

Selection of Tensor Representations For Multiview Forecasting

Student: *Nadezhda Alsahanova*
Skoltech Advisor: *Maxim Panov*

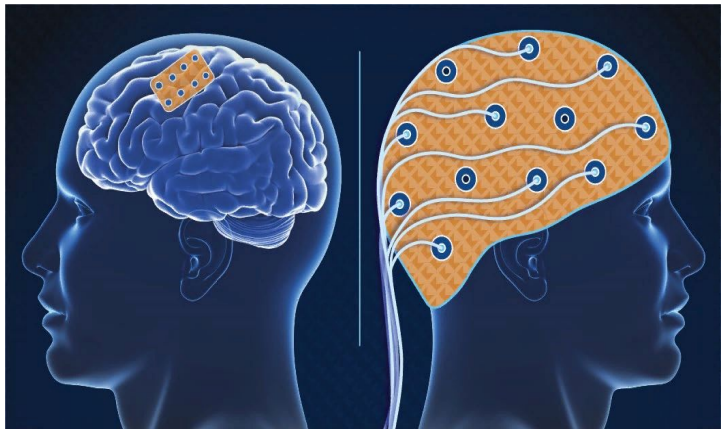
Area of research

Brain Computer Interfaces (BCI) help to restore communication and motor abilities. Data for BCI is acquired by electroencephalography (EEG) or electrocorticography (ECoG).

Problem:

Initial BCI data are redundant and highly correlated. It leads to instability of models. To solve this problem dimensionality reduction models are used. But dimensionality reduction causes loss of vital for prediction information.

Proposed solution is to use hankelization before dimensionality reduction. It helps to save information by revealing hidden correlations in time and space dimensions.



Source: GAO analysis (data). koya979/stock.adobe.com (images). | GAO-22-106118

Tensorization

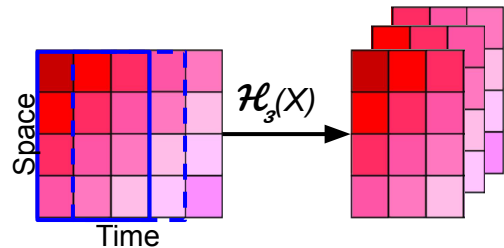
Tensorization can help to:

- find low-rank approximation with a high level of compression;
- reveal hidden correlations.

Hankelization is a natural data augmentation technique for time series to incorporate the intrinsic temporal correlation.

Hankelization connected with convolution for 1D:

$$x * h = \begin{pmatrix} x_1 & x_2 & \dots & x_k \\ x_2 & x_3 & \dots & x_{k+1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{I-k+1} & x_{I-k+2} & \dots & x_I \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ \vdots \\ h_k \end{pmatrix} = \mathcal{H}_{I-k+1}(x)h$$



Hankelization has not been used in BCI area.

Data in BCI are spatially and temporally correlated. Local correlations can be revealed by hankelization.

Aim and objectives

The thesis goal is to find optimal representations of feature and target tensors in latent space by combining hankelization and dimensionality reduction methods. These tensor representations should be optimal in terms of the forecasting quality of the target variables and the complexity of methods.

Objectives:

- determine whether hankelization along temporal mode improves quality of forecasting
- determine whether hankelization along spatial mode improves quality of forecasting
- determine whether hankelization along both modes improve quality of forecasting;

Task of multiview forecasting

Let $s(t), y(t)$ are time series, where $y(t)$ is target time series.

If there are several sources of initial time series $s(t), y(t)$, the dataset made from these time series can be presented as tensors:

$$\underline{\mathbf{X}} \in \mathbb{R}^{M \times I_1 \times \dots \times I_{D_x}}, \quad \underline{\mathbf{Y}} \in \mathbb{R}^{M \times J_1 \times \dots \times J_{D_y}}.$$

The task is to find an optimal model Φ for prediction $\underline{\mathbf{Y}}_m$ from an independent input object $\underline{\mathbf{X}}_m, m = 1, \dots, M$. The model is optimal, if it minimizes error functional \mathcal{L} :

$$\Phi^* = \arg \min_{\Theta} \mathcal{L}(\Phi(\underline{\mathbf{X}}, \Theta), \underline{\mathbf{Y}})$$

where Θ the parameters of model Φ .

Hankelization

Hankelization is an effective way to transform lower-order tensors to higher-order ones.

Hankelization of tensor $\underline{\mathbf{X}} \in \mathbb{R}^{I_1 \times \dots \times I_D}$ can be done by multi-way delay embedding transform (MDT) with use of matrices $\mathbf{S}_d, d = 1, \dots, D$:

$$\hat{\underline{\mathbf{X}}} = \mathcal{H}_\tau(\underline{\mathbf{X}}) = \text{Fold}_{(\mathbf{I}, \tau)}(\underline{\mathbf{X}} \times_1 \mathbf{S}_1 \dots \times_D \mathbf{S}_D),$$

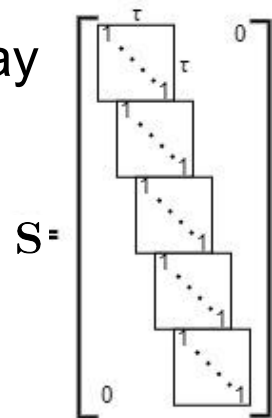
where $\hat{\underline{\mathbf{X}}} \in \mathbb{R}^{\tilde{I}_1 \times \dots \times \tilde{I}_D}, D < G$, and

$$\text{Fold}_{(\mathbf{I}, \tau)} : \mathbb{R}^{\tau_1(I_1 - \tau_1 + 1) \times \dots \times \tau_D(I_D - \tau_D + 1)} \rightarrow \mathbb{R}^{\tau_1 \times (I_1 - \tau_1 + 1) \times \dots \times \tau_D \times (I_D - \tau_D + 1)}.$$

The inverse MDT:

$$\underline{\mathbf{X}} = \mathcal{H}_\tau^{-1}(\hat{\underline{\mathbf{X}}}) = \text{Unfold}_{(\mathbf{I}, \tau)}(\hat{\underline{\mathbf{X}}}) \times_1 \mathbf{S}_1^\dagger \dots \times_N \mathbf{S}_N^\dagger,$$

where \dagger is the Moore-Penrose pseudo-inverse



Multilinear Principal Component Analysis

MPCA objective is to define a multilinear transformation that maps the original tensor space $\mathbb{R}^{I_1 \times \dots \times I_D}$ into a tensor subspace $\mathbb{R}^{L_1 \times \dots \times L_D}$ with $L_d < I_d$.

