^c \quad \text{target coverage} \quad \text{MPLE estimation error} \\
 & \uparrow \qquad \uparrow \qquad \nearrow \\
 & E \left[ \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{1-\alpha,s}\}} \right] \geq 1 - \alpha - \frac{2c_{K,\xi}}{(1-c)(1-q)} \sqrt{\frac{r^* \bar{d}}{d^*}} \\
 & \quad - \exp[-C_1 \bar{d} \log K] - \exp \left[ -\frac{c^2(1-q)d^* l_\xi}{2} \right]
 \end{aligned}$$

where  $q \in (0,1)$  is the train-calibration split ratio.

# Theoretical Results

## Coverage Guarantee

**Theorem 3** Assume the same assumptions as Theorem 2, then for any  $0 < c < 1$ :

$$E \left[ \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{1-\alpha,s}\}} \right] \geq 1 - \alpha - \frac{2c_{K,\xi}}{(1-c)(1-q)} \sqrt{\frac{r^* \bar{d}}{d^*}} - \exp[-C_1 \bar{d} \log K] - \exp \left[ -\frac{c^2(1-q)d^* l_\xi}{2} \right]$$

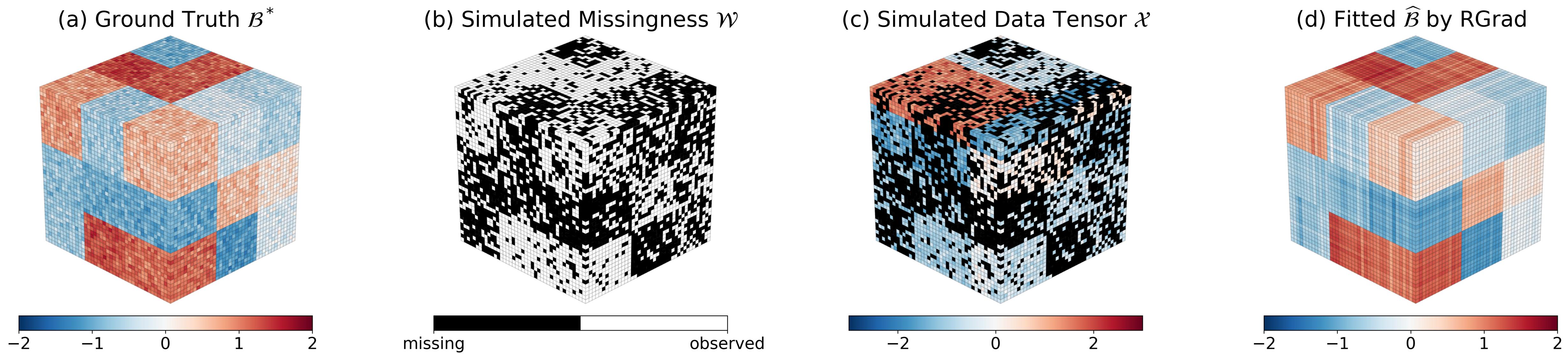
average coverage on  $\Omega^c$     target coverage    MPLE estimation error

where  $q \in (0,1)$  is the train-calibration split ratio .

sample size requirement for calibration

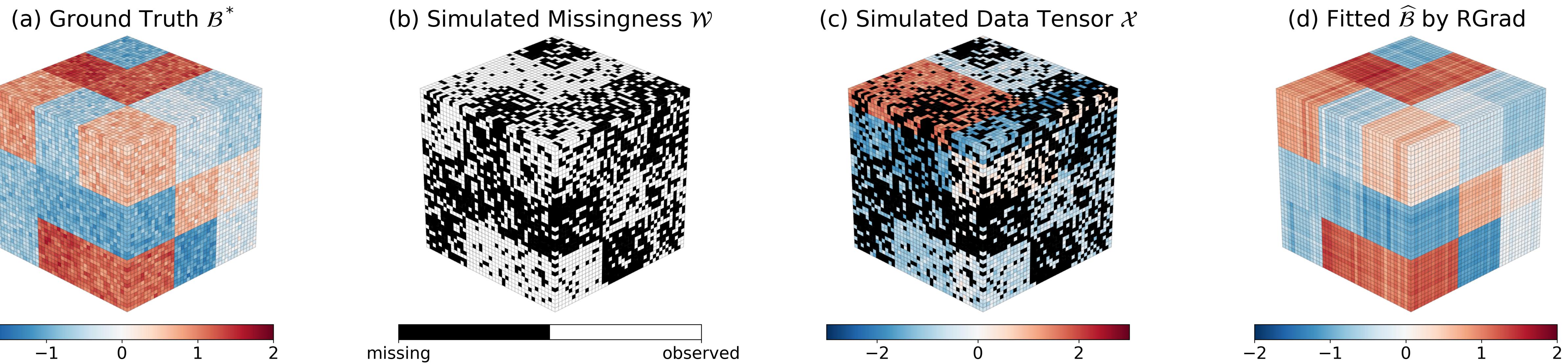
# Numerical Experiments

## Setup



# Numerical Experiments

## Setup



red: more likely to observe

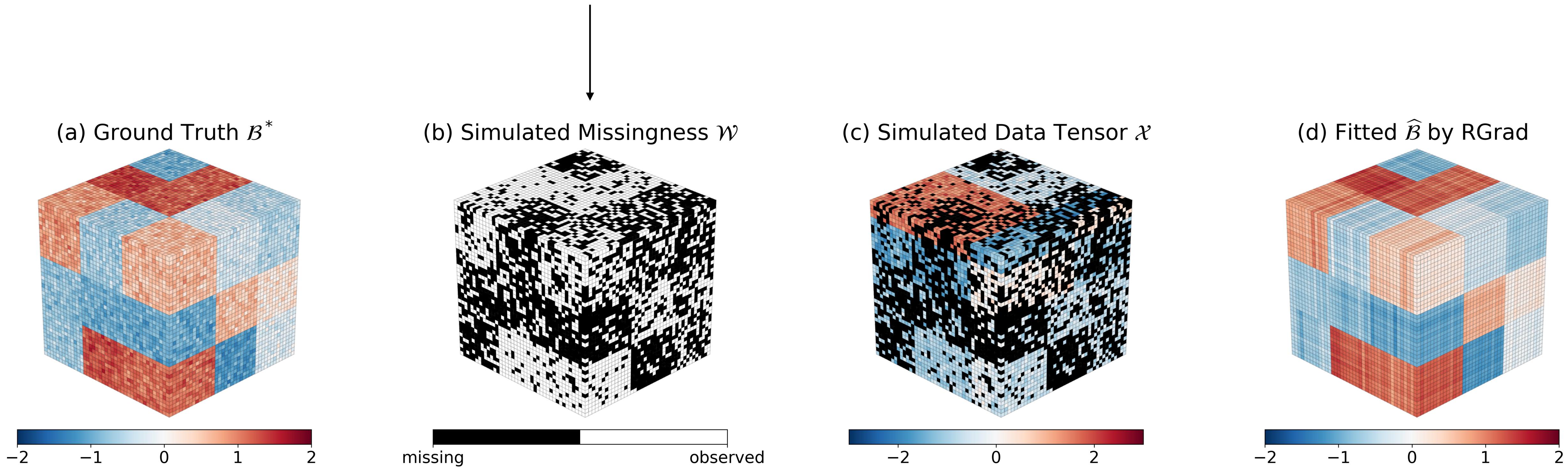
blue: more likely to miss

data tend to miss/observe together in blocks

# Numerical Experiments

## Setup

Simulated by MCMC with block-Gibbs sampler



red: more likely to observe

blue: more likely to miss

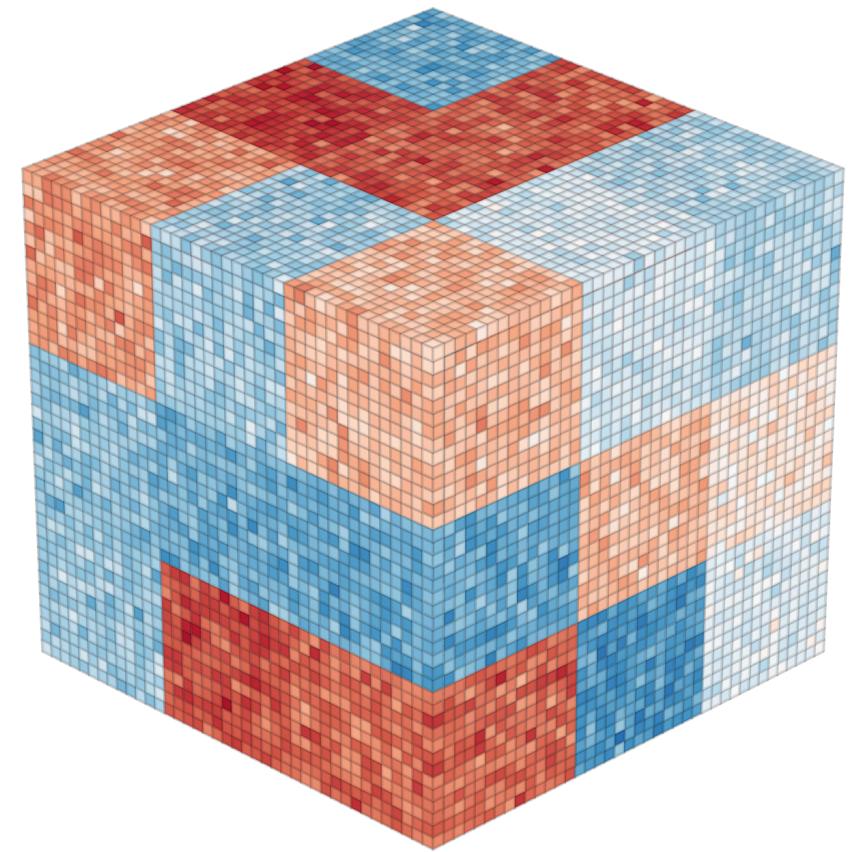
data tend to miss/observe together in blocks

# Numerical Experiments

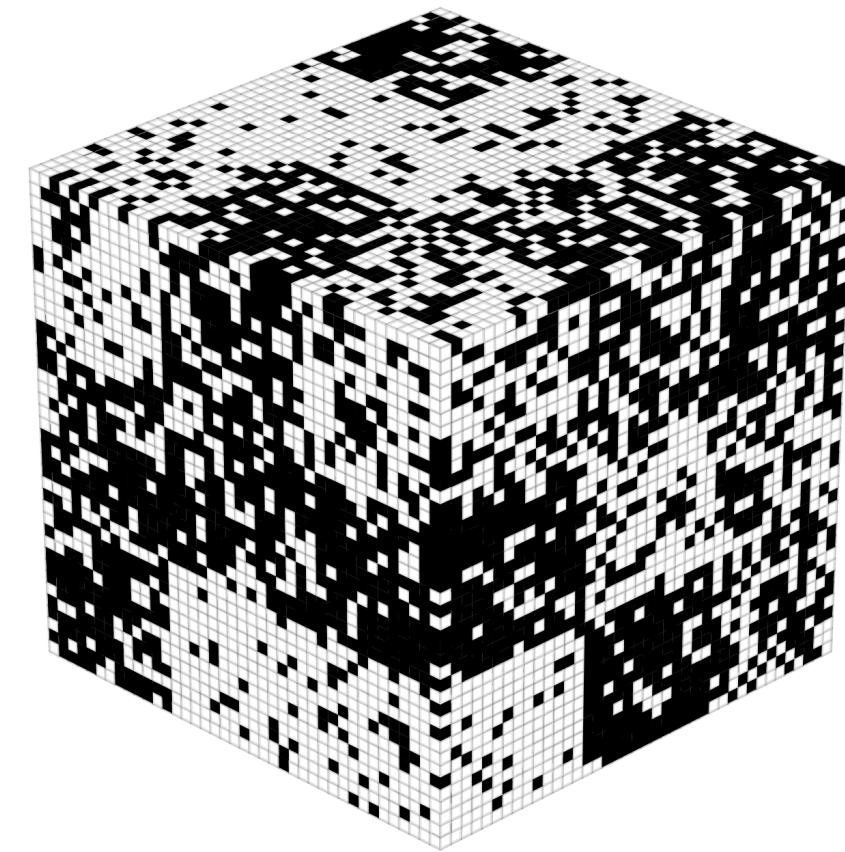
## Setup

Simulated by MCMC with block-Gibbs sampler

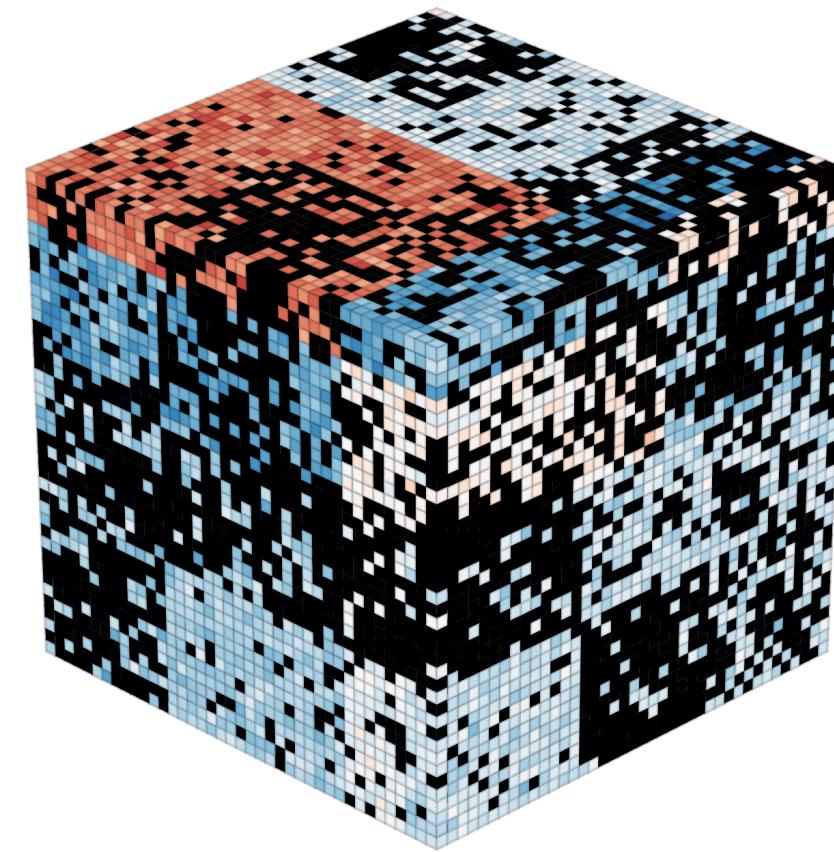
(a) Ground Truth  $\mathcal{B}^*$



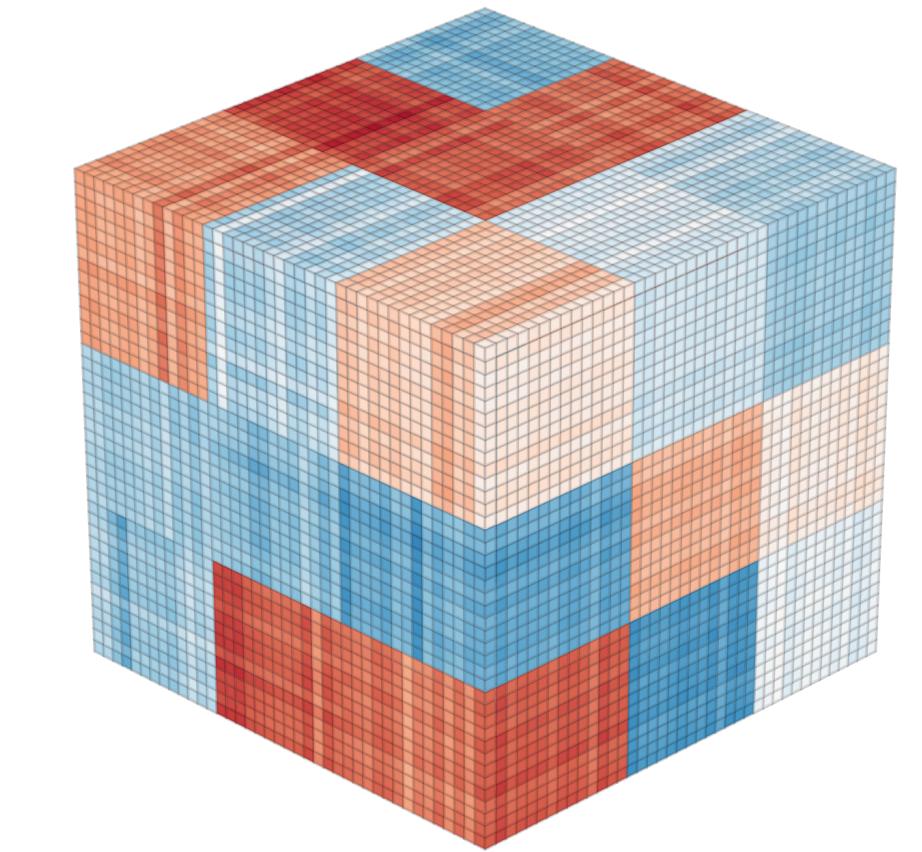
(b) Simulated Missingness  $\mathcal{W}$



(c) Simulated Data Tensor  $\mathcal{X}$



(d) Fitted  $\hat{\mathcal{B}}$  by RGrad



red: more likely to observe

blue: more likely to miss

data tend to miss/observe together in blocks

$$\mathcal{X} = \mathcal{X}^* + \mathcal{E}, \mathcal{X}^* \text{ is low rank}$$

# Numerical Experiments

## Setup

# Numerical Experiments

## Setup

- We consider two data missing patterns:

# Numerical Experiments

## Setup

- We consider two data missing patterns:

► **Bernoulli**:  $g(x, y) = 0, h(x) = x/2$ , so  $[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}\left(\frac{\exp([\mathcal{B}]_s)}{1 + \exp([\mathcal{B}]_s)}\right)$

# Numerical Experiments

## Setup

- We consider two data missing patterns:

- ▶ **Bernoulli**:  $g(x, y) = 0, h(x) = x/2$ , so  $[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}\left(\frac{\exp([\mathcal{B}]_s)}{1 + \exp([\mathcal{B}]_s)}\right)$
- ▶ **Ising**:  $g(x, y) = xy/15, h(x) = x/2$

# Numerical Experiments

## Setup

- We consider two data missing patterns:
  - ▶ **Bernoulli**:  $g(x, y) = 0, h(x) = x/2$ , so  $[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}\left(\frac{\exp([\mathcal{B}]_s)}{1 + \exp([\mathcal{B}]_s)}\right)$
  - ▶ **Ising**:  $g(x, y) = xy/15, h(x) = x/2$
- We consider two types of noise distributions:

# Numerical Experiments

## Setup

- We consider two data missing patterns:
  - ▶ **Bernoulli**:  $g(x, y) = 0, h(x) = x/2$ , so  $[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}\left(\frac{\exp([\mathcal{B}]_s)}{1 + \exp([\mathcal{B}]_s)}\right)$
  - ▶ **Ising**:  $g(x, y) = xy/15, h(x) = x/2$
- We consider two types of noise distributions:
  - **constant**:  $[\mathcal{E}]_s \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$

# Numerical Experiments

## Setup

- We consider two data missing patterns:
  - ▶ **Bernoulli**:  $g(x, y) = 0, h(x) = x/2$ , so  $[\mathcal{W}]_s \stackrel{ind.}{\sim} \text{Bern}\left(\frac{\exp([\mathcal{B}]_s)}{1 + \exp([\mathcal{B}]_s)}\right)$
  - ▶ **Ising**:  $g(x, y) = xy/15, h(x) = x/2$
- We consider two types of noise distributions:
  - **constant**:  $[\mathcal{E}]_s \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$
  - **adversarial**:  $[\mathcal{E}]_s \stackrel{ind.}{\sim} N(0, \sigma_s^2), \quad \sigma_s \propto [\mathcal{B}]_s^{-1}$  (higher uncertainty for missing entries)

# Numerical Experiments

## Evaluation Metrics & Benchmarks

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:
  - **Unweighted:** apply equal weight to all calibration entries;

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:
  - **Unweighted:** apply equal weight to all calibration entries;
  - **Oracle:** assume that the missing propensity is known;

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:
  - **Unweighted:** apply equal weight to all calibration entries;
  - **Oracle:** assume that the missing propensity is known;
  - **RGrad:** learn the missing propensity via tensor Ising model.