$$\hat{\underline{\mathbf{X}}}_m \approx \underline{\mathbf{X}}_m \times_1 \tilde{\mathbf{U}}^{(1)\top} \times_2 \tilde{\mathbf{U}}^{(2)\top} \dots \times_D \tilde{\mathbf{U}}^{(D)\top}, \quad m = 1, \dots, M,$$

such that $\hat{\underline{\mathbf{X}}}_m$ captures most of the variations observed in the original tensor objects. Therefore, the D projection matrices $\tilde{\mathbf{U}}^{(d)}, d = 1, \dots, D$ should maximize the total tensor scatter $\Upsilon_{\underline{\mathbf{X}}}$:

$$\{\tilde{\mathbf{U}}^{(d)} \in \mathbb{R}^{I_d \times L_d}, \quad d = 1, \dots, D\} = \arg \max_{\tilde{\mathbf{U}}^{(1)}, \dots, \tilde{\mathbf{U}}^{(D)}} \Upsilon_{\underline{\mathbf{X}}}.$$

$$\Upsilon_{\underline{\mathbf{X}}} = \sum_{m=1}^M \|\underline{\mathbf{X}}_m - \bar{\underline{\mathbf{X}}}\|_F^2$$

High-order partial least squares

HOPLS performs simultaneously constrained Tucker decompositions for an independent tensor $\underline{\mathbf{X}} \in \mathbb{R}^{M \times I_1 \times \dots \times I_D}$ and a dependent tensor $\underline{\mathbf{Y}} \in \mathbb{R}^{M \times J_1 \times \dots \times J_D}$ which have the same size in the first mode. The HOPLS model:

$$\underline{\mathbf{X}} = \underline{\mathbf{G}}_x \times_1 \mathbf{T} \times_2 \overline{\mathbf{P}}^{(1)} \dots \times_{D+1} \overline{\mathbf{P}}^{(D)} + \underline{\mathbf{E}}_R,$$

$$\underline{\mathbf{Y}} = \underline{\mathbf{G}}_y \times_1 \mathbf{T} \times_2 \overline{\mathbf{Q}}^{(1)} \dots \times_{D+1} \overline{\mathbf{Q}}^{(D)} + \underline{\mathbf{F}}_R,$$

where $\underline{\mathbf{E}}_R, \underline{\mathbf{F}}_R$ are the residuals, $\overline{\mathbf{P}}^{(d)}, \overline{\mathbf{Q}}^{(d)}$ are the mode-d loading matrices, \mathbf{T} is the latent matrix, and $\underline{\mathbf{G}}_x, \underline{\mathbf{G}}_y$ are the core tensors. The cross-covariance tensor is defined by

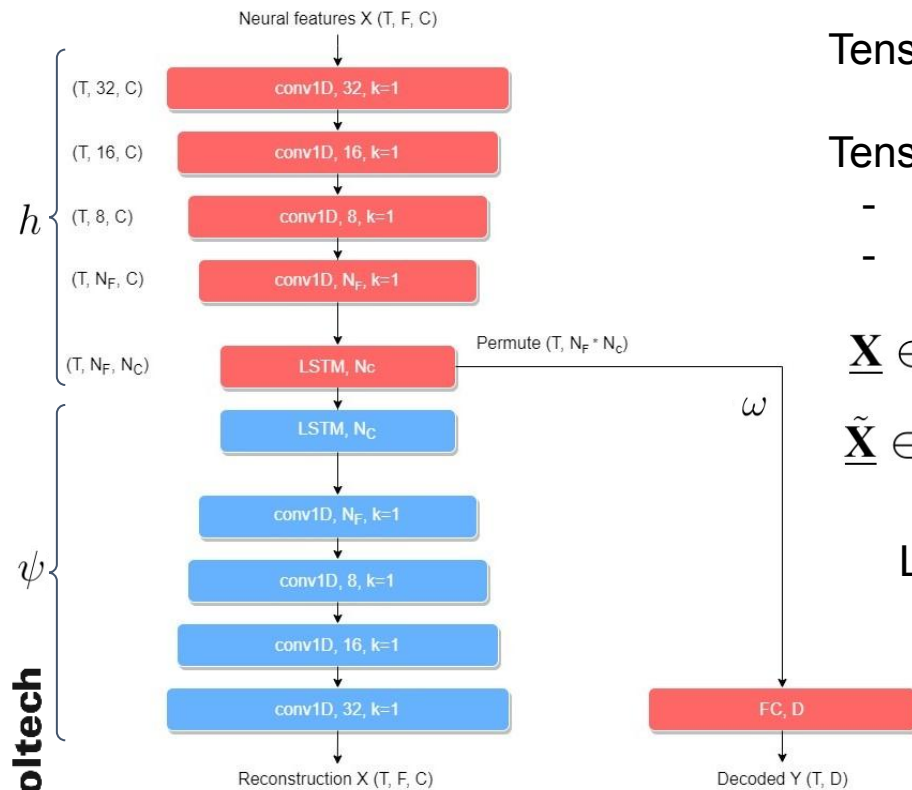
$$\underline{\mathbf{C}} = COV_{\{1,1\}}(\underline{\mathbf{X}}, \underline{\mathbf{Y}}) \in \mathbb{R}^{I_1 \times \dots \times I_D \times J_1 \times \dots \times J_D}$$

The optimization problem can be formulated as

$$\| \llbracket \underline{\mathbf{C}}; \mathbf{P}^{(1)\top}, \dots, \mathbf{P}^{(D)\top}, \mathbf{Q}^{(1)\top}, \dots, \mathbf{Q}^{(D)\top} \rrbracket_F \|^2 \rightarrow \max_{\{\mathbf{P}^{(d)}, \mathbf{Q}^{(d)}\}},$$

$$\text{s.t. } \mathbf{P}^{(d)\top} \mathbf{P}^{(d)} = \mathbf{I}_{L_d}, \mathbf{Q}^{(d)\top} \mathbf{Q}^{(d)} = \mathbf{I}_{K_d}.$$

TensorReducedNet



Example of 3D autoencoder

TensorReducedNet is a SOTA model.

TensorReducedNet helps to

- keep tensor structure of initial data;
- align features with target data.

$$\underline{\mathbf{X}} \in \mathbb{R}^{M \times I_1 \times \dots \times I_G} \xrightarrow{\text{conv1D blocks}} \tilde{\underline{\mathbf{X}}} \in \mathbb{R}^{M \times I_1 \times L_2 \times \dots \times L_G}$$

$$\tilde{\underline{\mathbf{X}}} \in \mathbb{R}^{M \times I_1 \times L_2 \times \dots \times L_G} \xrightarrow{\text{LSTM block}} \hat{\underline{\mathbf{X}}} \in \mathbb{R}^{M \times L_1 \times L_2 \times \dots \times L_G}$$

Loss: $\mathcal{L}_1 = \mathcal{L}_{rec} + \alpha \cdot \mathcal{L}_{dec},$

$$\mathcal{L}_{rec} = \frac{1}{M} \sum_{m=1}^M \|\underline{\mathbf{X}}_m - h \circ \psi(\underline{\mathbf{X}}_m)\|^2$$

$$\mathcal{L}_{dec} = \frac{1}{M} \sum_{m=1}^M \|\mathbf{Y}_m - \omega \circ h(\underline{\mathbf{X}}_m)\|^2$$

Tensor Regression

We are trying to avoid matricization on every step of decoding target variables. So, tensor regression can be defined as:

$$\mathbf{y}_m = \langle \underline{\mathbf{X}}_m | \underline{\mathbf{W}} \rangle + \varepsilon$$

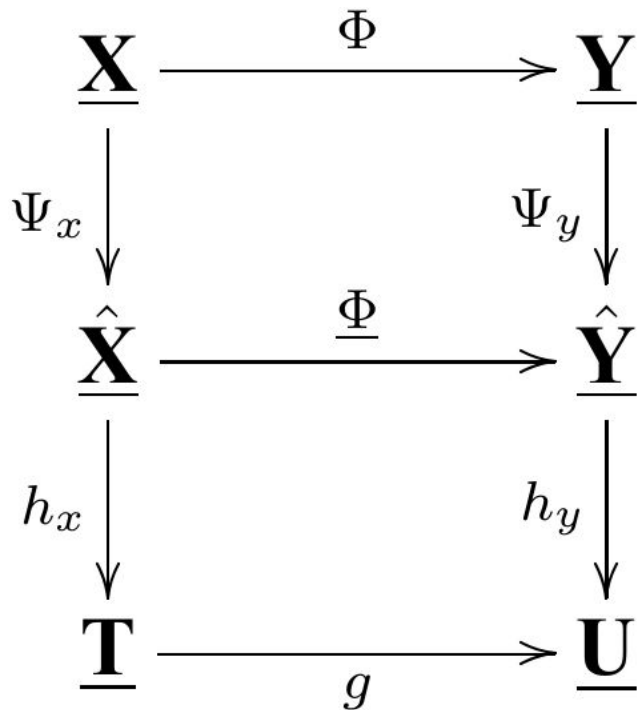
Where $\langle \underline{\mathbf{X}}_m | \underline{\mathbf{W}} \rangle$ denotes a tensor contraction along the first D modes:

$$\langle \underline{\mathbf{X}}_m | \underline{\mathbf{W}} \rangle_j = \sum_{i_1=1}^{I_1} \cdots \sum_{i_D=1}^{I_D} x_{i_1, \dots, i_D} w_{i_1, \dots, i_D, k}$$

In practice, for very large scale problems, tensors are expressed approximately in tensor network formats. For example, with the application of Tucker multilinear rank tensor representation:

$$\underline{\mathbf{W}} \approx \underline{\mathbf{G}} \times_1 \mathbf{U}^{(1)} \cdots \times_D \mathbf{U}^{(D)} \times_{D+1} \mathbf{U}^{(D+1)}$$

Algorithm



Ψ_x, Ψ_y are tensorization methods:

- without tensorization
- hankelization along time
- hankelization along space
- hankelization along both dimensions

h_x, h_y are dimensionality reduction models

g is regression model in latent space

Φ can be presented by PLS, HOPLS

$\underline{\mathbf{T}}, \underline{\mathbf{U}}$ are latent tensors

$$\Phi = \Psi_x \circ h_x \circ g \circ h_y^{-1} \circ \Psi_y^{-1}$$

NeuroTycho food-tracking dataset

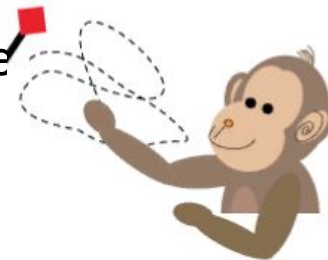
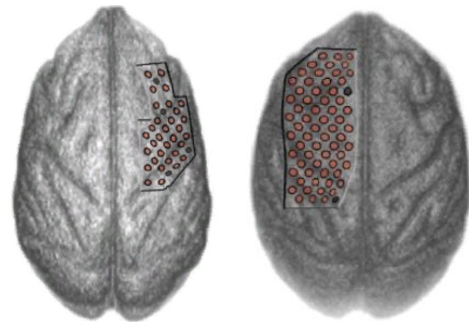
ECoG signals were obtained from 32 channels. Moreover, frequency-domain features were obtained with wavelet transform with 27 frequencies.

$$\underline{\mathbf{X}} \in \mathbb{R}^{T \times 32 \times 27}, \underline{\mathbf{Y}} \in \mathbb{R}^{T \times 3}.$$

After hankelization along temporal and spatial modes:

$$\hat{\underline{\mathbf{X}}} \in \mathbb{R}^{\hat{T} \times 10 \times 27 \times 31 \times 2}, \hat{\underline{\mathbf{Y}}} \in \mathbb{R}^{\hat{T} \times 10 \times 3}.$$

Additionally, we considered the data from other subjects with 64 channels ECoG.