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:
  - **Unweighted:** apply equal weight to all calibration entries;
  - **Oracle:** assume that the missing propensity is known;
  - **RGrad:** learn the missing propensity via tensor Ising model.
- We evaluate each method by the average mis-coverage % in  $\mathbb{Q}$ :

# Numerical Experiments

## Evaluation Metrics & Benchmarks

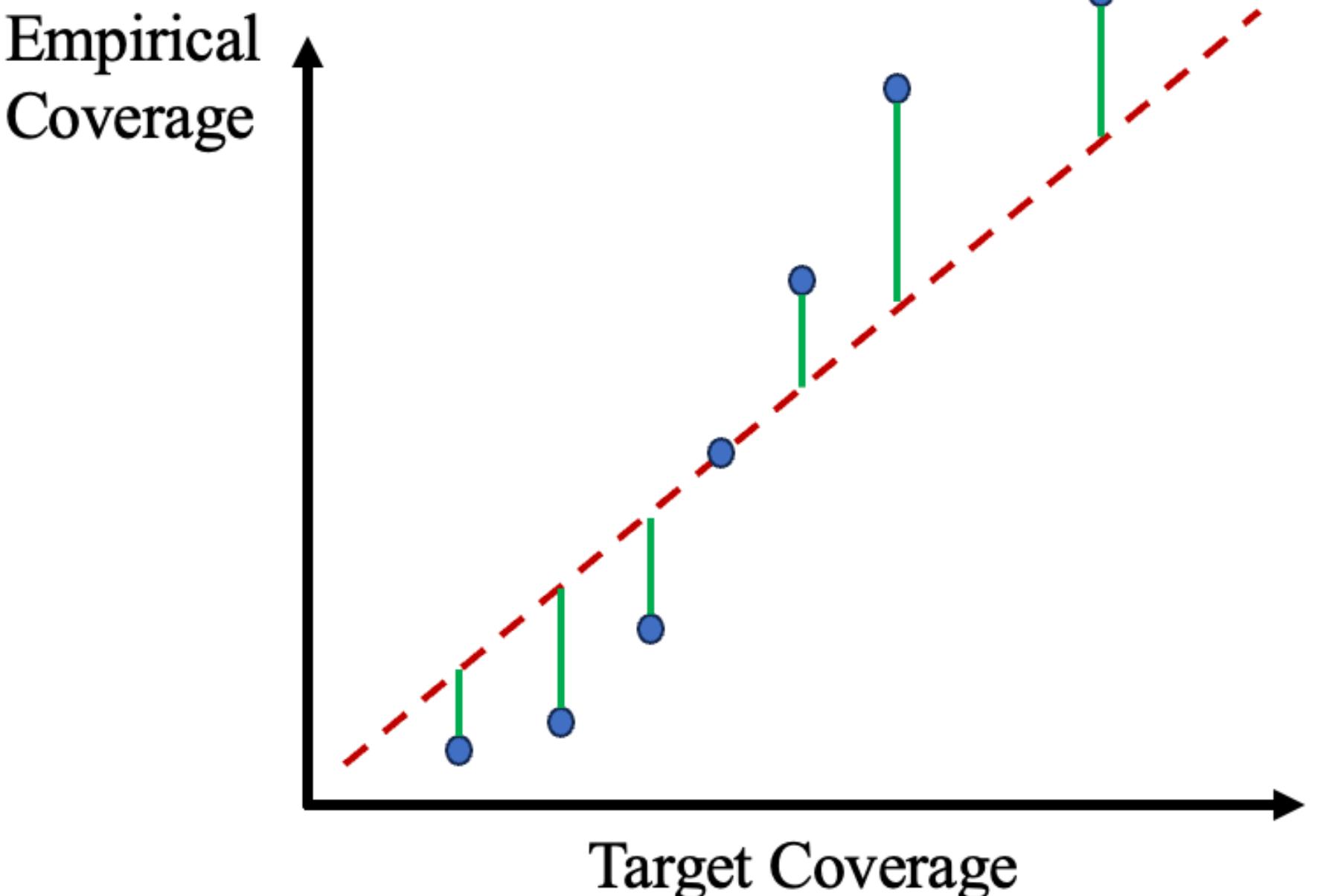
- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:
  - **Unweighted:** apply equal weight to all calibration entries;
  - **Oracle:** assume that the missing propensity is known;
  - **RGrad:** learn the missing propensity via tensor Ising model.
- We evaluate each method by the average mis-coverage % in  $\mathbb{Q}$ :

$$\frac{1}{|\mathbb{Q}|} \sum_{\tau \in \mathbb{Q}} \left| \frac{\text{target coverage}_{\tau} - \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{\tau,s}\}}}{\text{empirical coverage at } \tau} \right|$$

# Numerical Experiments

## Evaluation Metrics & Benchmarks

- We choose 20 coverage levels:  $\mathbb{Q} = \{0.80, 0.81, \dots, 0.99\}$
- We consider three versions of conformal prediction:
  - **Unweighted**: apply equal weight to all calibration entries;
  - **Oracle**: assume that the missing propensity is known;
  - **RGrad**: learn the missing propensity via tensor Ising model.
- We evaluate each method by the average mis-coverage % in  $\mathbb{Q}$ :

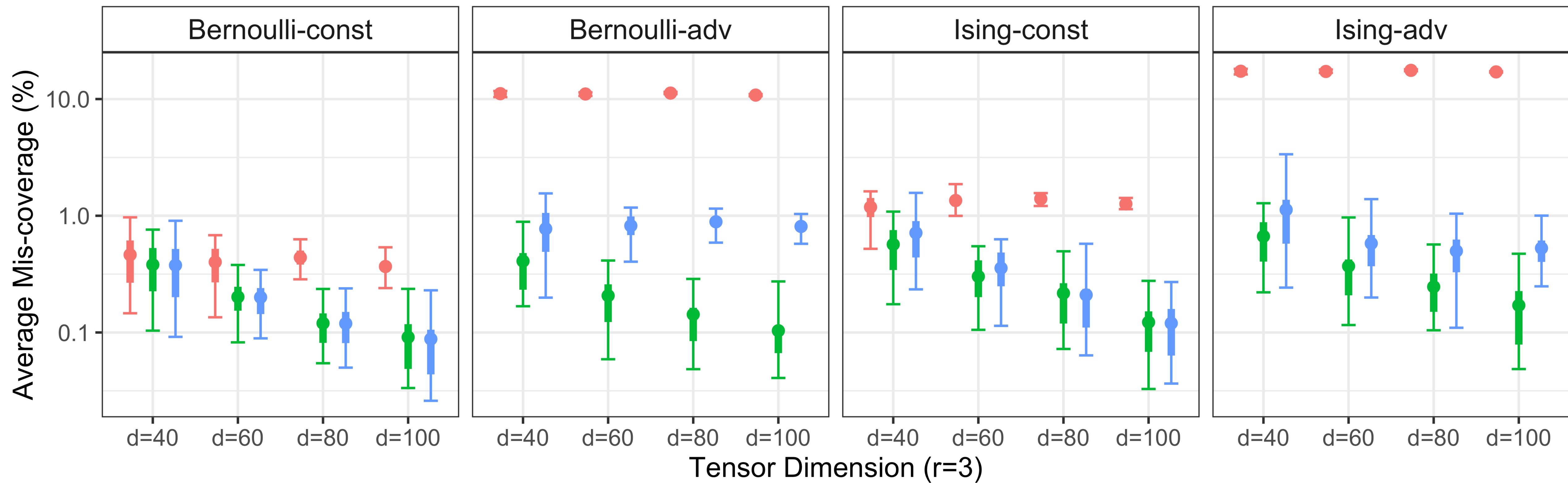


$$\frac{1}{|\mathbb{Q}|} \sum_{\tau \in \mathbb{Q}} \left| \underbrace{\text{target coverage} - \frac{1}{|\Omega^c|} \sum_{s \in \Omega^c} \mathbb{I}_{\{[\mathcal{X}]_s \in \widehat{\mathcal{C}}_{\tau,s}\}}}_{\text{empirical coverage at } \tau} \right|$$

# Numerical Experiments

## Conformal Inference Validation ( $\text{rank}(\mathcal{B}^*) = (3,3)$ )

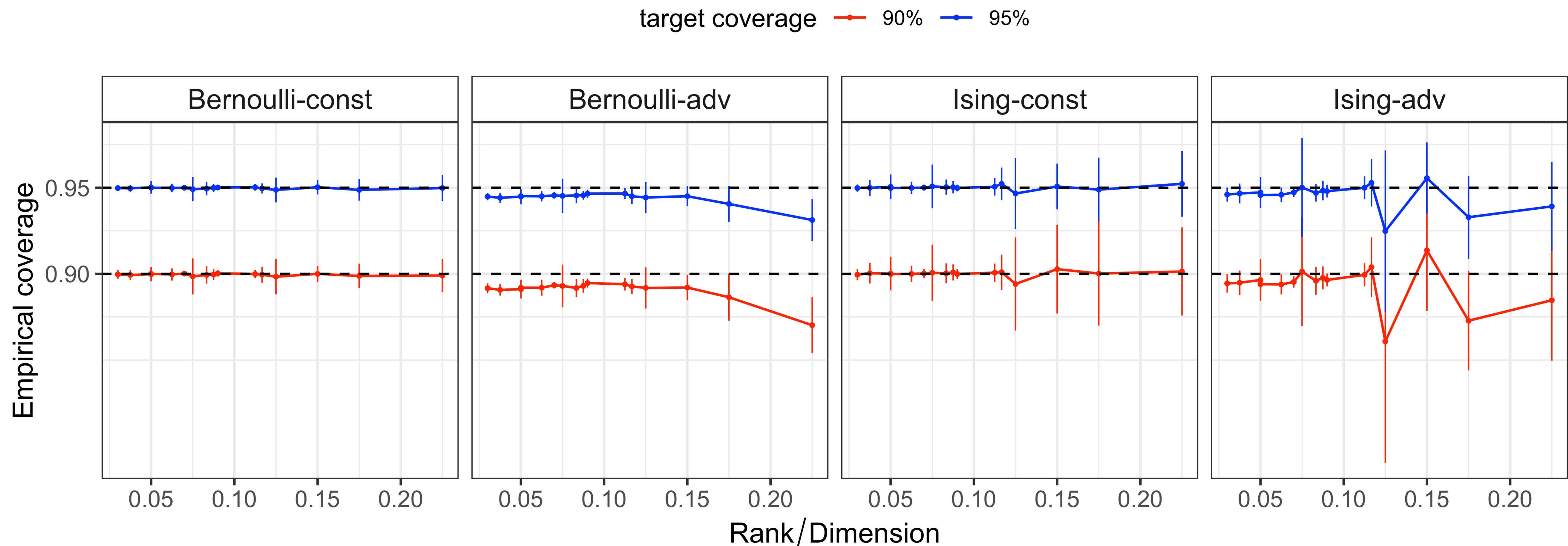
weight unweighted oracle RGrad



const:  $[\mathcal{E}]_s \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ ; adv:  $[\mathcal{E}]_s \stackrel{ind.}{\sim} N(0, \sigma_s^2)$ ,  $\sigma_s \propto [\mathcal{B}]_s^{-1}$  (higher uncertainty for missing entries)

# Numerical Experiments

## Conformal Inference Validation (90%, 95% CI)



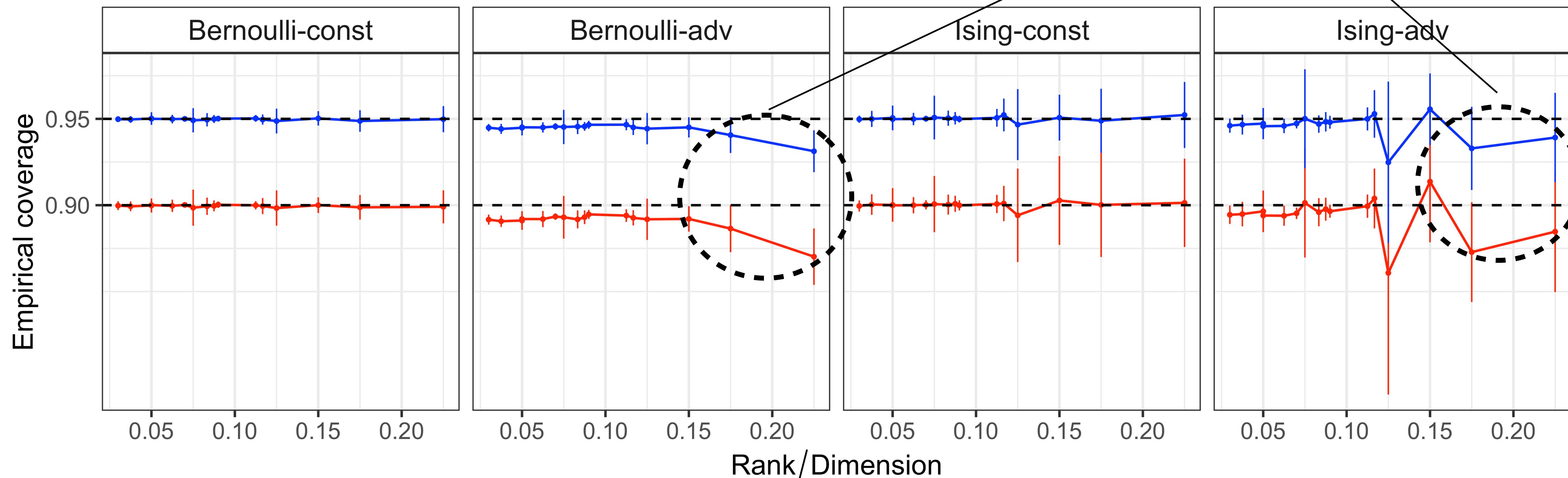
x-axis: r/d of the tensor parameter  $\mathcal{B}^*$