Quality criteria

- The scaled Root Mean Squared Error (sRMSE):

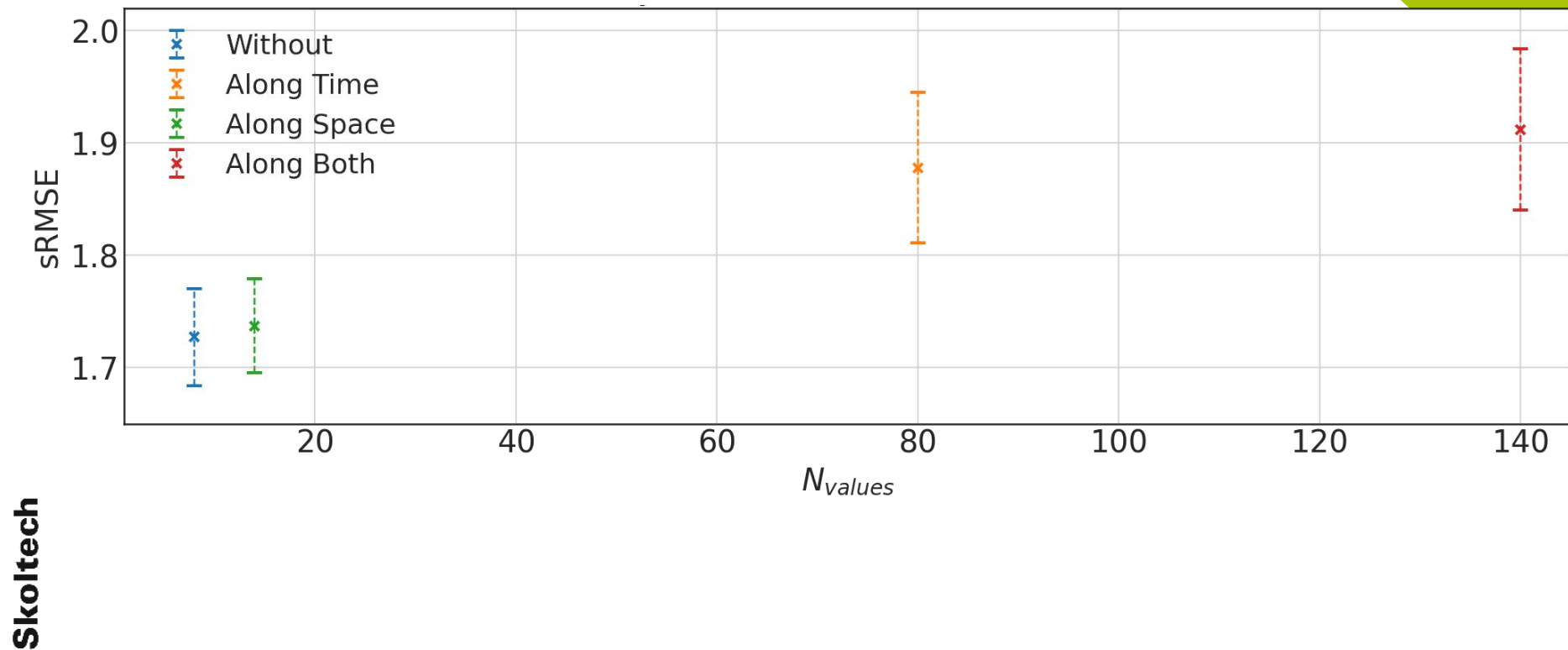
$$\text{sRMSE}(\underline{\mathbf{Y}}, \hat{\underline{\mathbf{Y}}}_{\mathbf{a}}) = \sqrt{\frac{\text{MSE}(\underline{\mathbf{Y}}, \hat{\underline{\mathbf{Y}}}_{\mathbf{a}})}{\text{MSE}(\underline{\mathbf{Y}}, \underline{\overline{\mathbf{Y}}})}} = \frac{\|\underline{\mathbf{Y}} - \hat{\underline{\mathbf{Y}}}_{\mathbf{a}}\|_2}{\|\underline{\mathbf{Y}} - \underline{\overline{\mathbf{Y}}}\|_2},$$

where $\hat{\underline{\mathbf{Y}}}_{\mathbf{a}} = \Phi(\underline{\mathbf{X}}, \underline{\Theta})$ is a model prediction and $\underline{\overline{\mathbf{Y}}}$ is an average prediction.

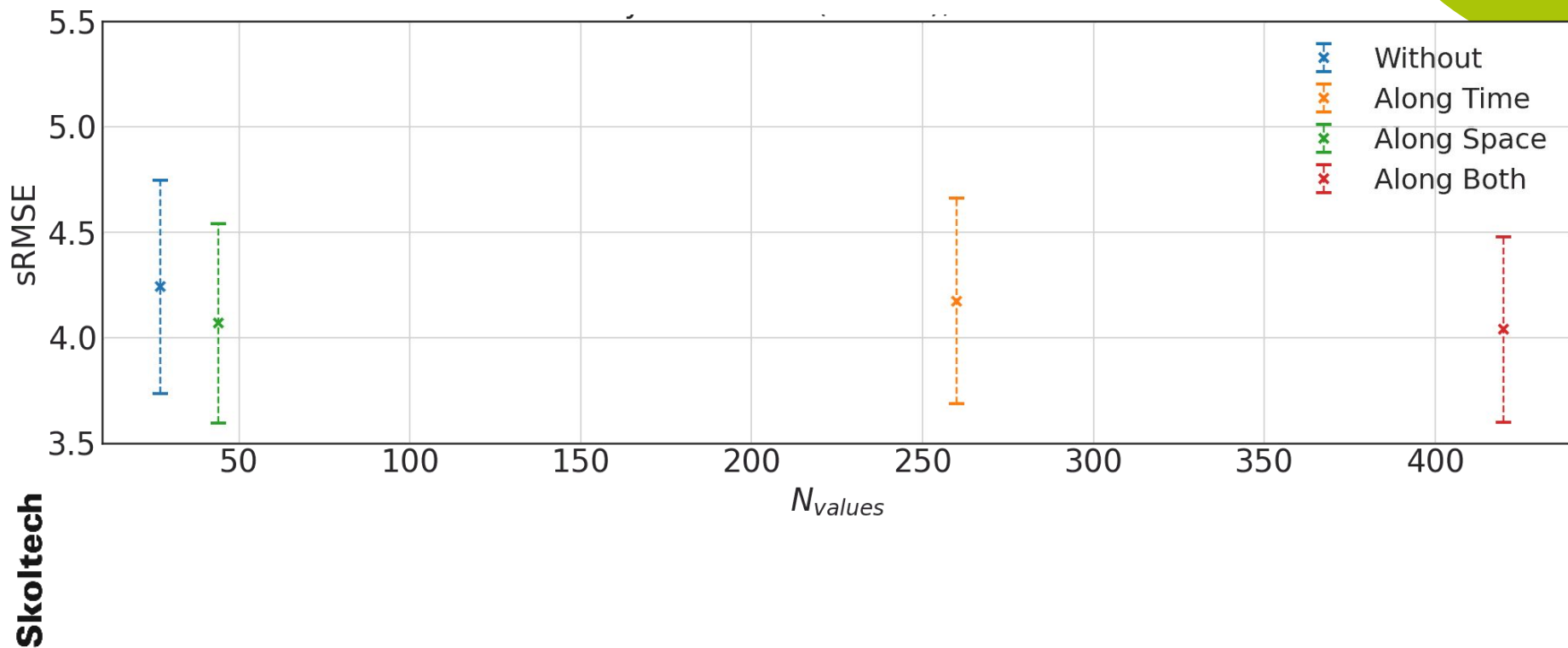
- Number of values in the latent representation of $\underline{\mathbf{X}}, \underline{\mathbf{T}} \in \mathbb{R}^{M \times L_1 \times \dots \times L_{G_x}} :$

$$N_{\text{values}} = L_1 \cdot \dots \cdot L_{G_x}$$

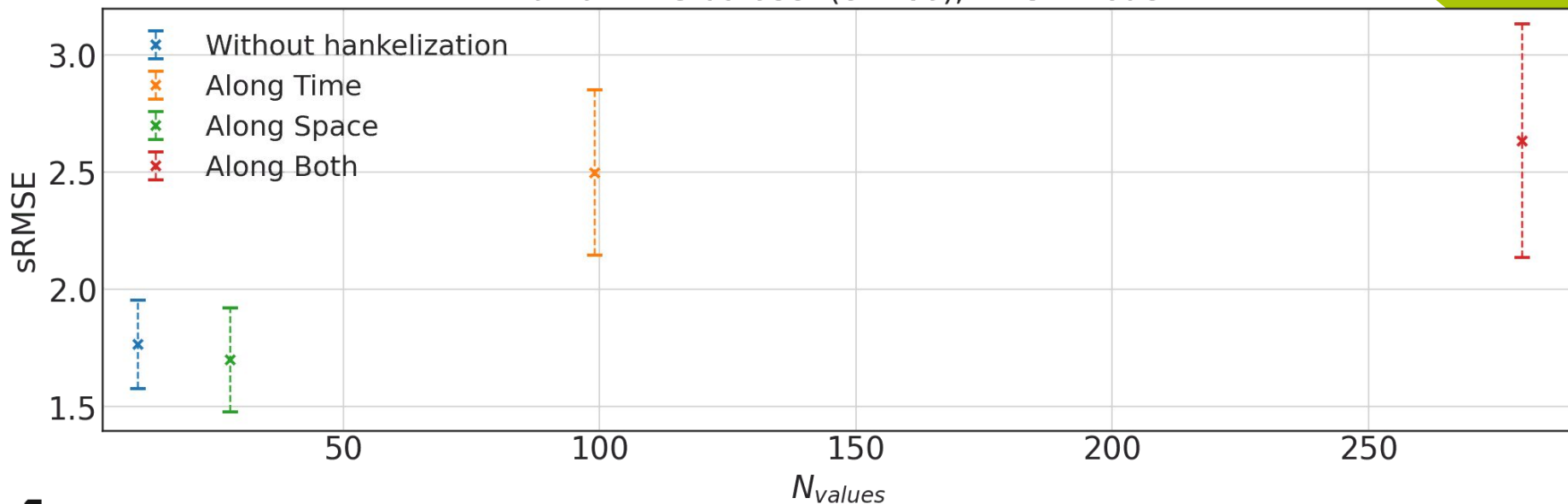
Results for NeuroTycho (ch=32) with MPCA



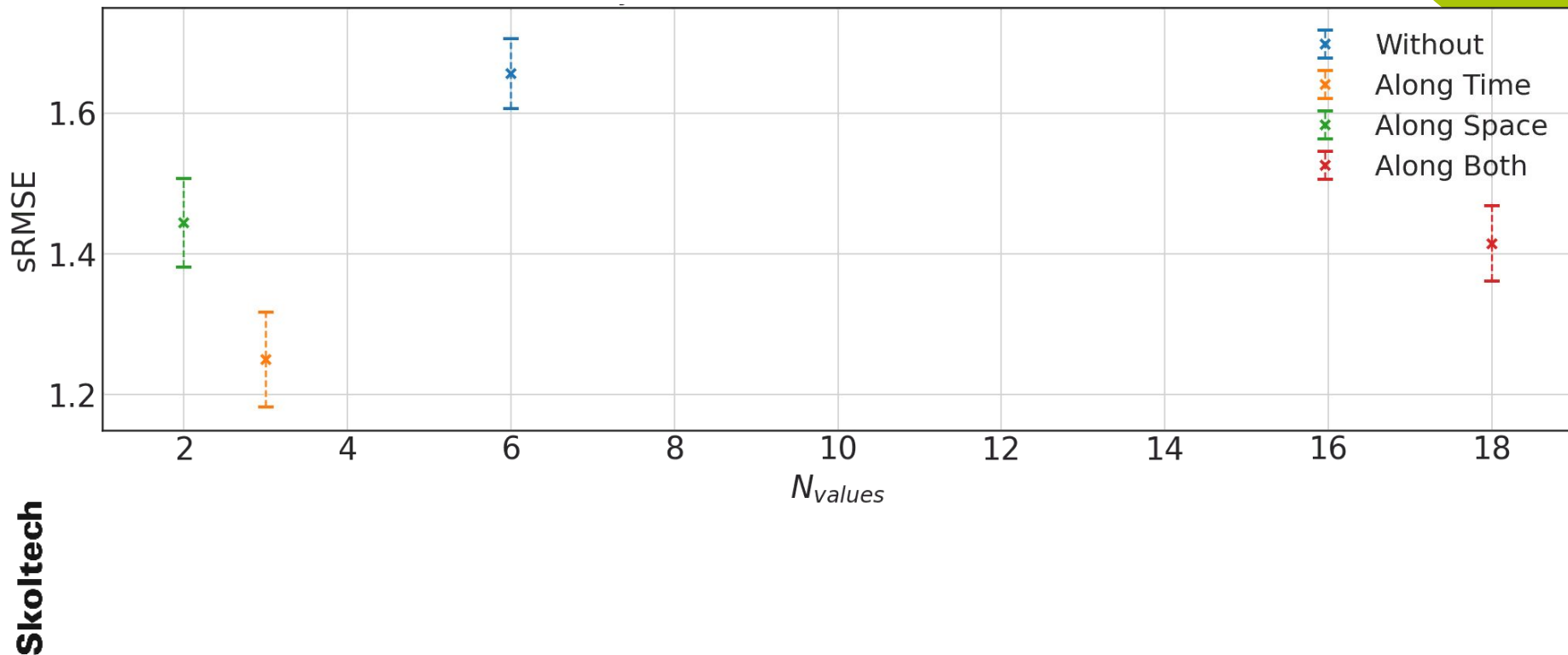
Results for NeuroTycho (ch=64) with MPCA