# Numerical Experiments

## Conformal Inference Validation (90%, 95% CI)

higher r/d leads to higher error of the MPLE

target coverage    —●— 90%    —●— 95%

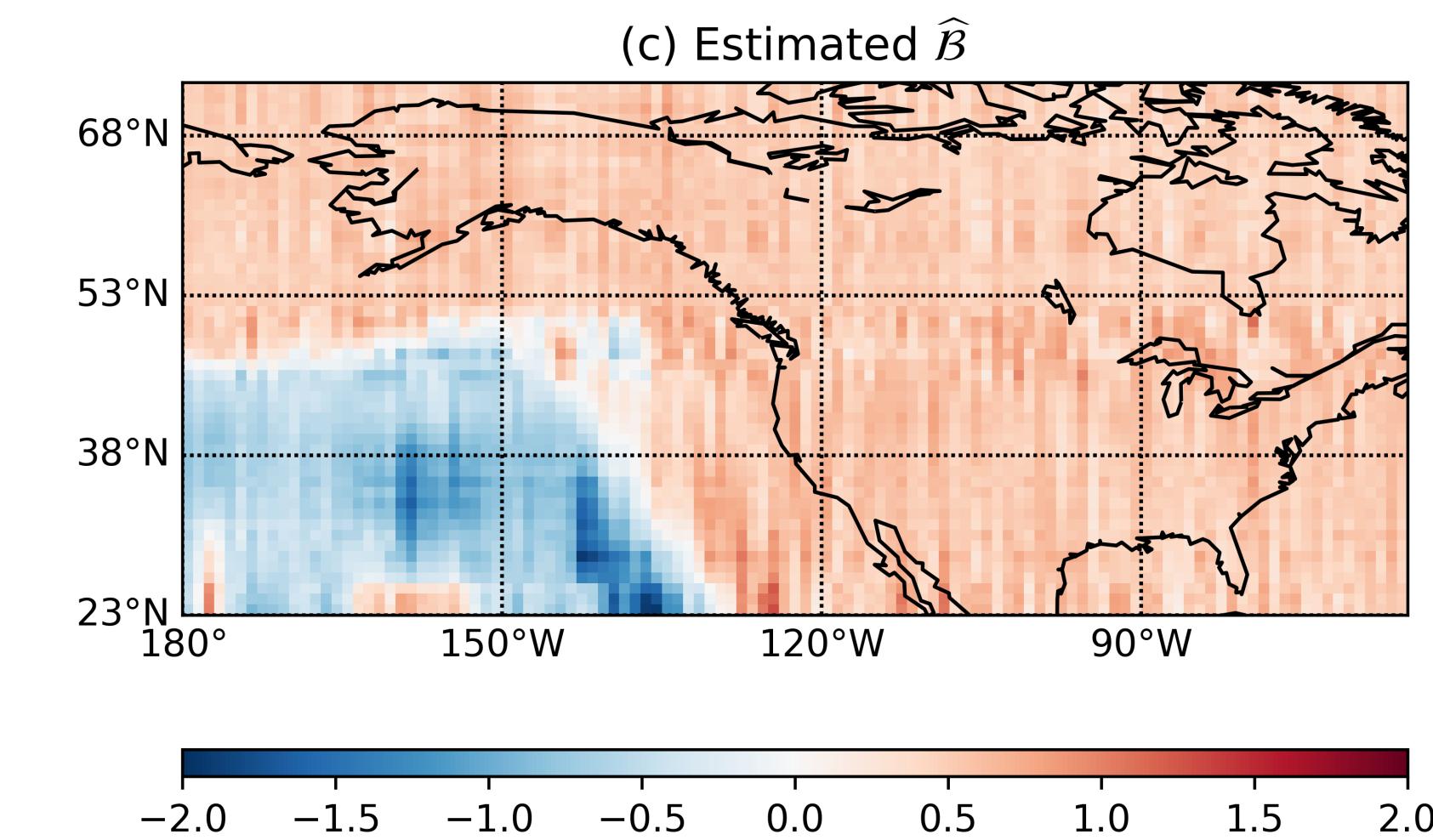
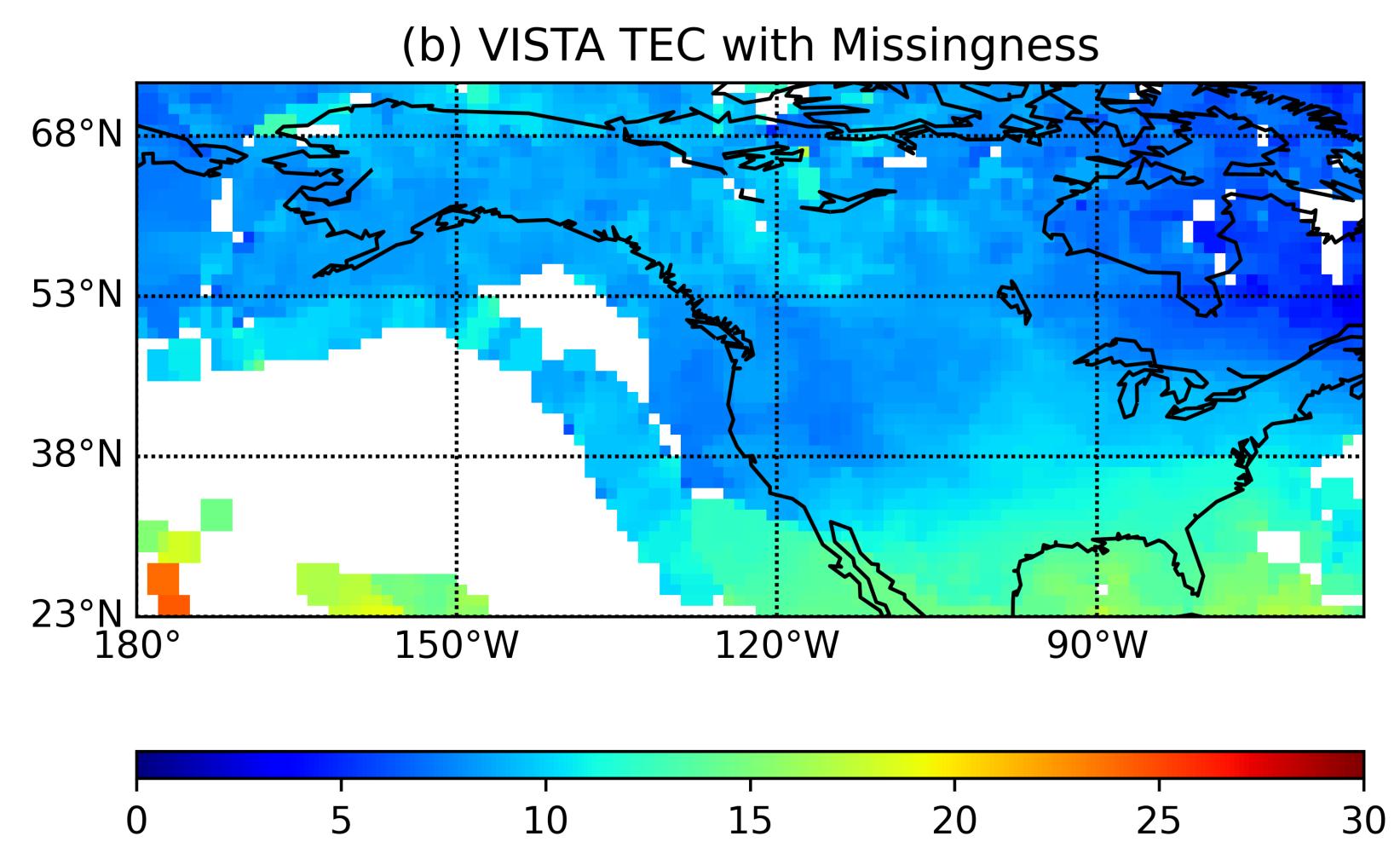
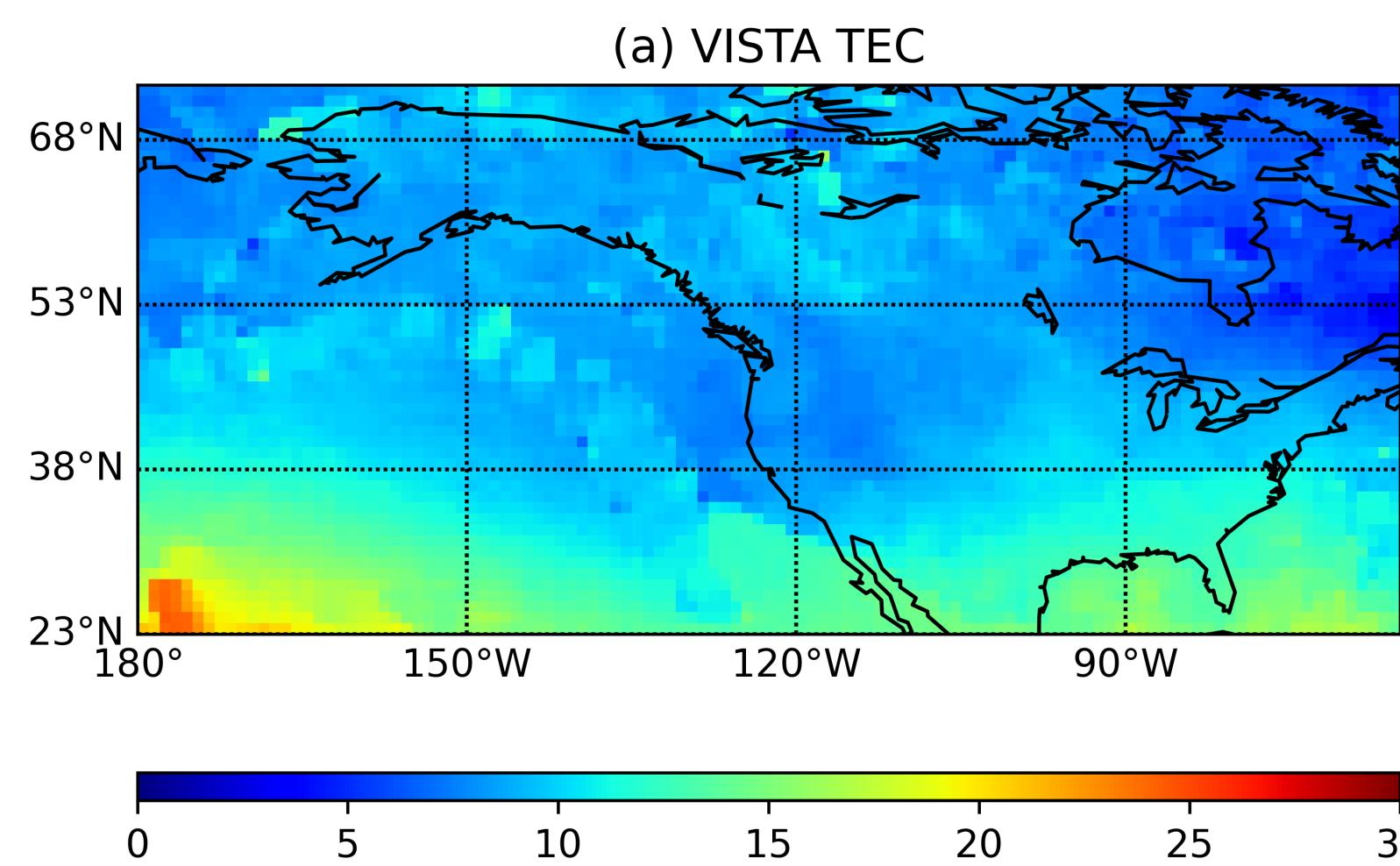


x-axis: r/d of the tensor parameter  $\mathcal{B}^*$

# Data Application

## Regional TEC Reconstruction

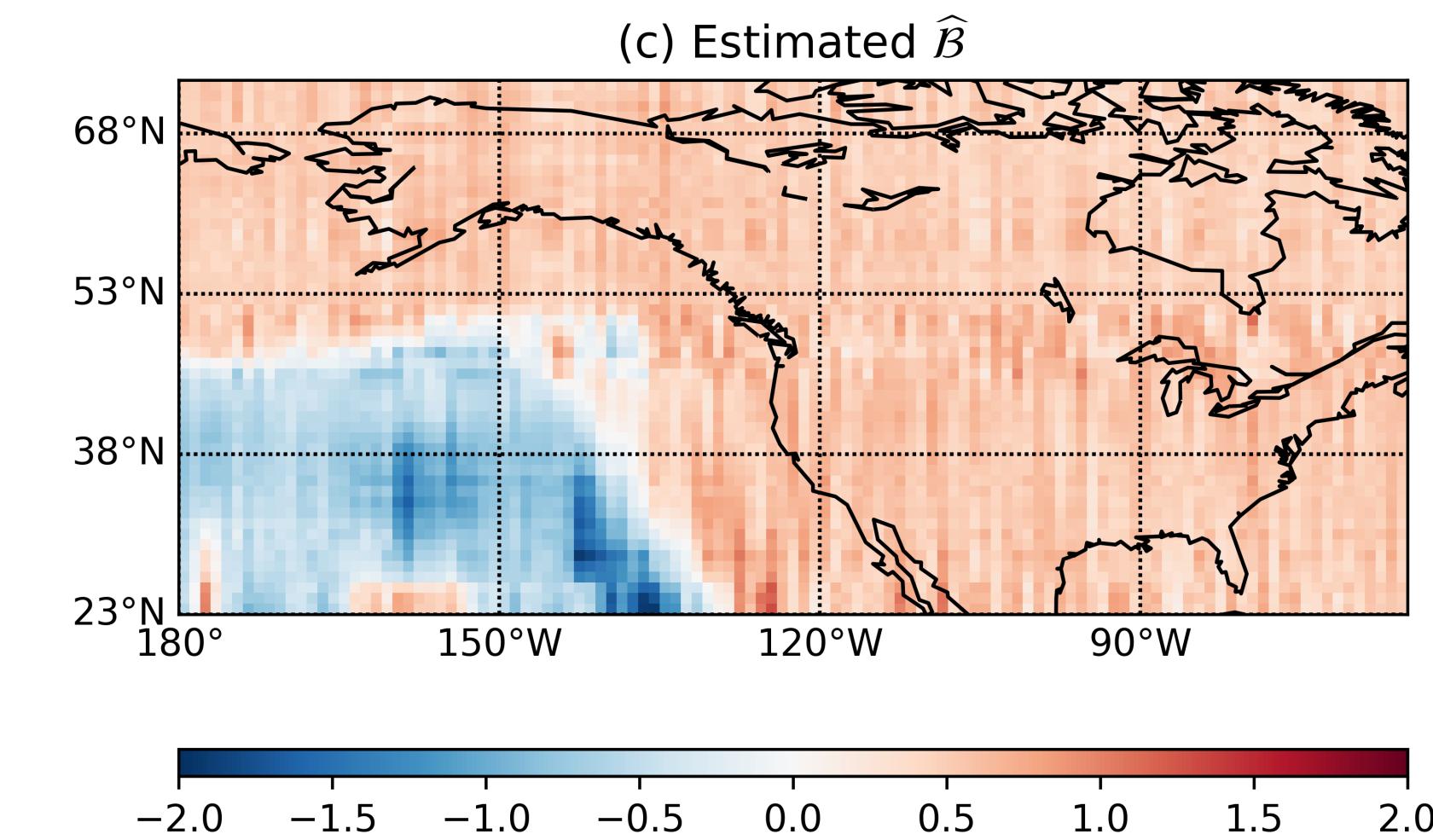
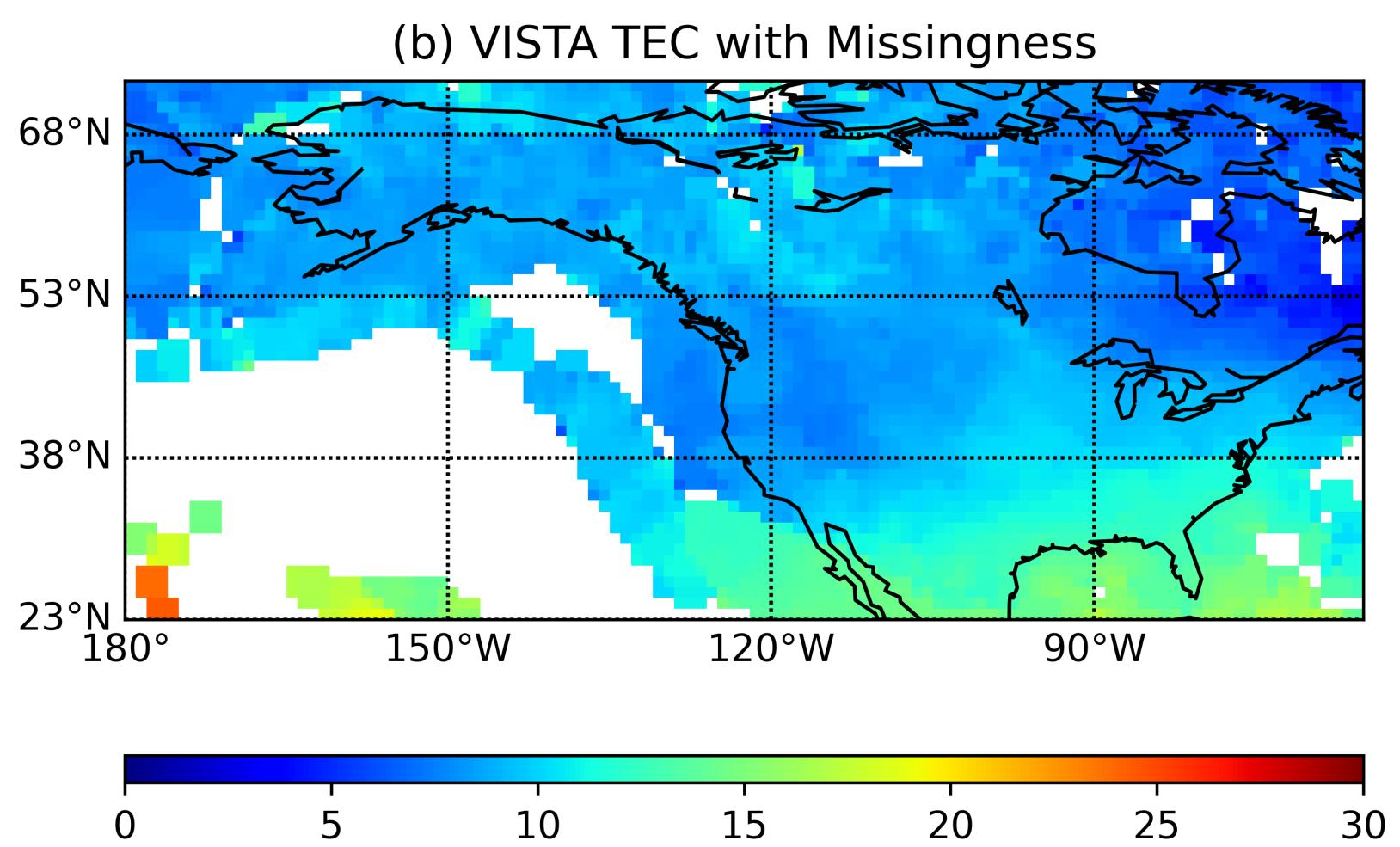
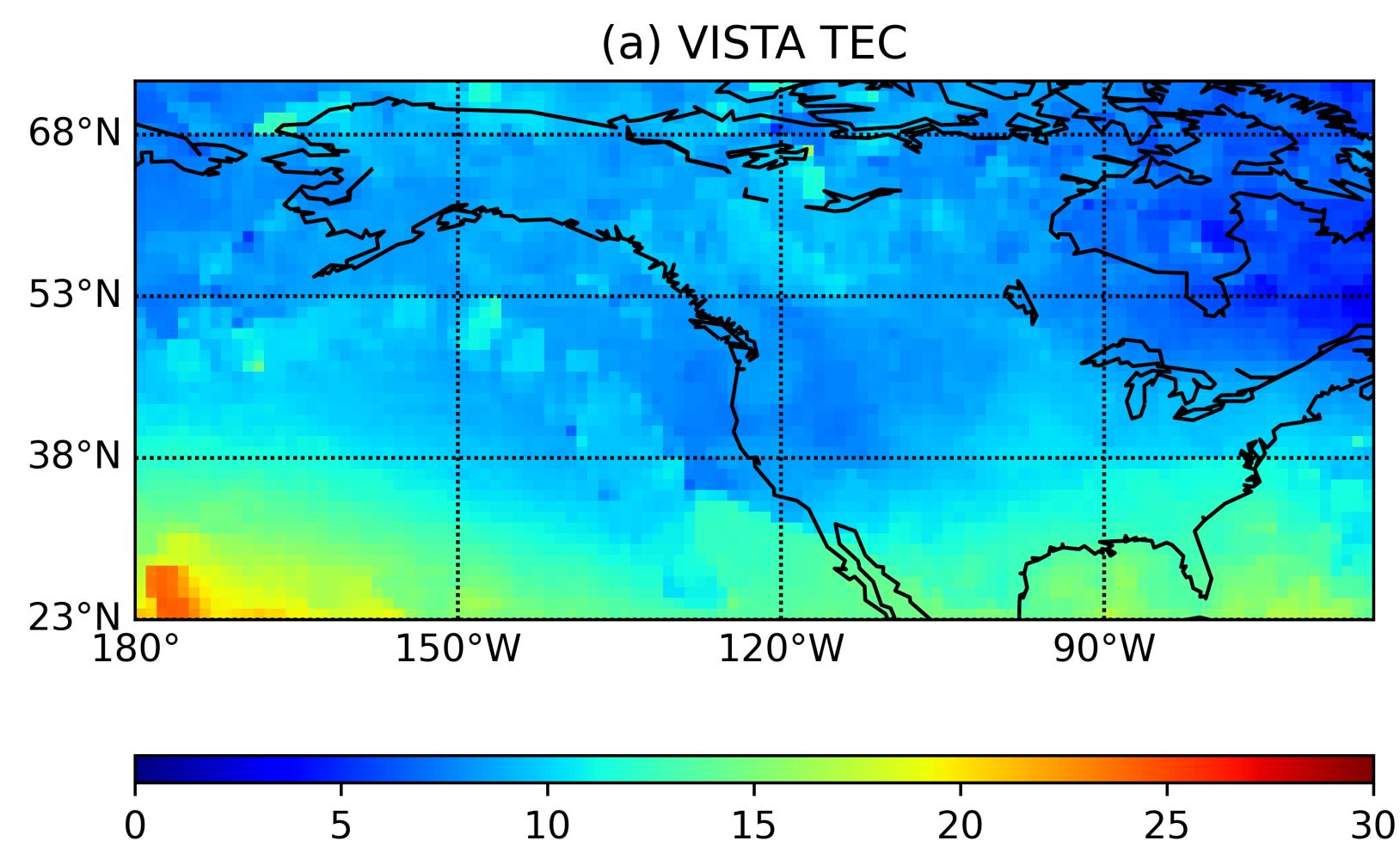
- We apply our method to reconstruct the total electron content (TEC) over the US and Canada for 15 days in Sept, 2017.