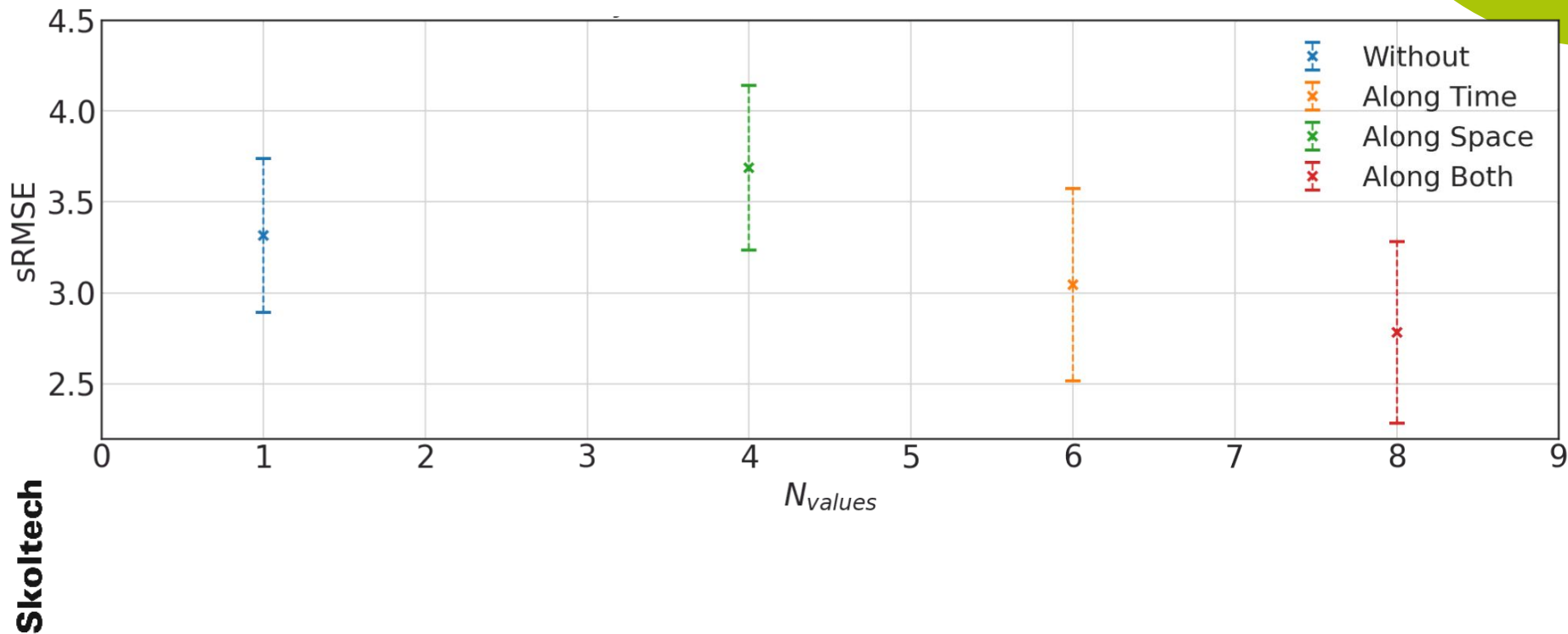
Results for EEG (ch=60) with MPCA



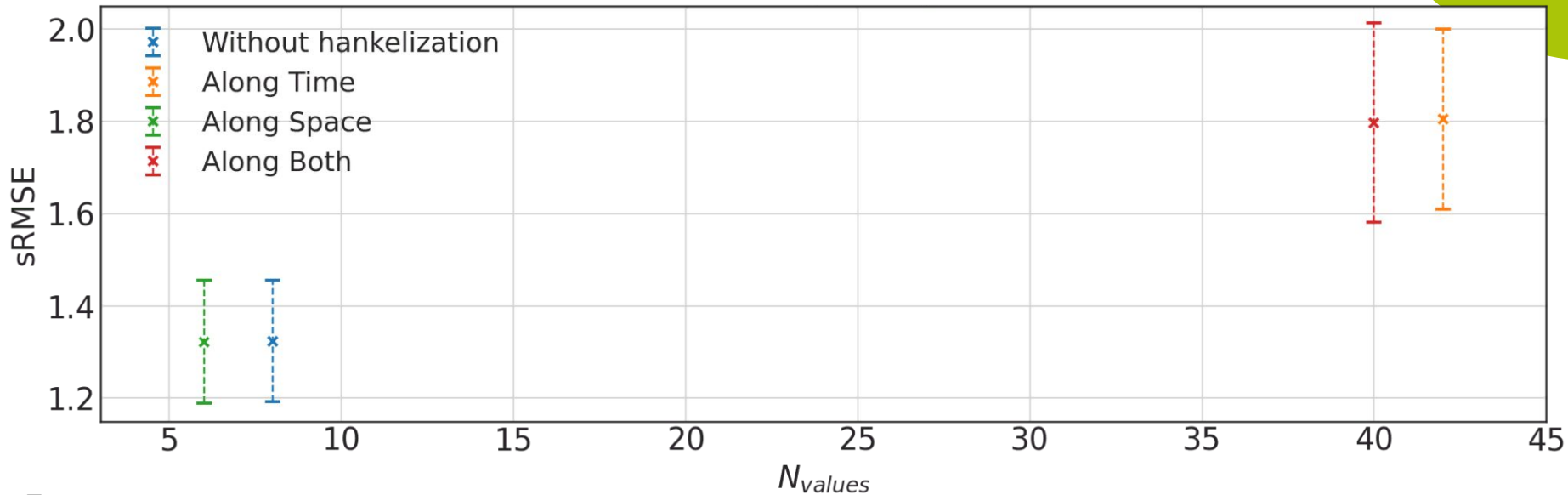
Results for NeuroTycho (ch=32) with HOPLS