# Data Application

## Regional TEC Reconstruction

- We apply our method to reconstruct the total electron content (TEC) over the US and Canada for 15 days in Sept, 2017.

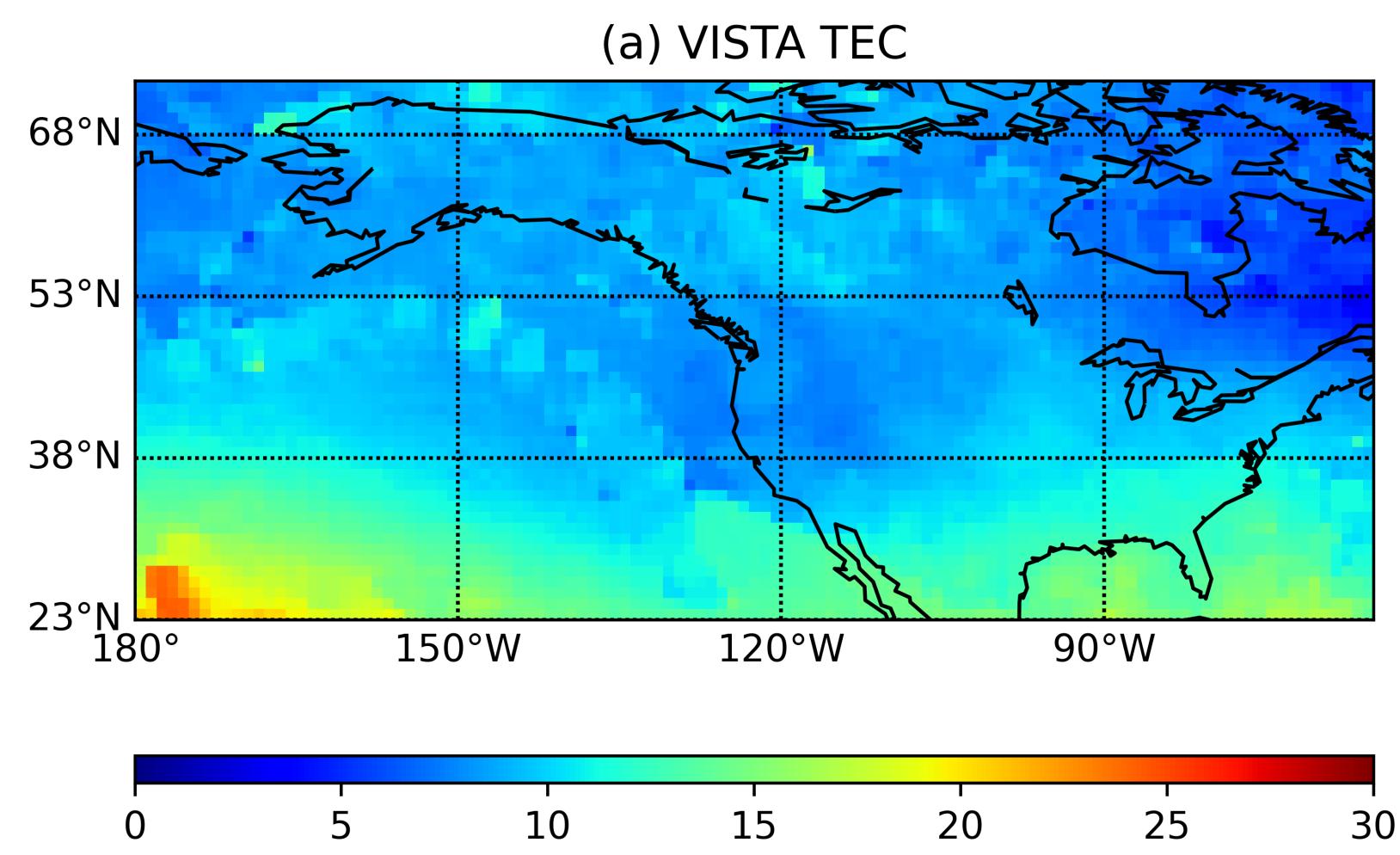


data without missingness

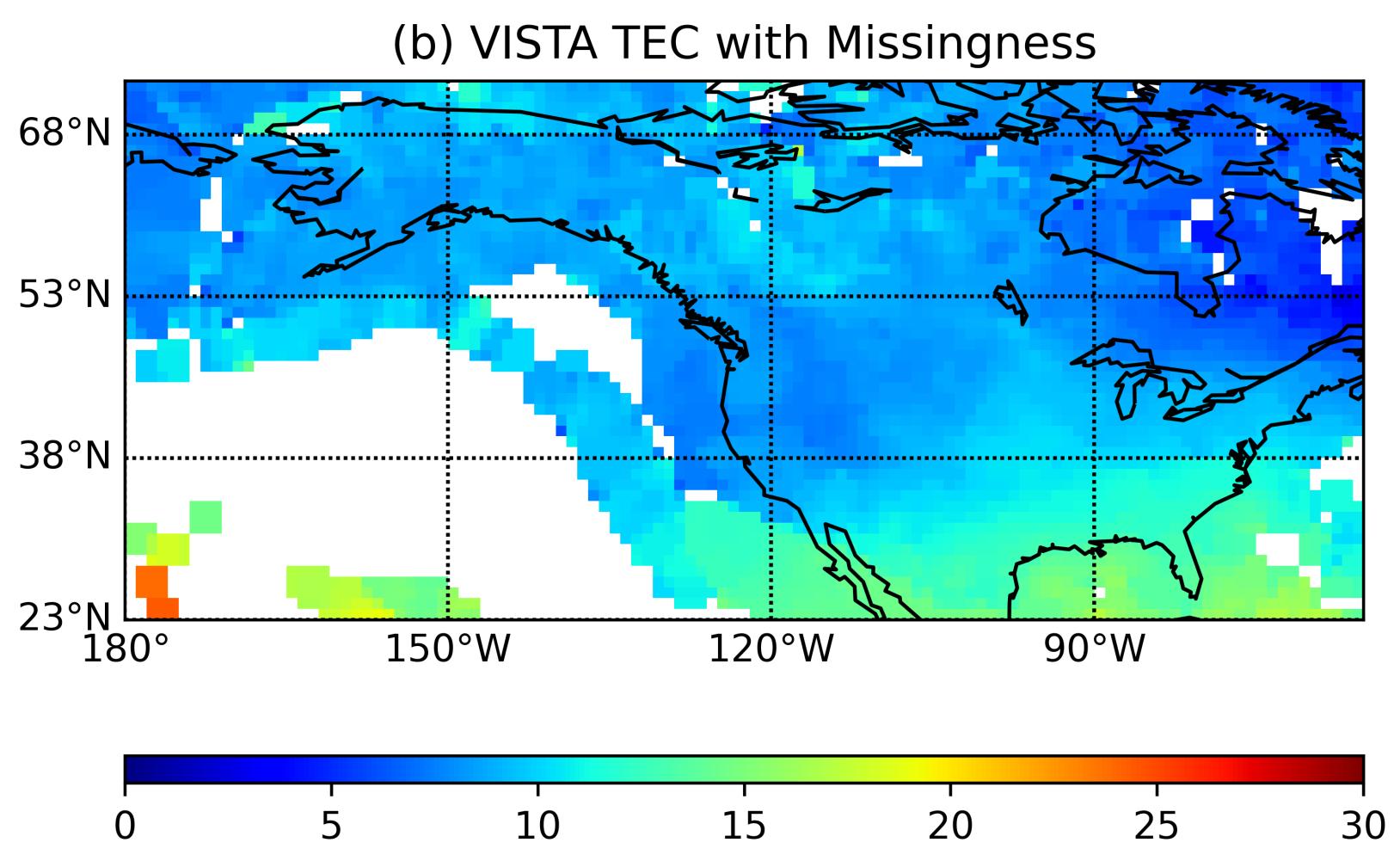
# Data Application

## Regional TEC Reconstruction

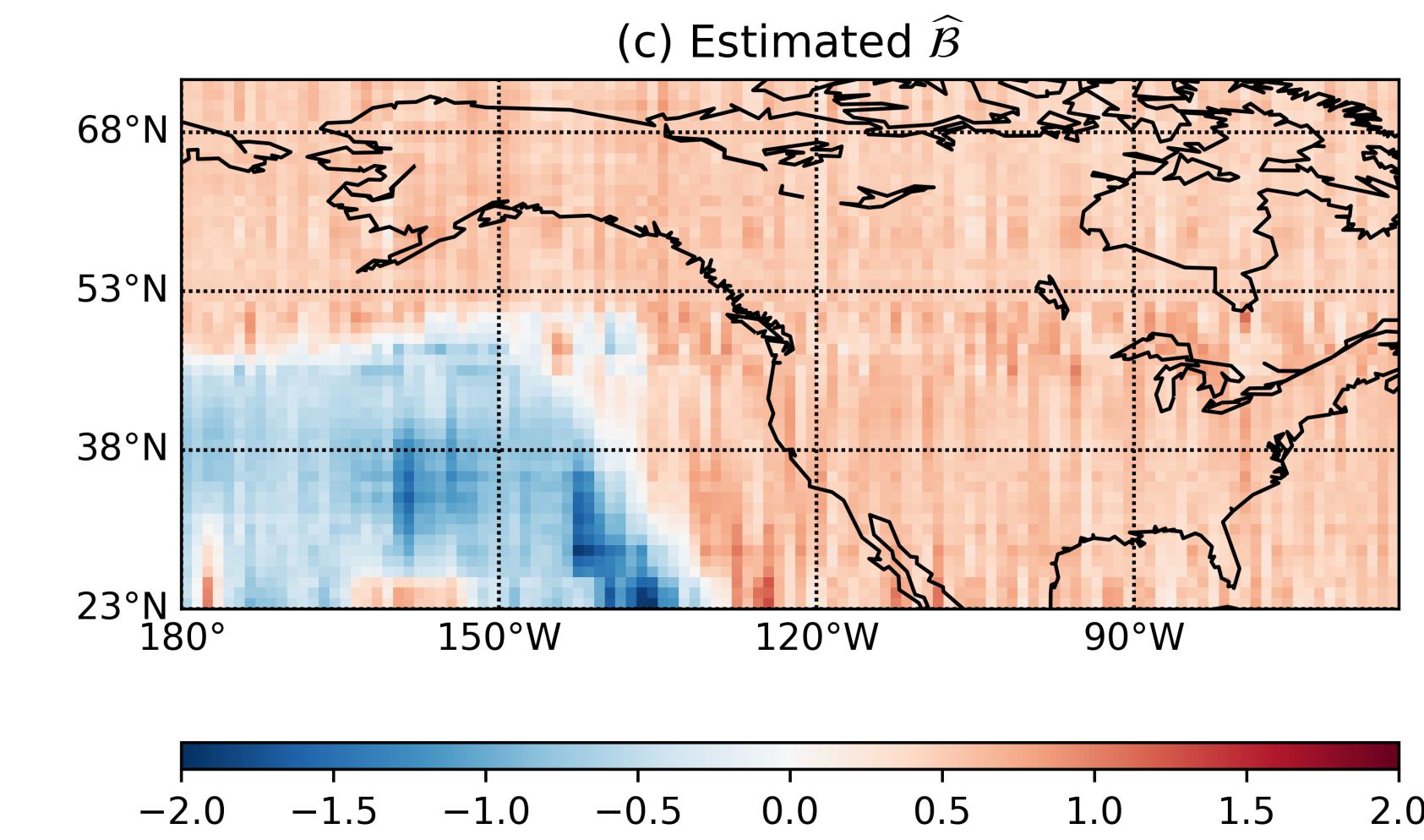
- We apply our method to reconstruct the total electron content (TEC) over the US and Canada for 15 days in Sept, 2017.