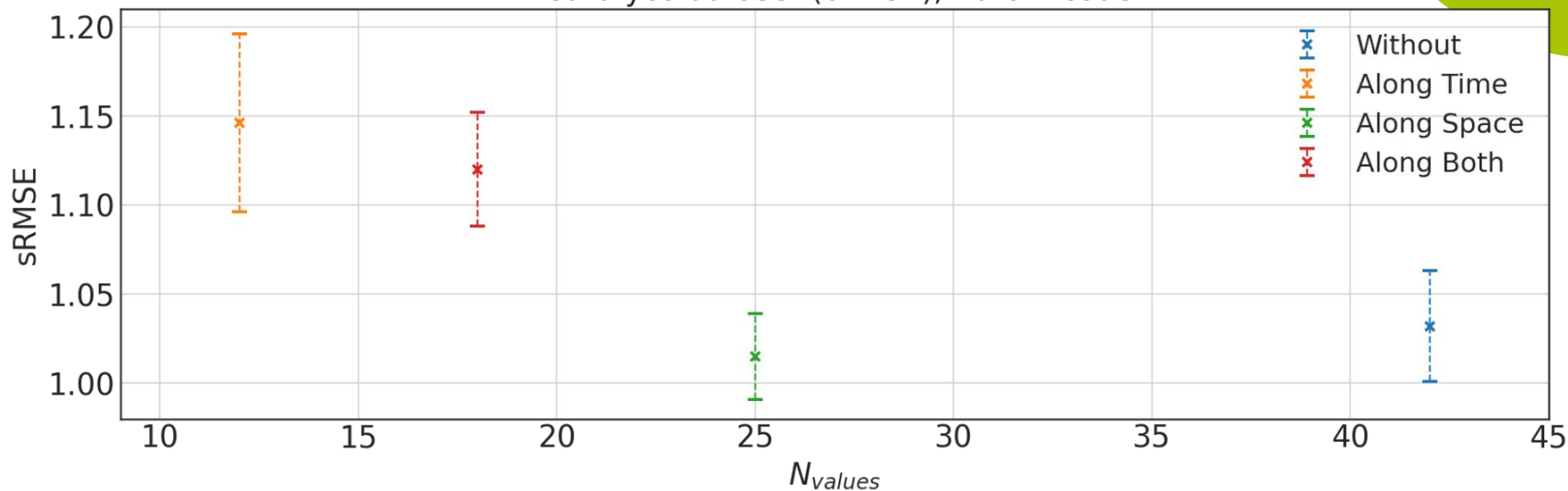
Results for NeuroTycho (ch=64) with HOPLS