data without missingness



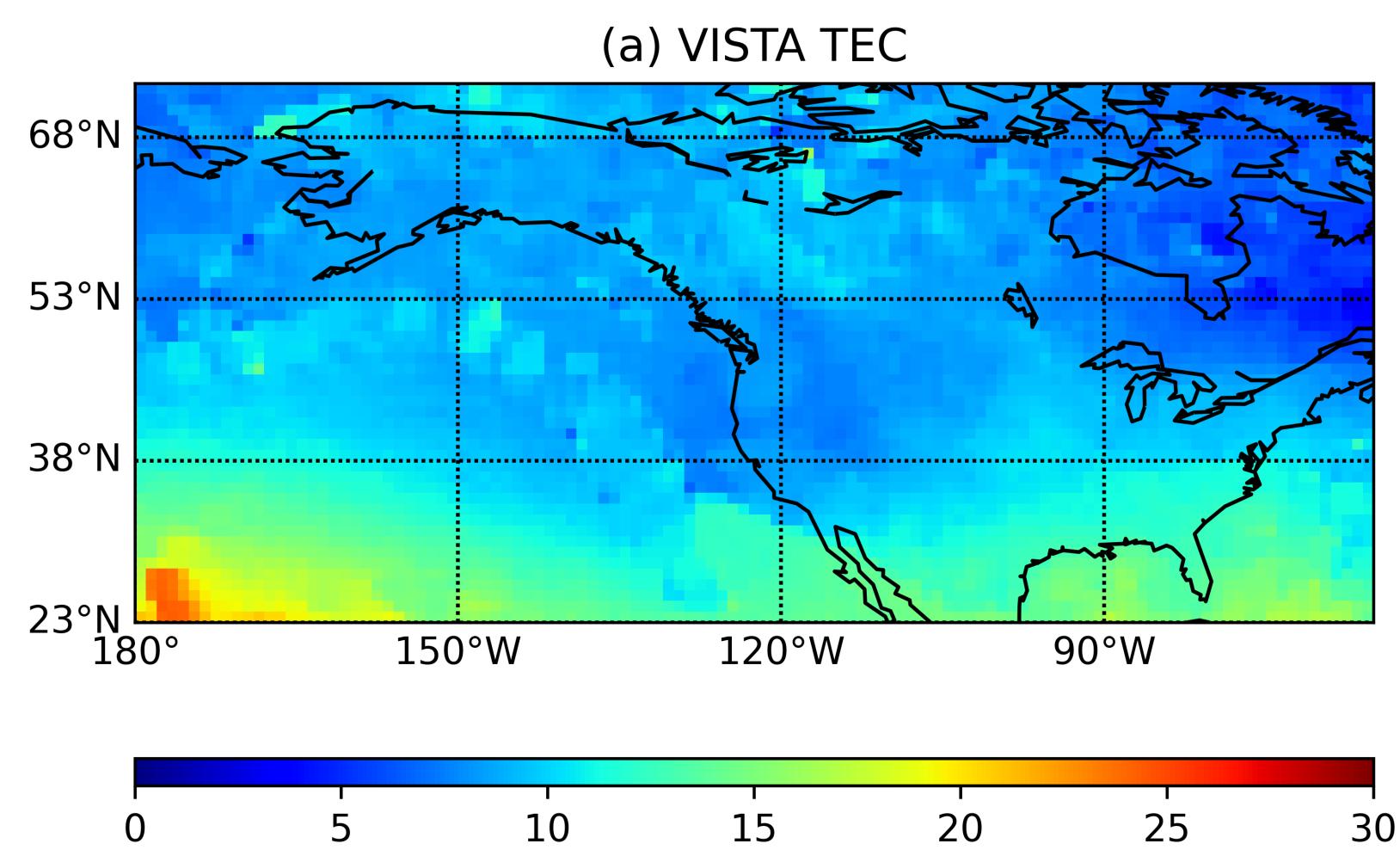
data with missingness



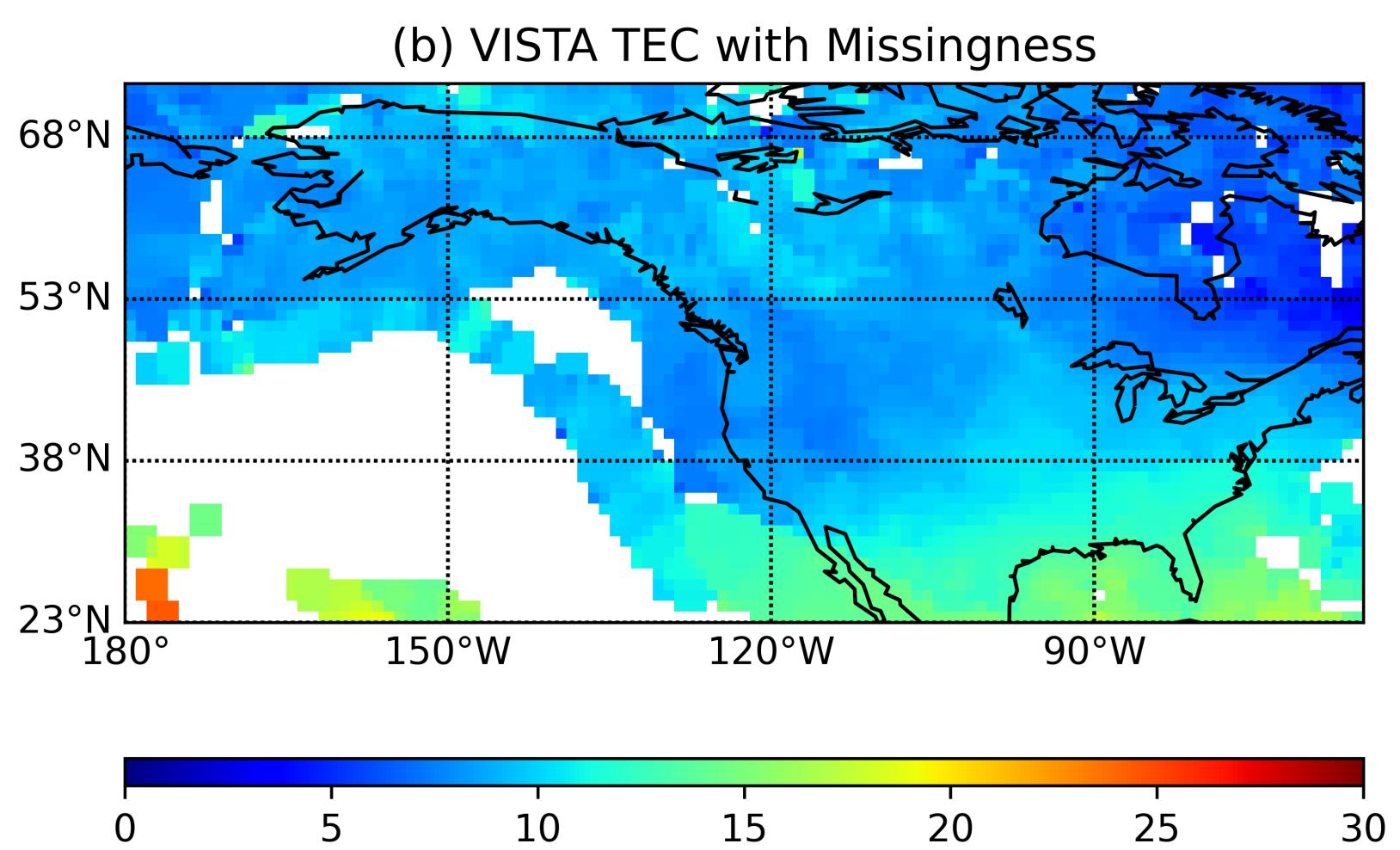
# Data Application

## Regional TEC Reconstruction

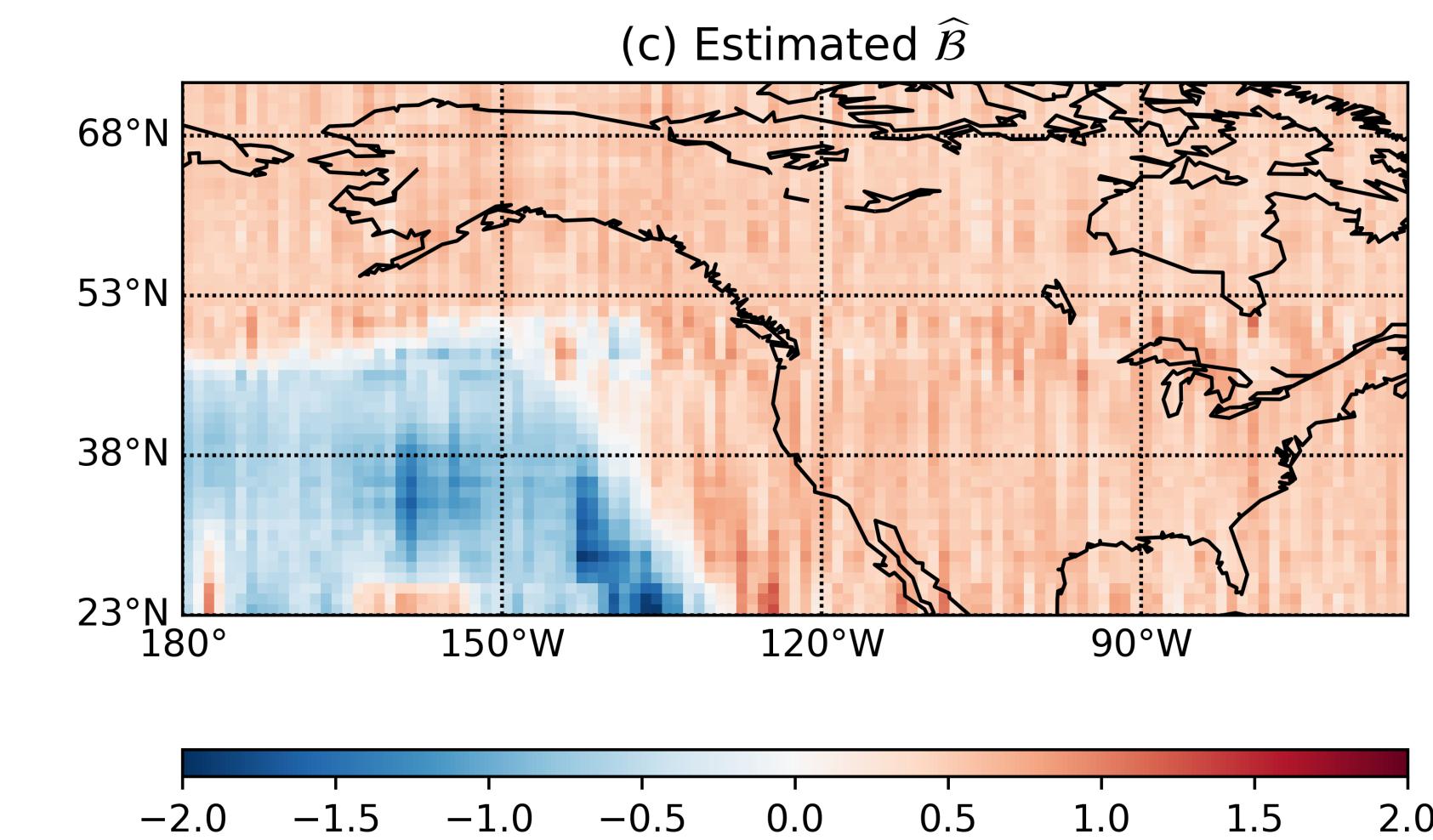
- We apply our method to reconstruct the total electron content (TEC) over the US and Canada for 15 days in Sept, 2017.



data without missingness



data with missingness



fitted  $\widehat{\mathcal{B}}$   
blue: more likely to miss

# Data Application

## Coverage Performance

- We **compare** our method under both Ising model and Bernoulli model with the unweighted conformal prediction:

# Data Application

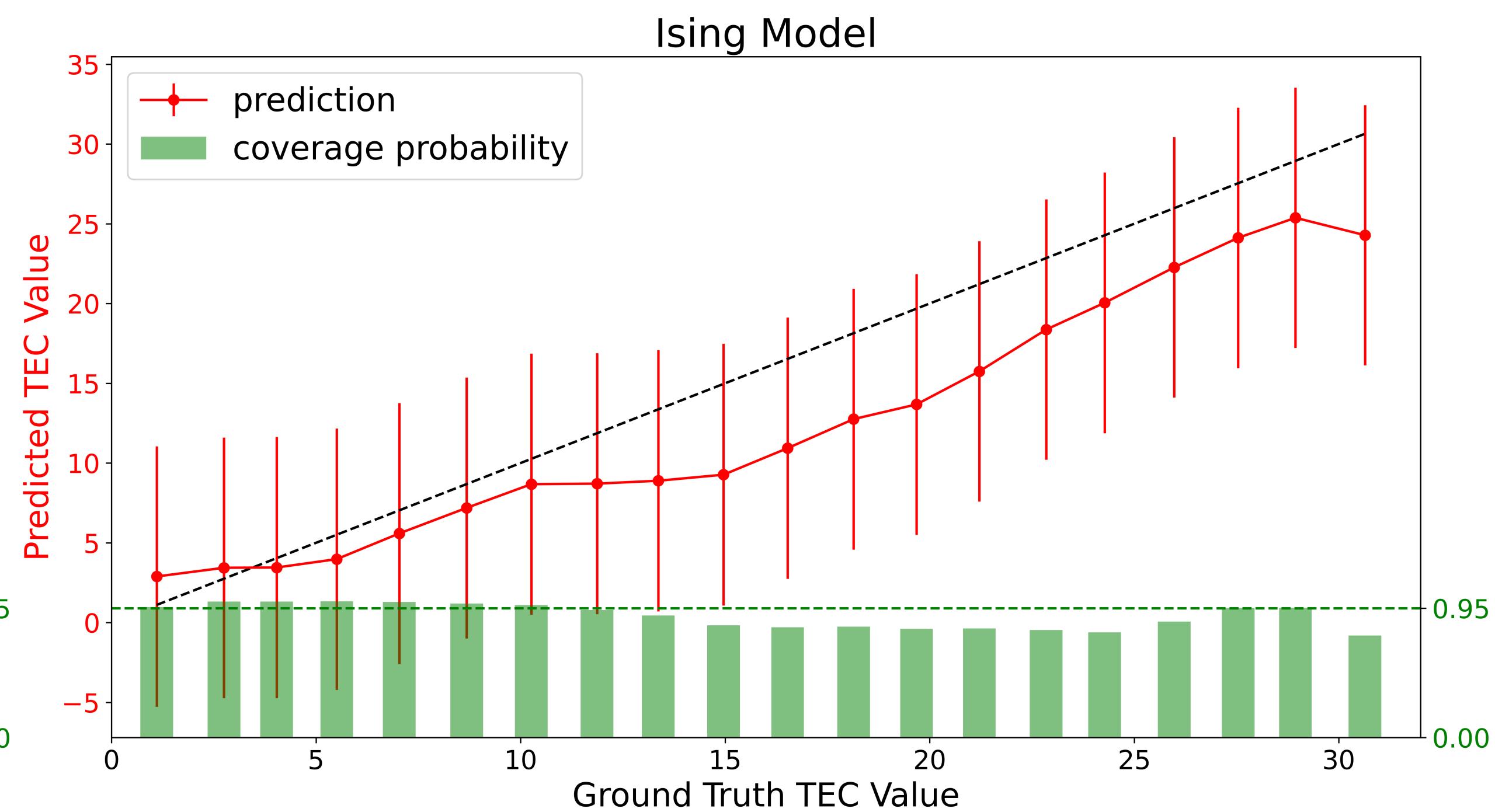
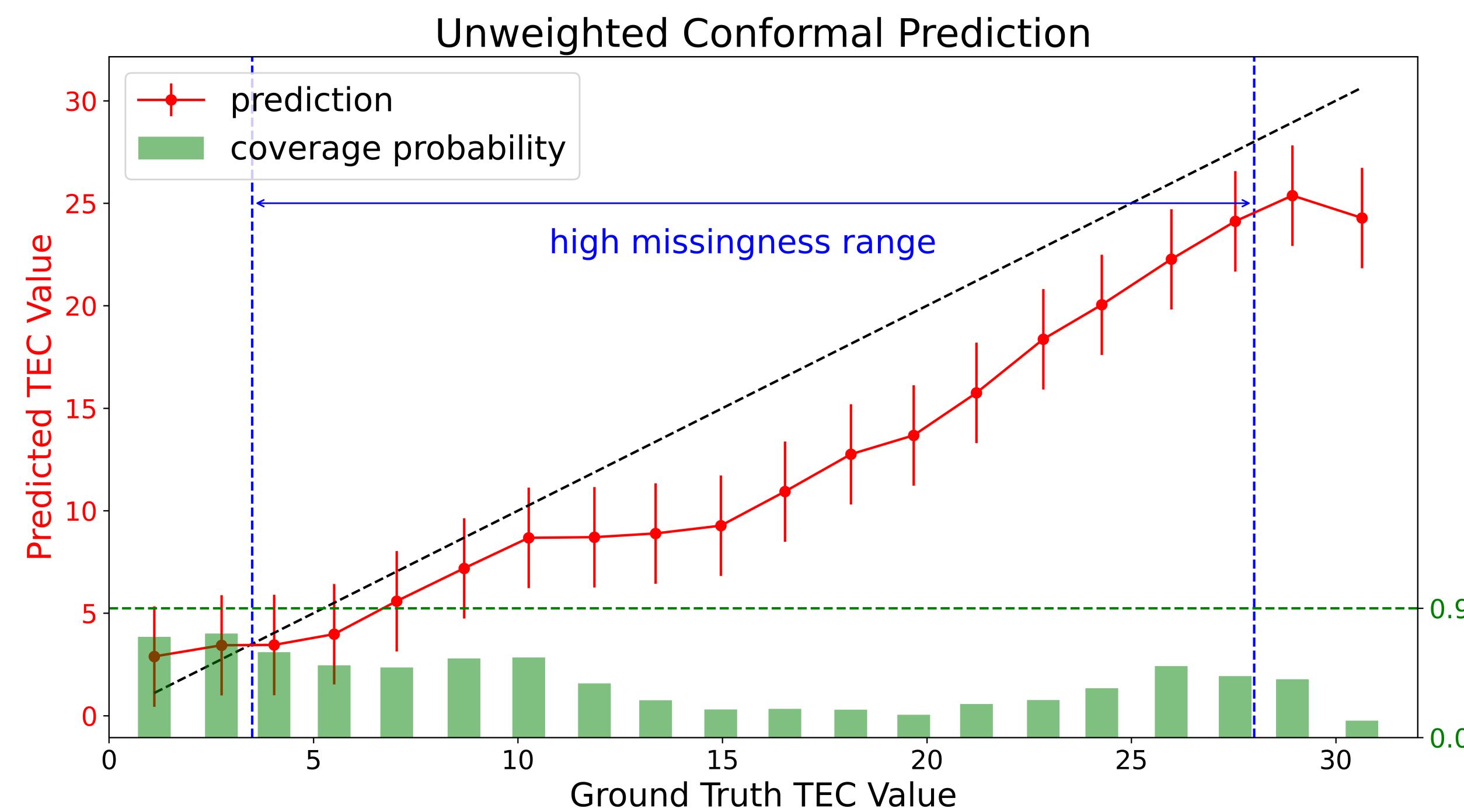
## Coverage Performance

- We **compare** our method under both Ising model and Bernoulli model with the unweighted conformal prediction:

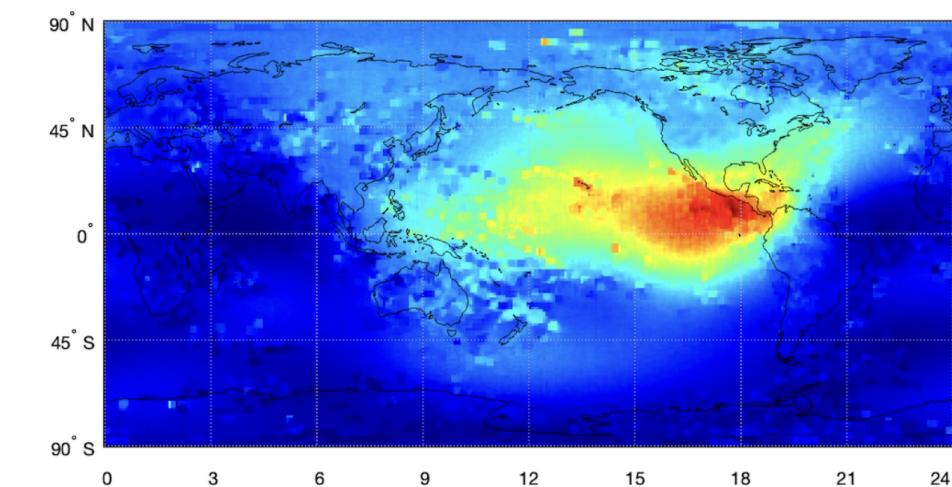
method	mis-coverage %	90% CI coverage %	95% CI coverage %
unweighted	42.1(6.49)	46.3(6.58)	52.3(7.23)
Bernoulli	23.1(5.34)	64.6(5.97)	76.8(5.03)
Ising	6.01(2.45)	90.0(6.06)	94.2(3.74)

# Data Application

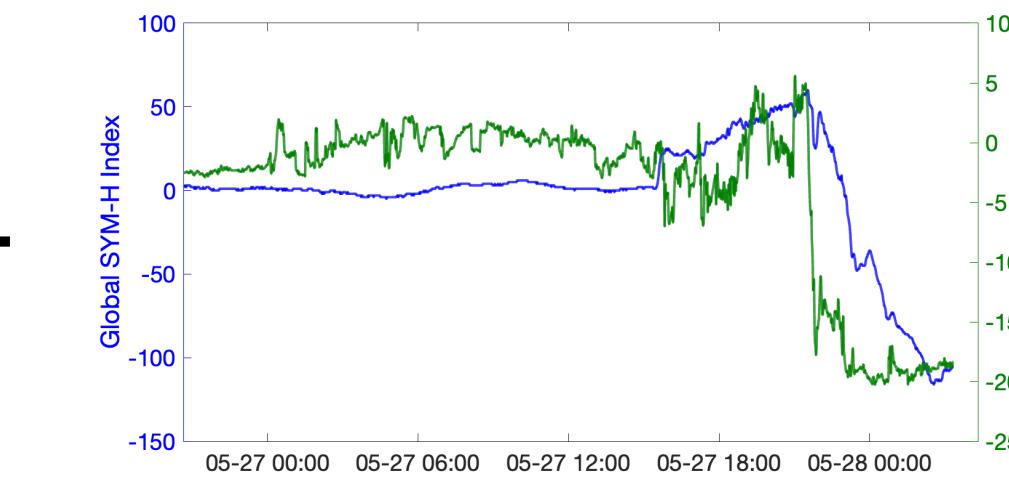
A Single Day Example (Sept 6, 2017)



# III. Matrix Autoregression with Auxiliary Data



+

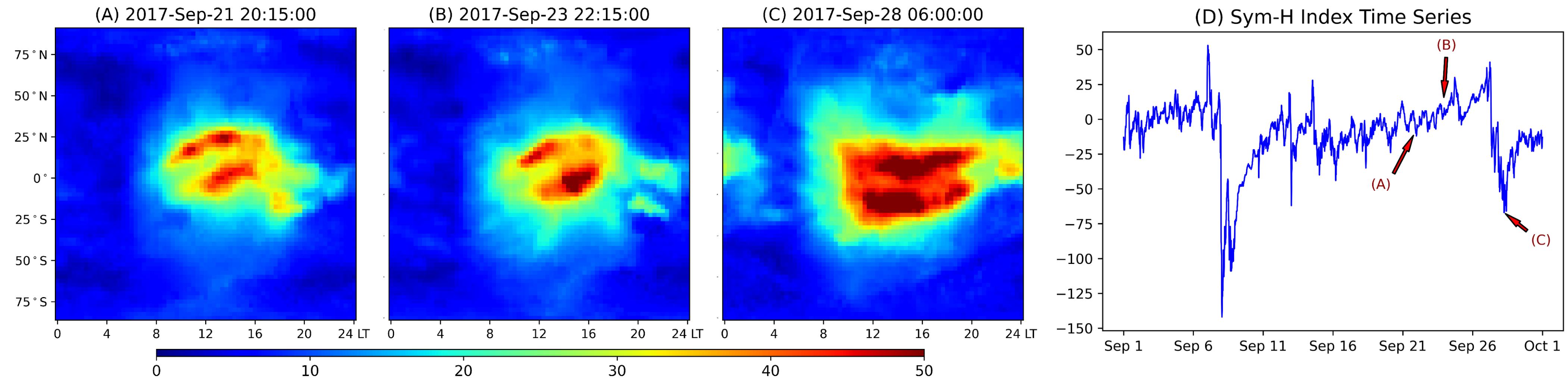


Spatial Matrix Time Series   Non-spatial Vector Time Series

How to forecast future matrix data?