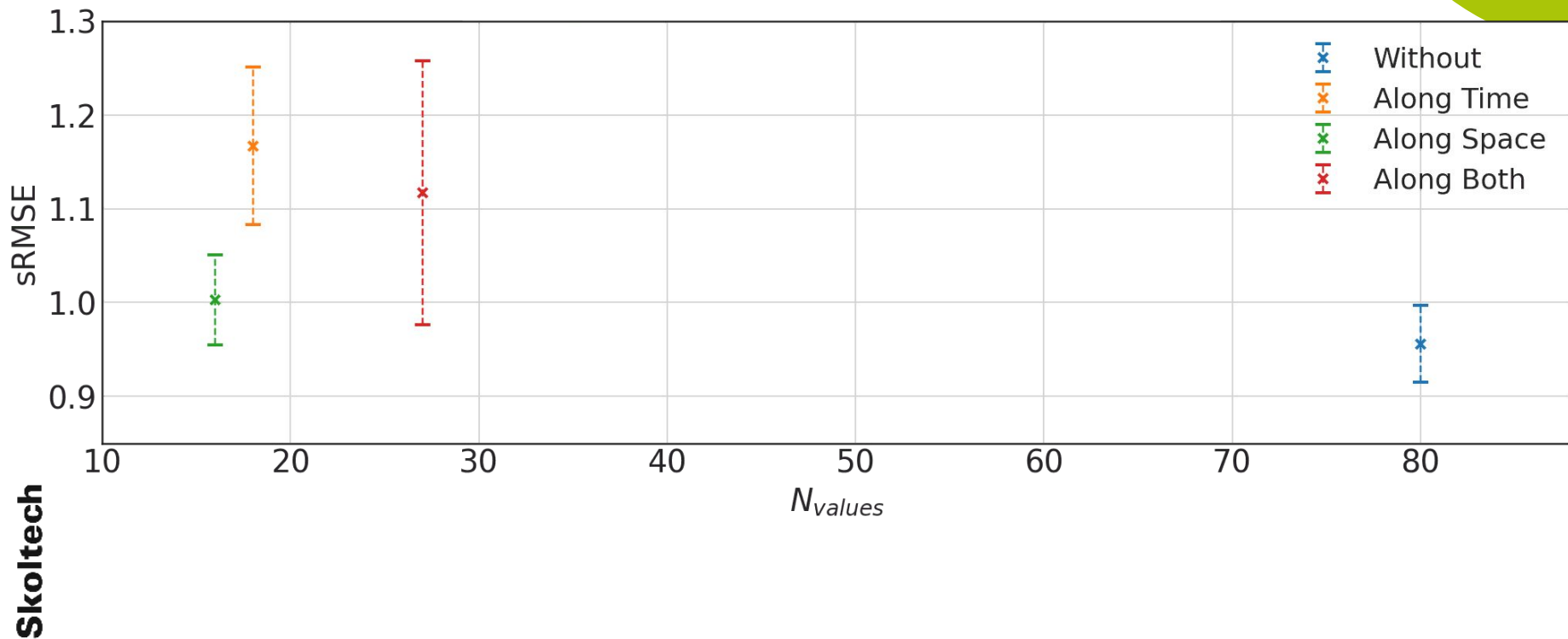
Results for EEG (ch=60) with HOPLS



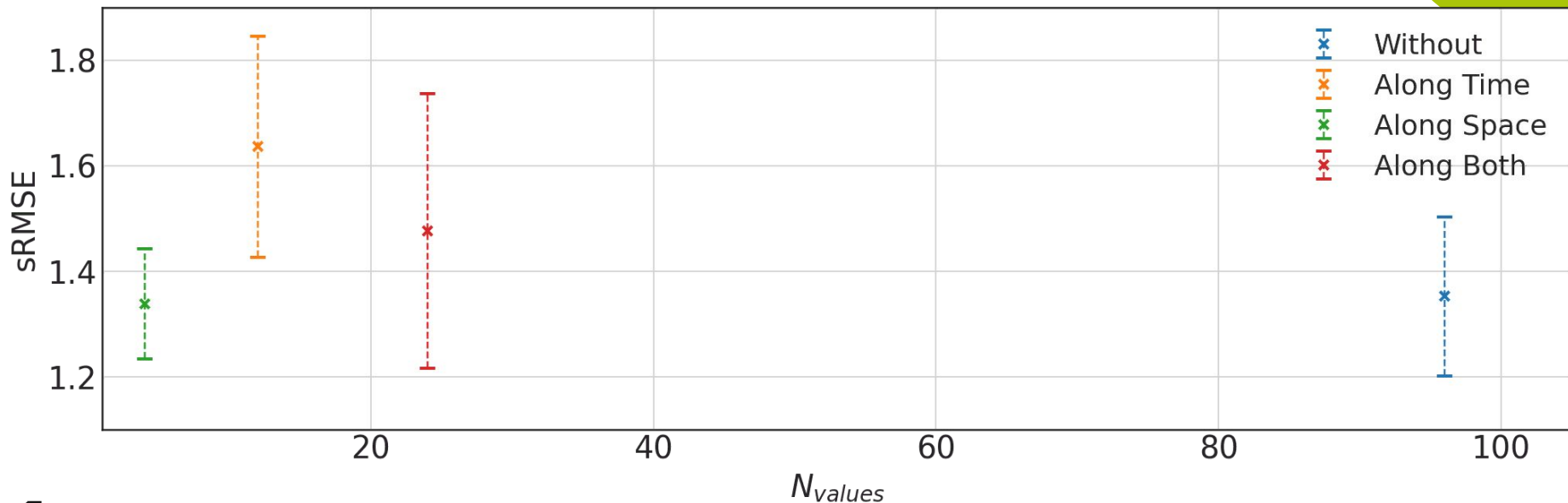
Results for NeuroTycho (ch=32) with AutoEncoder



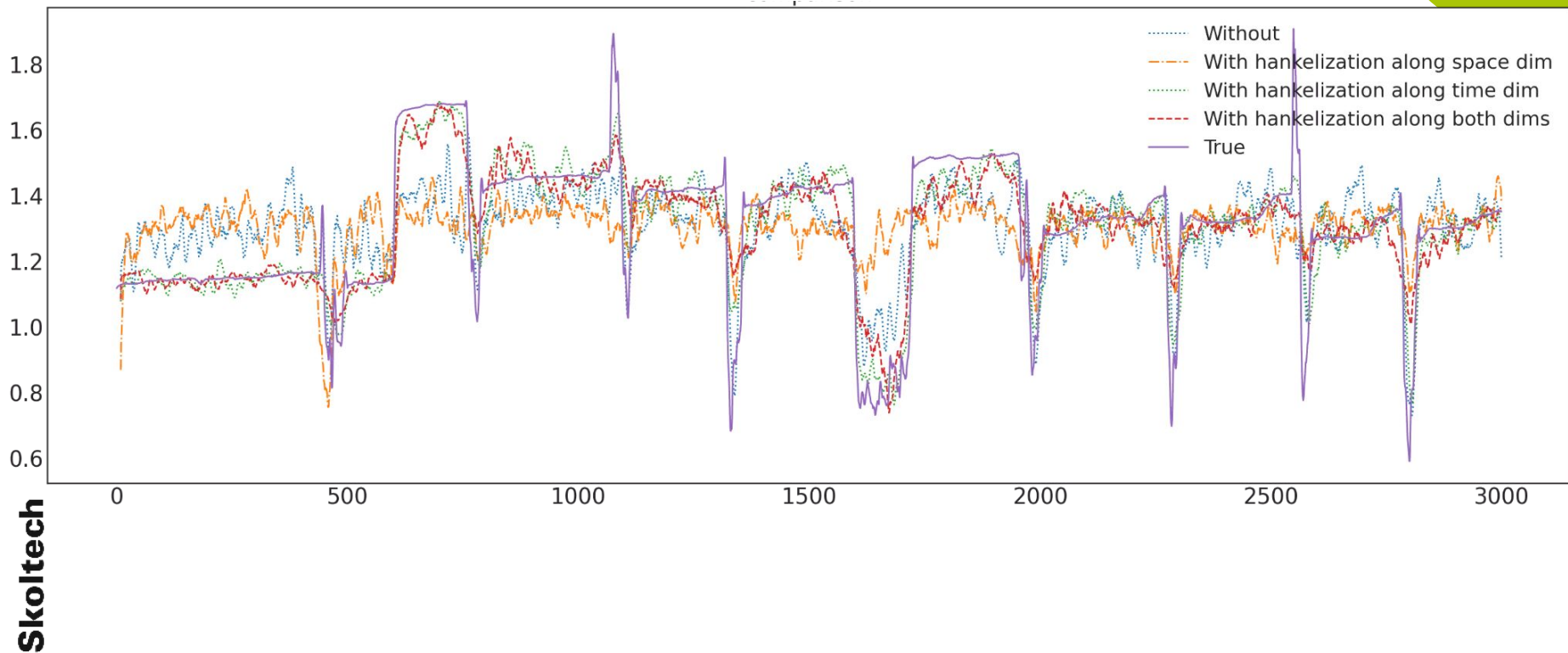
Results for NeuroTycho (ch=64) with AutoEncoder



Results for EEG (ch=60) with AutoEncoder



Example of predictions for AutoEncoder



Summary of the results

sRMSE:

- Hankelization along space dimension decreased sRMSE for the EEG dataset and for autoencoder for all datasets;
- Hankelization along time dimension decreased sRMSE for HOPLS for the NeuroTycho dataset with 32 channels for MPCA and HOPLS;
- Hankelization along both dimension decreased sRMSE for the NeuroTycho dataset with 64 channels for MPCA and HOPLS;

N_{values} :

- In half of the cases, hankelization gives increased dimensionality of the latent tensor;
- For HOPLS and AutoEncoder , hankelization mostly gives smaller latent tensors;

Discussion

- Previously, it was shown that the forecasting of time series not for BCI is better with hankelization only along temporal dimension.
- For BCI:
 - hankelization along spatial dimension works the best way for one of the datasets and for autoencoder.
 - hankelization along temporal and spatial dimensions works better than only along temporal dimension for the other dataset.
- It can be because of high correlation between data from different electrodes.



Thank you for attention!