# Motivating Example

## Global TEC Forecast with Solar Wind Parameters



- Left 3 panels: Global TEC, measurement of the Earth's ionospheric activity
- Right panel: Sym-H index, measurement of the impact of solar eruption on Earth

# Our Model

**Matrix AutoRegression with Auxiliary Covariates (MARAC)**

# Our Model

## Matrix AutoRegression with Auxiliary Covariates (MARAC)

- Denote  $\{\mathbf{X}_t\}_{t=1}^T$  as the matrix time series, and  $\{\mathbf{z}_t\}_{t=1}^T$  as the vector time series, then our MARAC(P,Q) model is:

# Our Model

## Matrix AutoRegression with Auxiliary Covariates (MARAC)

- Denote  $\{\mathbf{X}_t\}_{t=1}^T$  as the matrix time series, and  $\{\mathbf{z}_t\}_{t=1}^T$  as the vector time series, then our MARAC(P,Q) model is:

$$\mathbf{X}_t = \underbrace{\sum_{p=1}^P \mathbf{A}_p \mathbf{X}_{t-p} \mathbf{B}_p^\top}_{\text{autoregressive part}} + \underbrace{\sum_{q=1}^Q \mathcal{G}_q \bar{\times} \mathbf{z}_{t-q}}_{\text{auxiliary part}} + \mathbf{E}_t, \quad \text{vec}(\mathbf{E}_t) \sim \mathbf{N}(\mathbf{0}, \Sigma_c \otimes \Sigma_r)$$

# Our Model

## Matrix AutoRegression with Auxiliary Covariates (MARAC)

- Denote  $\{\mathbf{X}_t\}_{t=1}^T$  as the matrix time series, and  $\{\mathbf{z}_t\}_{t=1}^T$  as the vector time series, then our MARAC(P,Q) model is:

$$\mathbf{X}_t = \underbrace{\sum_{p=1}^P \mathbf{A}_p \mathbf{X}_{t-p} \mathbf{B}_p^\top}_{\text{autoregressive part}} + \underbrace{\sum_{q=1}^Q \mathcal{G}_q \bar{\times} \mathbf{z}_{t-q} + \mathbf{E}_t}_{\text{auxiliary part}}, \quad \text{vec}(\mathbf{E}_t) \sim \mathbf{N}(\mathbf{0}, \Sigma_c \otimes \Sigma_r)$$

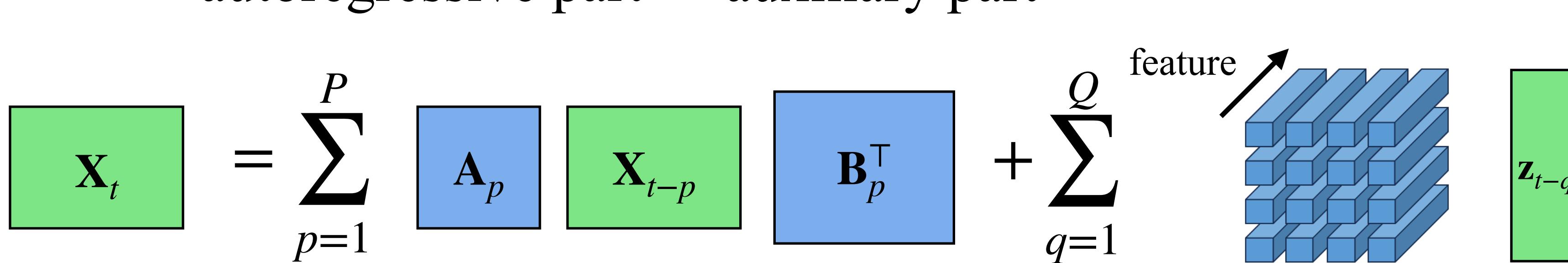
$$\boxed{\mathbf{X}_t} = \sum_{p=1}^P \boxed{\mathbf{A}_p} \boxed{\mathbf{X}_{t-p}} \boxed{\mathbf{B}_p^\top}$$

# Our Model

## Matrix AutoRegression with Auxiliary Covariates (MARAC)

- Denote  $\{\mathbf{X}_t\}_{t=1}^T$  as the matrix time series, and  $\{\mathbf{z}_t\}_{t=1}^T$  as the vector time series, then our MARAC(P,Q) model is:

$$\mathbf{X}_t = \underbrace{\sum_{p=1}^P \mathbf{A}_p \mathbf{X}_{t-p} \mathbf{B}_p^\top}_{\text{autoregressive part}} + \underbrace{\sum_{q=1}^Q \mathcal{G}_q \bar{\times} \mathbf{z}_{t-q} + \mathbf{E}_t}_{\text{auxiliary part}}, \quad \text{vec}(\mathbf{E}_t) \sim \mathbf{N}(\mathbf{0}, \Sigma_c \otimes \Sigma_r)$$

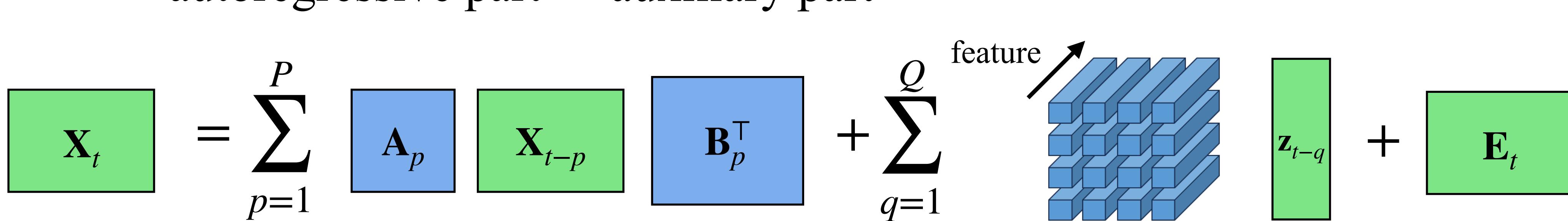


# Our Model

## Matrix AutoRegression with Auxiliary Covariates (MARAC)

- Denote  $\{\mathbf{X}_t\}_{t=1}^T$  as the matrix time series, and  $\{\mathbf{z}_t\}_{t=1}^T$  as the vector time series, then our MARAC(P,Q) model is:

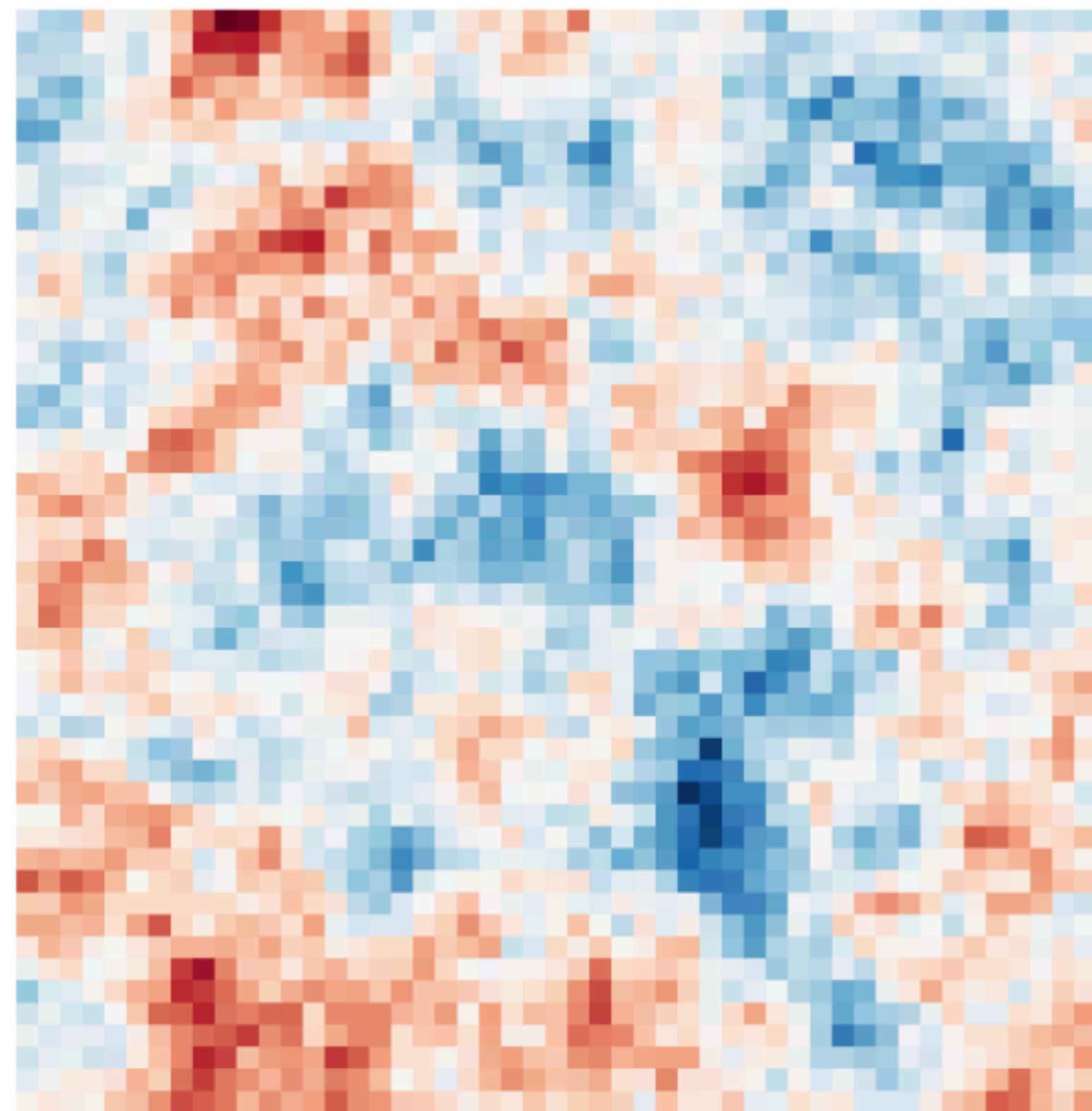
$$\mathbf{X}_t = \underbrace{\sum_{p=1}^P \mathbf{A}_p \mathbf{X}_{t-p} \mathbf{B}_p^\top}_{\text{autoregressive part}} + \underbrace{\sum_{q=1}^Q \mathcal{G}_q \bar{\times} \mathbf{z}_{t-q} + \mathbf{E}_t}_{\text{auxiliary part}}, \quad \text{vec}(\mathbf{E}_t) \sim \mathbf{N}(\mathbf{0}, \Sigma_c \otimes \Sigma_r)$$



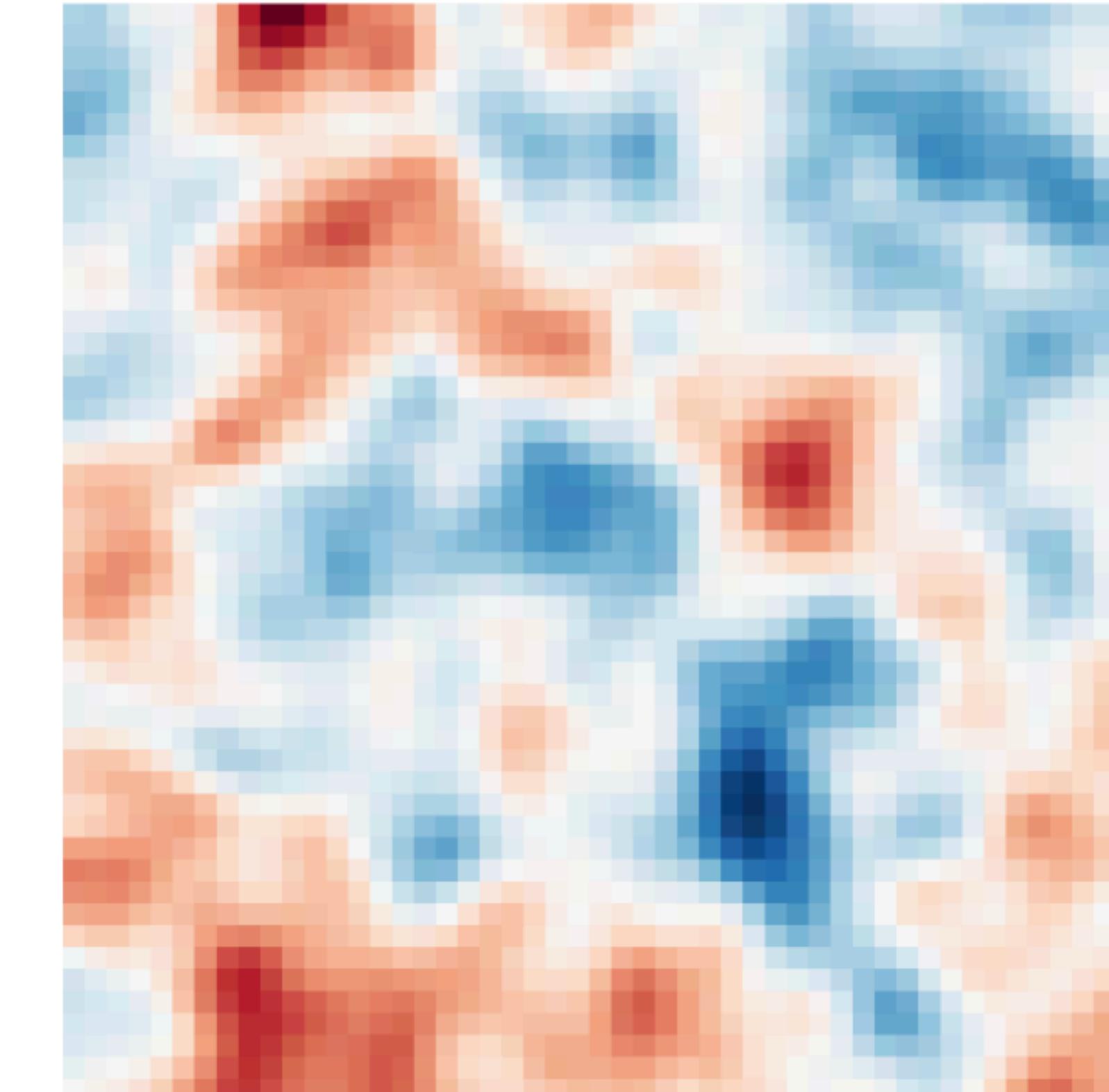
# Model

## Spatially-Smooth Regression Coefficient

Slice of Tensor Coefficient  $\mathcal{G}_q(:,:,d)$

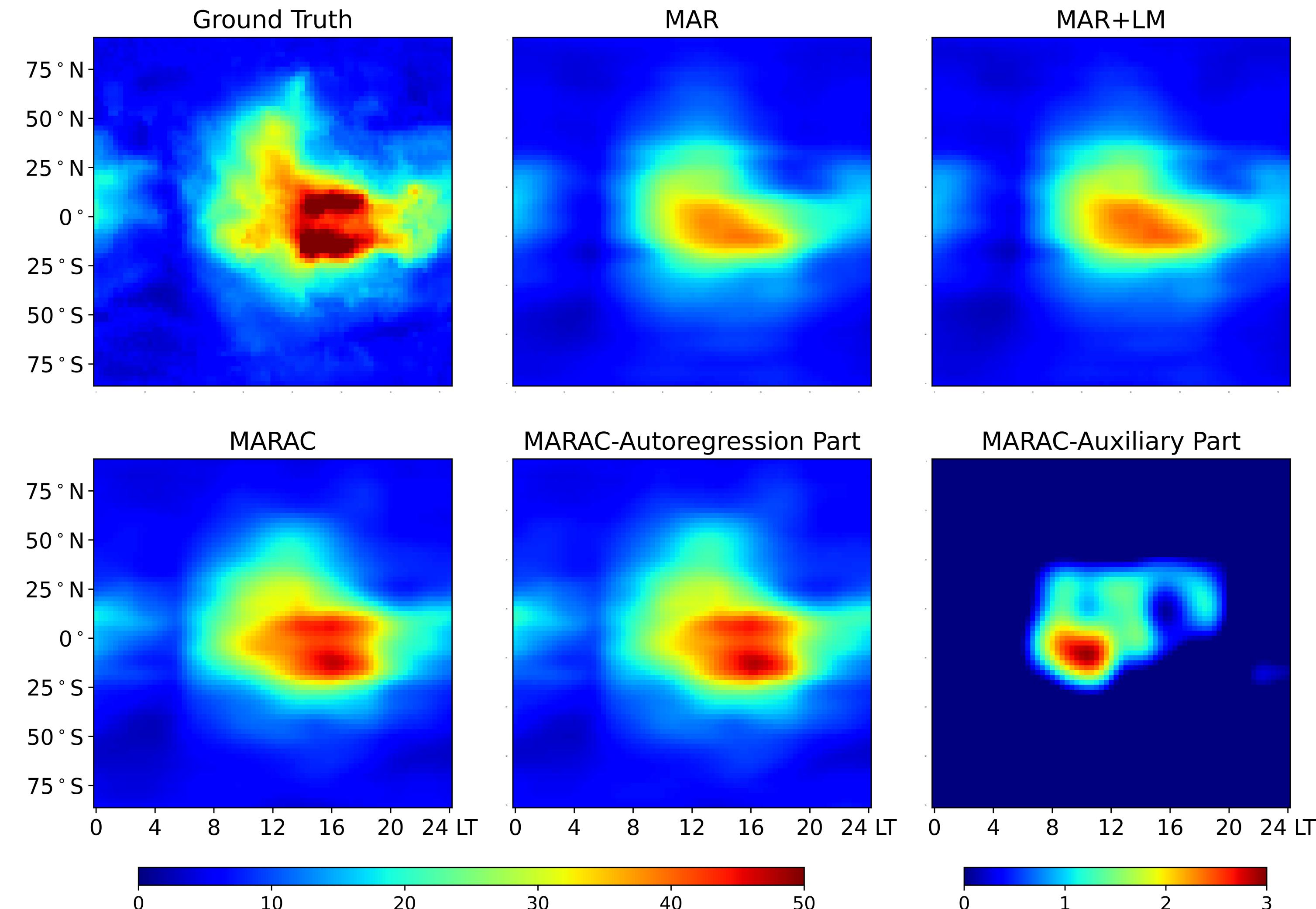


True Functional Parameter  $g_{q,d}(\cdot)$



# Real Data Application

## Global Total Electron Content (TEC) Prediction



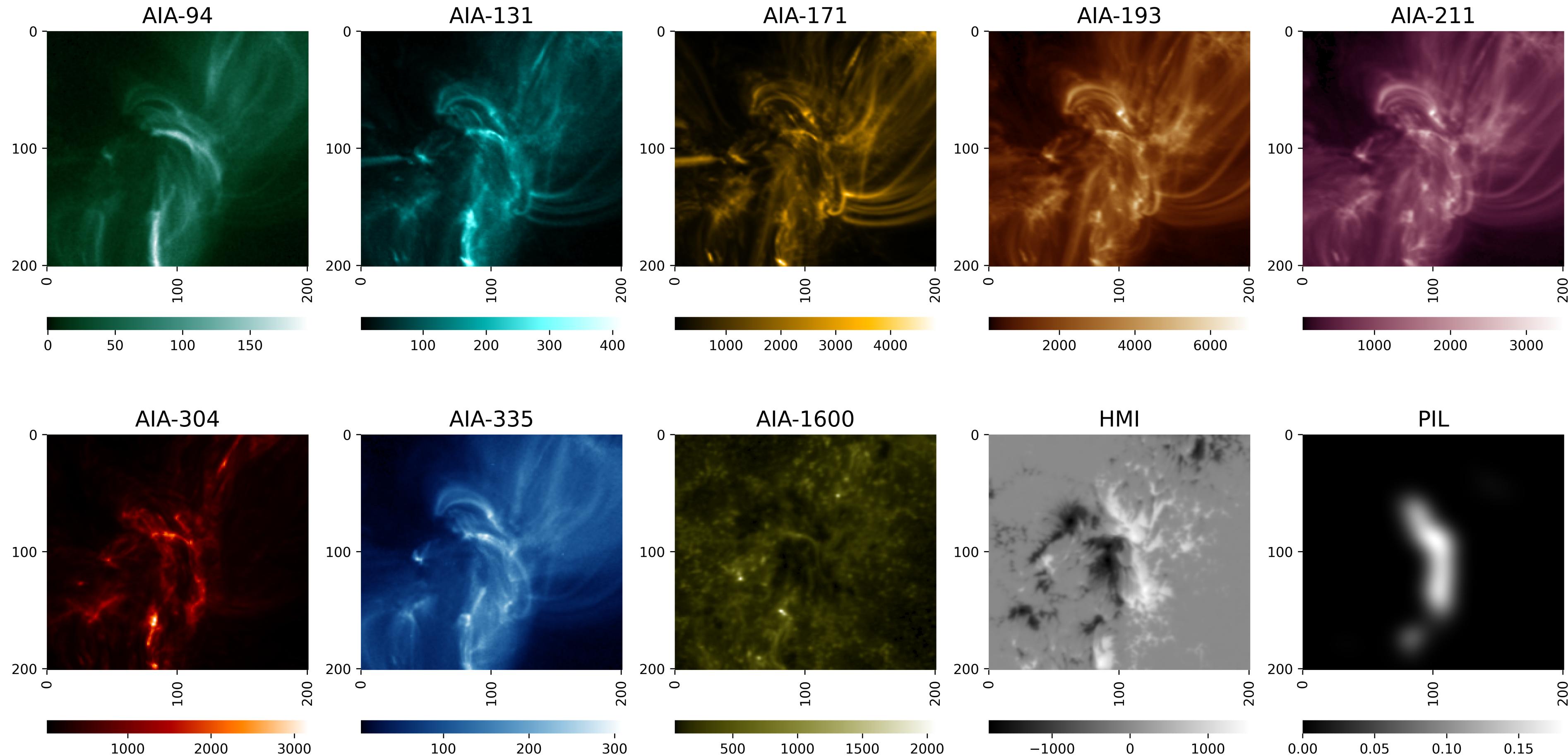
# IV. Scalar-on-Tensor Gaussian Process with Contraction

$$y = f(\begin{array}{c} \text{Modality} \\ \uparrow \\ \text{Space} \\ \xrightarrow{\hspace{1cm}} \end{array}) + \epsilon$$

**scalar label    multi-channel imaging covariate**

# Motivation

Covariate: multi-modal solar imaging data



# **Proposed Method**

## **Tensor Gaussian Process with Contraction**

# Proposed Method

## Tensor Gaussian Process with Contraction

Let  $y \in \mathbb{R}$  be the scalar label and  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$  be the tensor covariate, then we have:

# Proposed Method

## Tensor Gaussian Process with Contraction

Let  $y \in \mathbb{R}$  be the scalar label and  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$  be the tensor covariate, then we have:

$$y = g \circ h(\mathcal{X}) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

# Proposed Method

## Tensor Gaussian Process with Contraction

Let  $y \in \mathbb{R}$  be the scalar label and  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$  be the tensor covariate, then we have:

$$y = g \circ h(\mathcal{X}) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

$$\mathcal{Z}_{h \times w \times C} = h(\mathcal{X}) = \mathcal{X}_{H \times W \times C} \times_1 \mathbf{A}_{h \times H} \times_2 \mathbf{B}_{w \times W} \quad (\text{dimension reduction})$$

# Proposed Method

## Tensor Gaussian Process with Contraction

Let  $y \in \mathbb{R}$  be the scalar label and  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$  be the tensor covariate, then we have:

$$y = g \circ h(\mathcal{X}) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

$$\mathcal{Z}_{h \times w \times C} = h(\mathcal{X}) = \mathcal{X}_{H \times W \times C} \times_1 \mathbf{A}_{h \times H} \times_2 \mathbf{B}_{w \times W} \quad (\text{dimension reduction})$$

$$g(\cdot) \sim \mathbf{GP}(0, k(\cdot, \cdot)) \quad (\text{tensor Gaussian process})$$

# Proposed Method

## Tensor Gaussian Process with Contraction

Let  $y \in \mathbb{R}$  be the scalar label and  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$  be the tensor covariate, then we have:

$$y = g \circ h(\mathcal{X}) + \epsilon, \epsilon \sim N(0, \sigma^2)$$

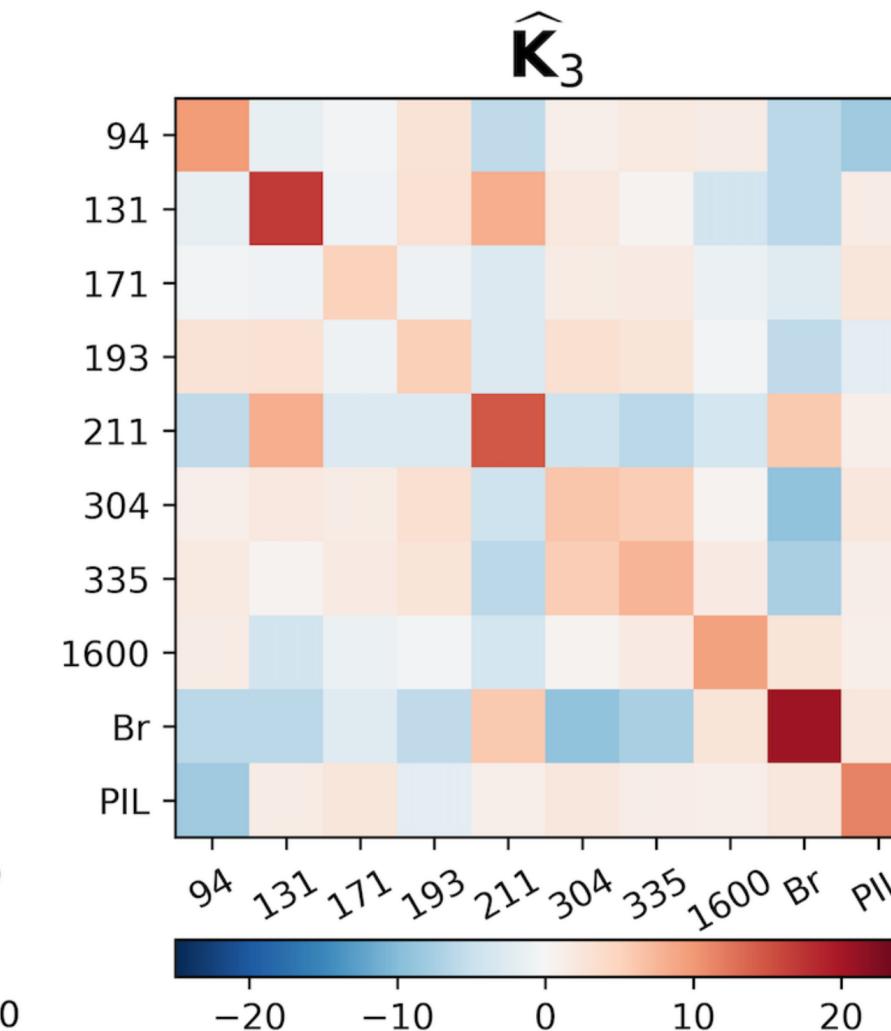
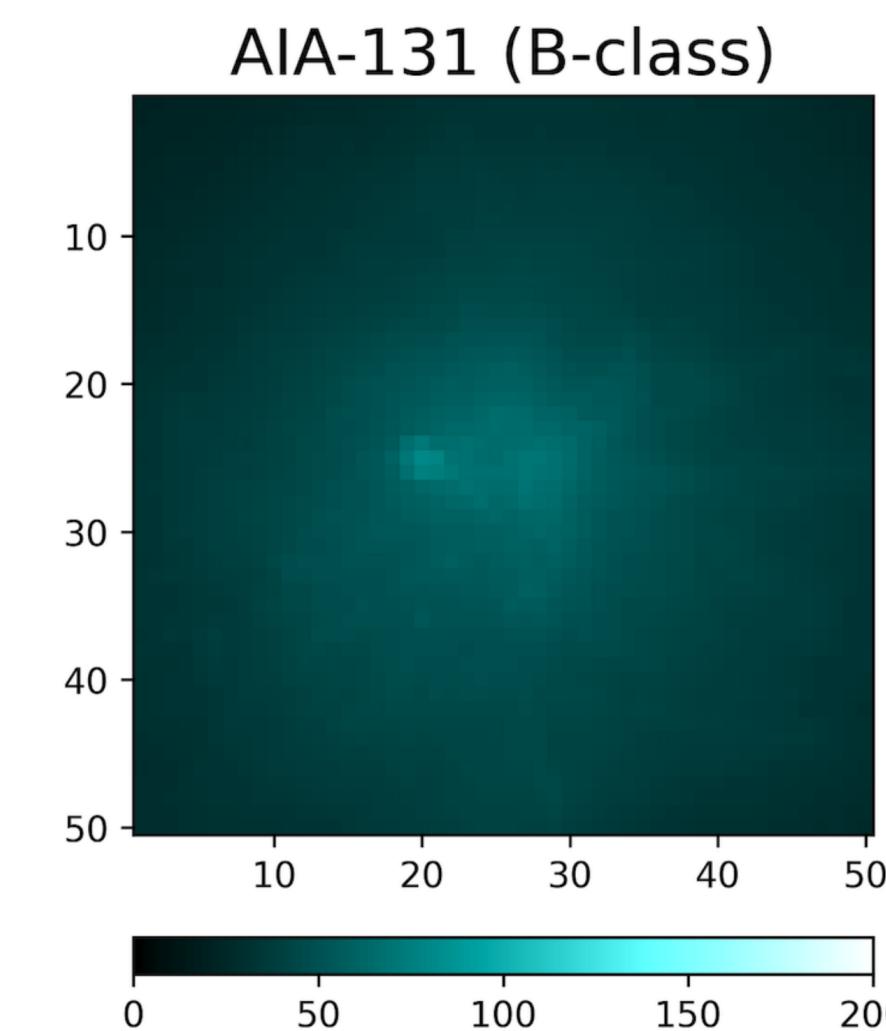
$$\mathcal{Z}_{h \times w \times C} = h(\mathcal{X}) = \mathcal{X}_{H \times W \times C} \times_1 \mathbf{A}_{h \times H} \times_2 \mathbf{B}_{w \times W} \quad (\text{dimension reduction})$$

$$g(\cdot) \sim \mathbf{GP}(0, k(\cdot, \cdot)) \quad (\text{tensor Gaussian process})$$

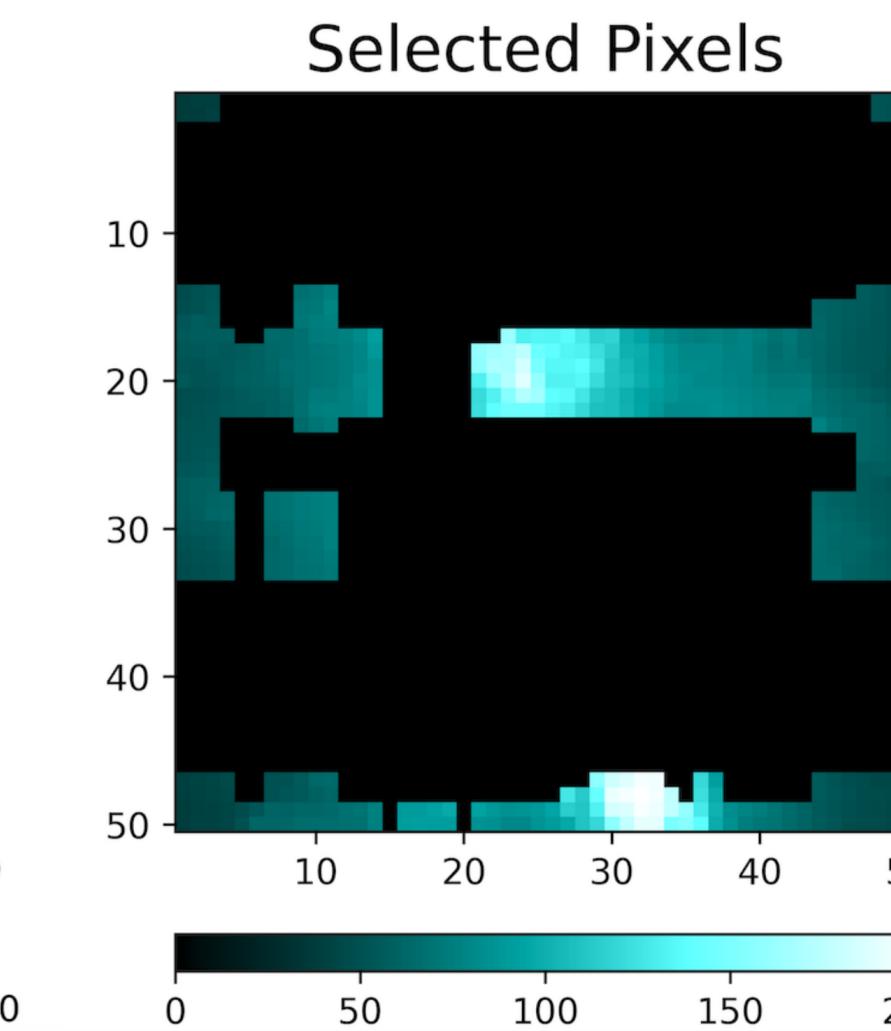
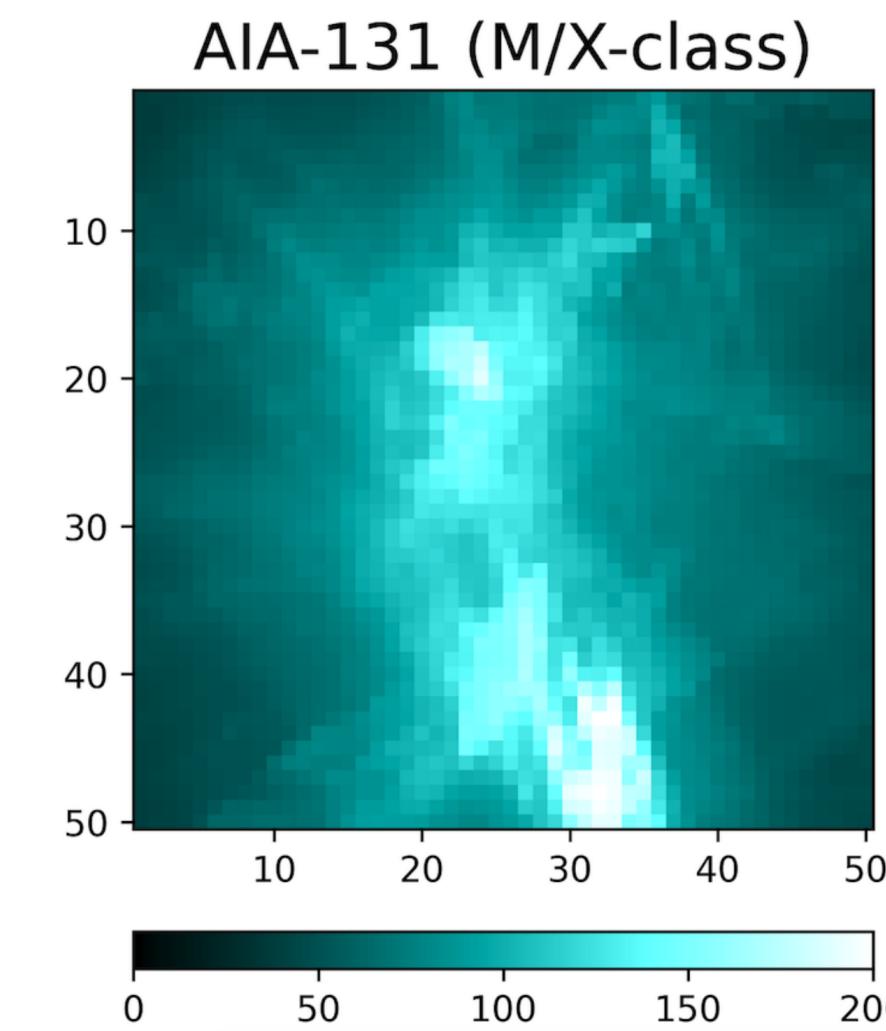
$$k(\mathcal{Z}_1, \mathcal{Z}_2) = \text{vec}(\mathcal{Z}_1)^\top \left[ \otimes_{j=1}^3 \mathbf{K}_j \right] \text{vec}(\mathcal{Z}_2)$$

# Real Data Application

## Solar Flare Intensity Forecast



Modality pairwise interaction



Feature Selection

# Final Remarks

# Summary

Roadmaps of all works

# Summary

## Roadmaps of all works

I. Tensor Completion with  
Spatio-Temporal Smoothness  
(2022, AOAS)

II. Conformalized Tensor  
Completion  
(2024, under review)

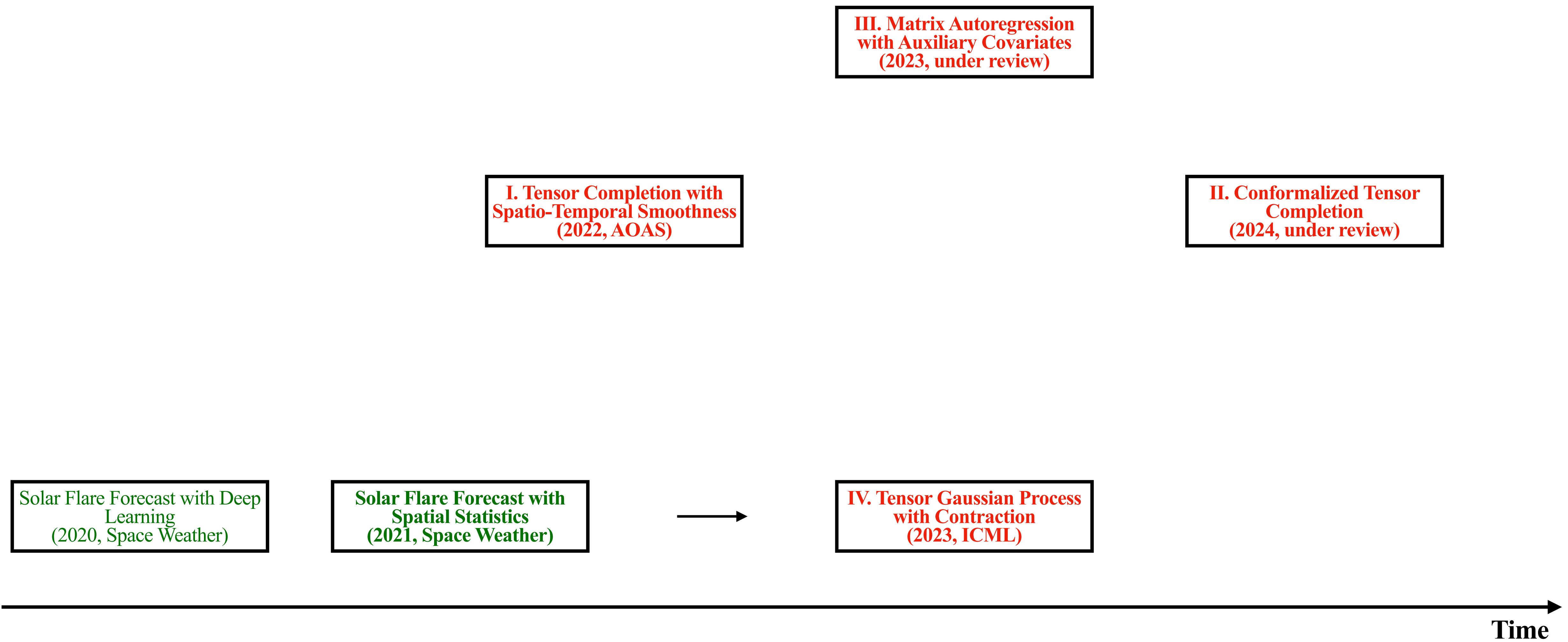
III. Matrix Autoregression  
with Auxiliary Covariates  
(2023, under review)

IV. Tensor Gaussian Process  
with Contraction  
(2023, ICML)

Time

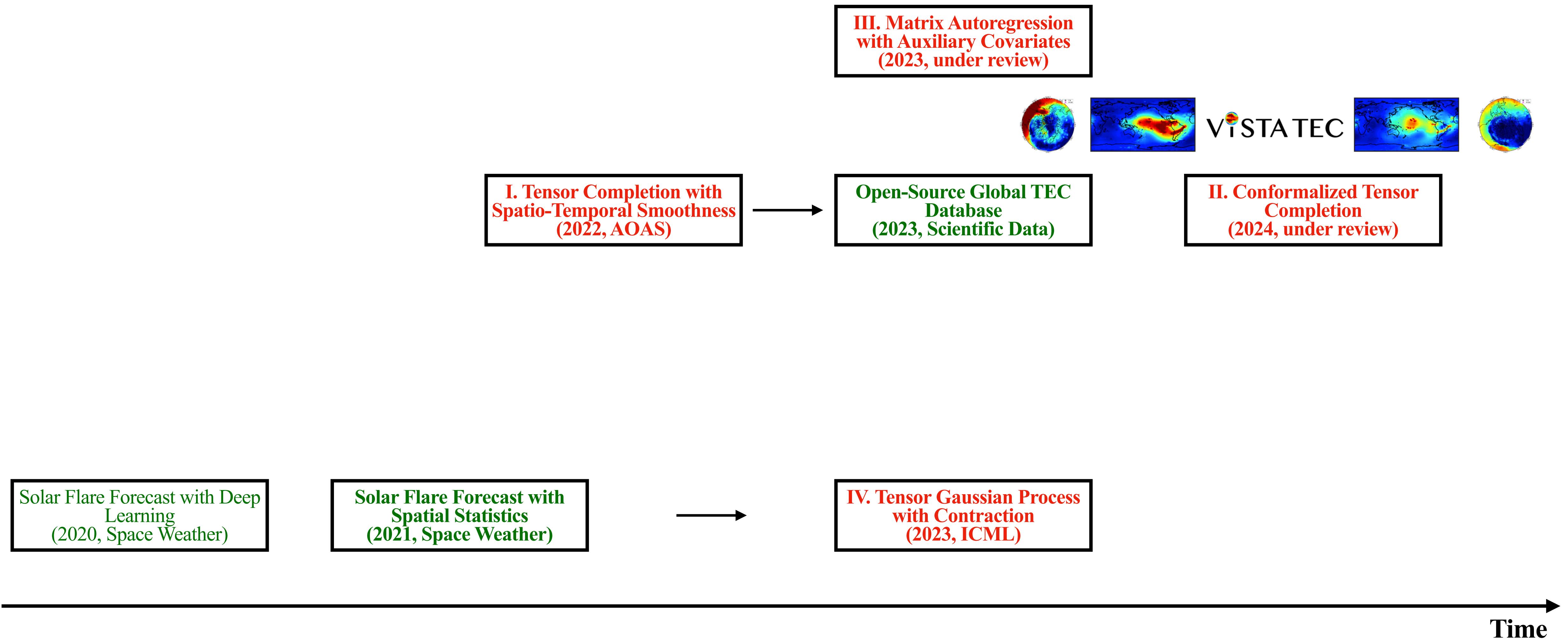
# Summary

## Roadmaps of all works



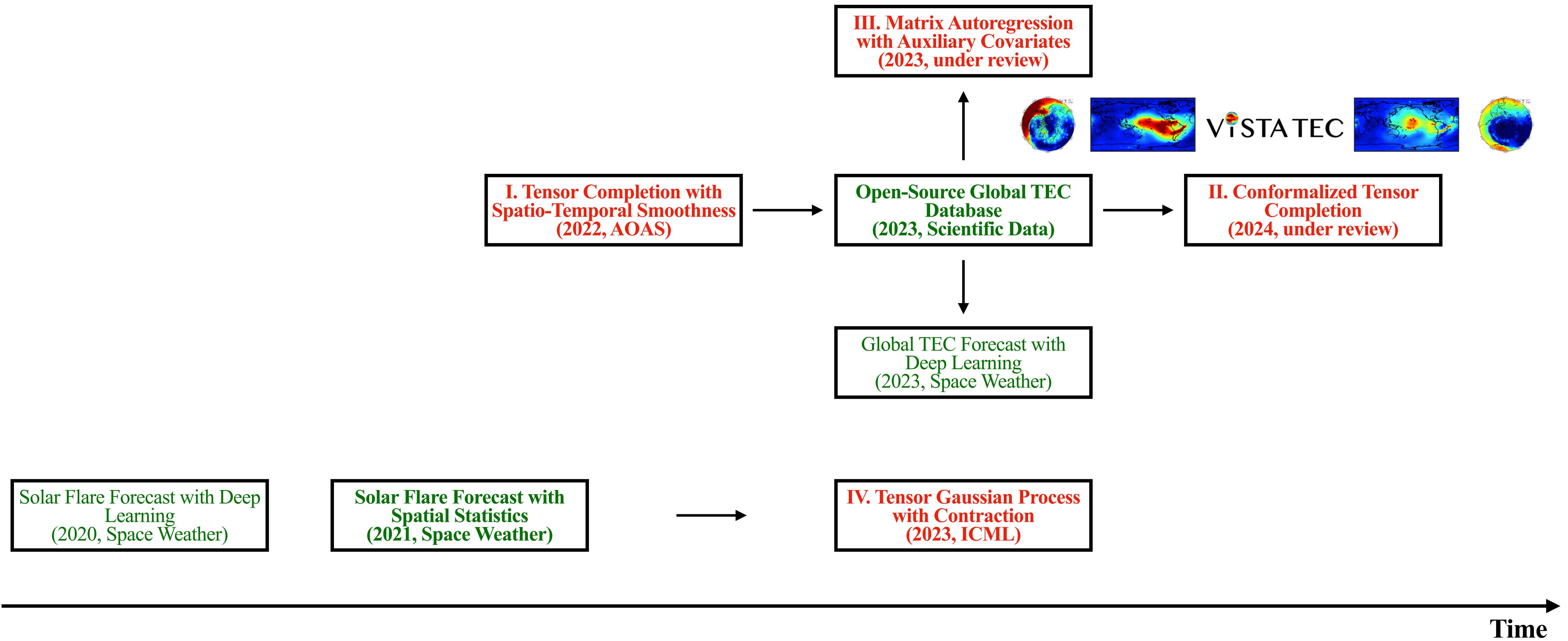
# Summary

## Roadmaps of all works



# Summary

## Roadmaps of all works



# Future Directions

# Future Directions

## Theory

# Future Directions

## Theory

- Concentration bounds of locally-dependent tensors

# Future Directions

## Theory

- Concentration bounds of locally-dependent tensors

## Methodology

# Future Directions

## Theory

- Concentration bounds of locally-dependent tensors

## Methodology

- Multi-modal tensor inference

# Future Directions

## Theory

- Concentration bounds of locally-dependent tensors

## Methodology

- Multi-modal tensor inference
- Online spatio-temporal tensor learning

# Future Directions

## Theory

- Concentration bounds of locally-dependent tensors

## Methodology

- Multi-modal tensor inference
- Online spatio-temporal tensor learning

## Computation

# Future Directions

## Theory

- Concentration bounds of locally-dependent tensors

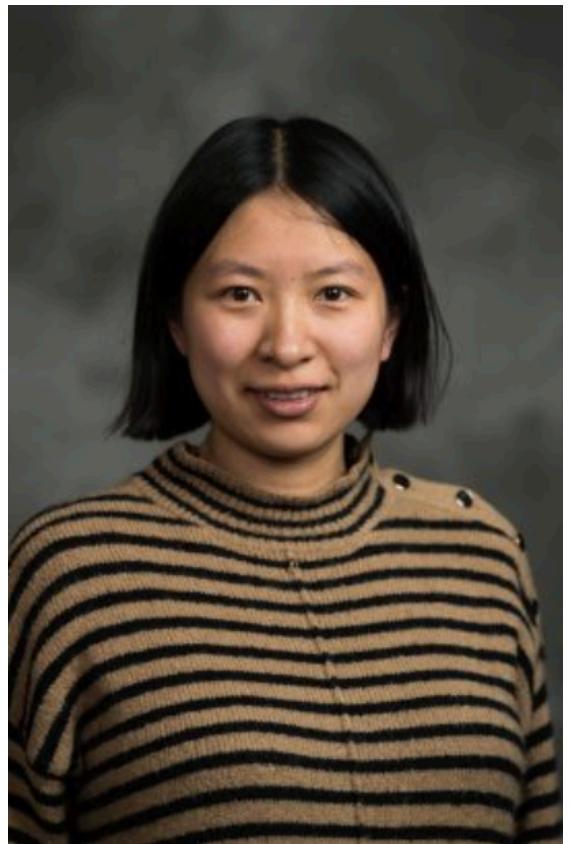
## Methodology

- Multi-modal tensor inference
- Online spatio-temporal tensor learning

## Computation

- Efficient hyper-parameter tuning

# Acknowledgements



- **Statistics:** Prof. Yang Chen, Prof. Zuofeng Shang, Prof. Yuekai Sun, Dr. Yu Wang, Zhenbang Jiao, Yurui Chang
- **EECS:** Prof. Alfred Hero, Zeyu Sun
- **CLASP:** Prof. Ward Manchester, Prof. Shasha Zou, Prof. Tamas Gombosi, Dr. Xiantong Wang, Dr. Zihan Wang, Dr. Jiaen Ren
- **Outside of UM:** Dr. Meng Jin, Dr. Yang Liu, Dr. Anthea Coster, Dr. Monica Bobra

# References I

1. Bi, X., Qu, A., & Shen, X. (2018). Multilayer tensor factorization with applications to recommender systems. *Annals of Statistics*, 46(6B), 3308-3333.
2. Wei, B., Peng, L., Guo, Y., Manatunga, A., & Stevens, J. (2023). Tensor response quantile regression with neuroimaging data. *Biometrics*, 79(3), 1947-1958.
3. Chen, Y., Fan, J., Ma, C., & Yan, Y. (2019). Inference and uncertainty quantification for noisy matrix completion. *Proceedings of the National Academy of Sciences*, 116(46), 22931-22937.
4. Gui, Y., Barber, R., & Ma, C. (2023). Conformalized matrix completion. *Advances in Neural Information Processing Systems*, 36, 4820-4844.
5. Cai, C., Poor, H. V., & Chen, Y. (2022). Uncertainty quantification for nonconvex tensor completion: Confidence intervals, heteroscedasticity and optimality. *IEEE Transactions on Information Theory*, 69(1), 407-452.
6. Oseledets, I. V. (2011). Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5), 2295-2317.
7. Holtz, S., Rohwedder, T., & Schneider, R. (2012). On manifolds of tensors of fixed TT-rank. *Numerische Mathematik*, 120(4), 701-731.
8. Hastie, T., Mazumder, R., Lee, J. D., & Zadeh, R. (2015). Matrix completion and low-rank SVD via fast alternating least squares. *The Journal of Machine Learning Research*, 16(1), 3367-3402.

# References II

10. Chen, R., Xiao, H., & Yang, D. (2021). Autoregressive models for matrix-valued time series. *Journal of Econometrics*, 222(1), 539-560.
11. Hsu, N. J., Huang, H. C., & Tsay, R. S. (2021). Matrix autoregressive spatio-temporal models. *Journal of Computational and Graphical Statistics*, 30(4), 1143-1155.
12. Li, L., & Zhang, X. (2017). Parsimonious tensor response regression. *Journal of the American Statistical Association*, 112(519), 1131-1146.
13. Carpentier, A., Eisert, J., Gross, D., & Nickl, R. (2019). Uncertainty quantification for matrix compressed sensing and quantum tomography problems. In *High Dimensional Probability VIII: The Oaxaca Volume* (pp. 385-430). Cham: Springer International Publishing.
14. Carpentier, A., Klopp, O., Löffler, M., & Nickl, R. (2018). Adaptive confidence sets for matrix completion. *Bernoulli*, 24(4A), 2429-2460.
15. Farias, V., Li, A. A., & Peng, T. (2022, May). Uncertainty quantification for low-rank matrix completion with heterogeneous and sub-exponential noise. In *International Conference on Artificial Intelligence and Statistics* (pp. 1179-1189). PMLR.
16. Xia, D., & Yuan, M. (2021). Statistical inferences of linear forms for noisy matrix completion. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 83(1), 58-77.

# References III

17. Mai, T. T., & Alquier, P. (2015). A Bayesian approach for noisy matrix completion: Optimal rate under general sampling distribution. *Electronic Journal of Statistics*, 9, 823-841.
18. Yuchi, H. S., Mak, S., & Xie, Y. (2023). Bayesian uncertainty quantification for low-rank matrix completion. *Bayesian Analysis*, 18(2), 491-518.
19. Kasalicky, P., Ledent, A., & Alves, R. (2023, September). Uncertainty-adjusted Inductive Matrix Completion with Graph Neural Networks. In *Proceedings of the 17th ACM Conference on Recommender Systems* (pp. 1169-1174).
20. Zeldes, Y., Theodorakis, S., Solodnik, E., Rotman, A., Chamiel, G., & Friedman, D. (2017). Deep density networks and uncertainty in recommender systems. *arXiv preprint arXiv:1711.02487*.
21. Shao, M., & Zhang, Y. (2023). Distribution-free matrix prediction under arbitrary missing pattern. *arXiv preprint arXiv:2305.11640*.
22. Tibshirani, R. J., Foygel Barber, R., Candes, E., & Ramdas, A. (2019). Conformal prediction under covariate shift. *Advances in neural information processing systems*, 32.
23. Cipra, B. A. (1987). An introduction to the Ising model. *The American Mathematical Monthly*, 94(10), 937-959.

# Reference IV

24. Ravikumar, P., Wainwright, M. J., & Lafferty, J. D. (2010). High-dimensional Ising model selection using  $\ell_1$ -regularized logistic regression. *Annals of Statistics*, 38(1), 1287-1319.
25. Liu, T., Mukherjee, S., & Biswas, R. (2024). Tensor recovery in high-dimensional Ising models. *Journal of Multivariate Analysis*, 203, 105335.
26. Barber, R. F., & Drton, M. (2015). High-dimensional Ising model selection with Bayesian information criteria. *Electronic Journal of Statistics*, 9, 567-607.
27. Barber, R. F., Candes, E. J., Ramdas, A., & Tibshirani, R. J. (2023). Conformal prediction beyond exchangeability. *The Annals of Statistics*, 51(2), 816-845.