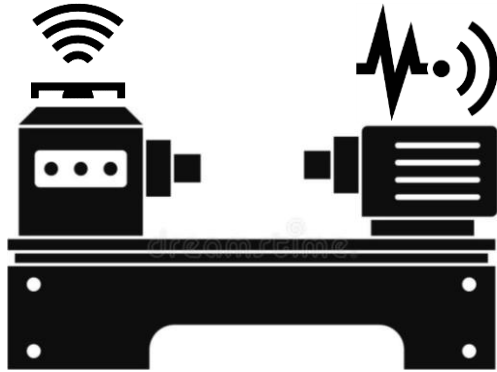


Scientific seminar
SUSU Big Data and Machine Learning Lab

Method of imputation missing values in multivariate streaming time series

Alexey Yurtin, Mikhail Zymbler

Processing of streaming time series



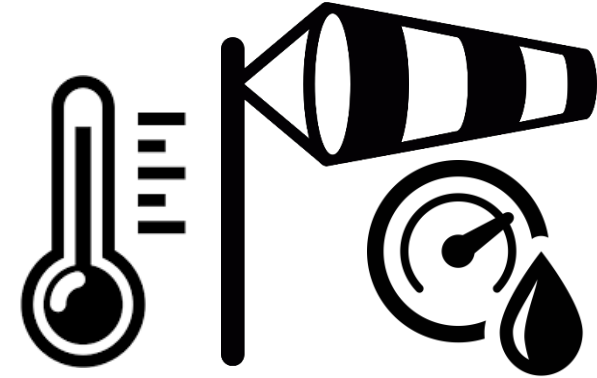
**Predictive maintenance,
smart manufacturing**



**Internet
of Things**



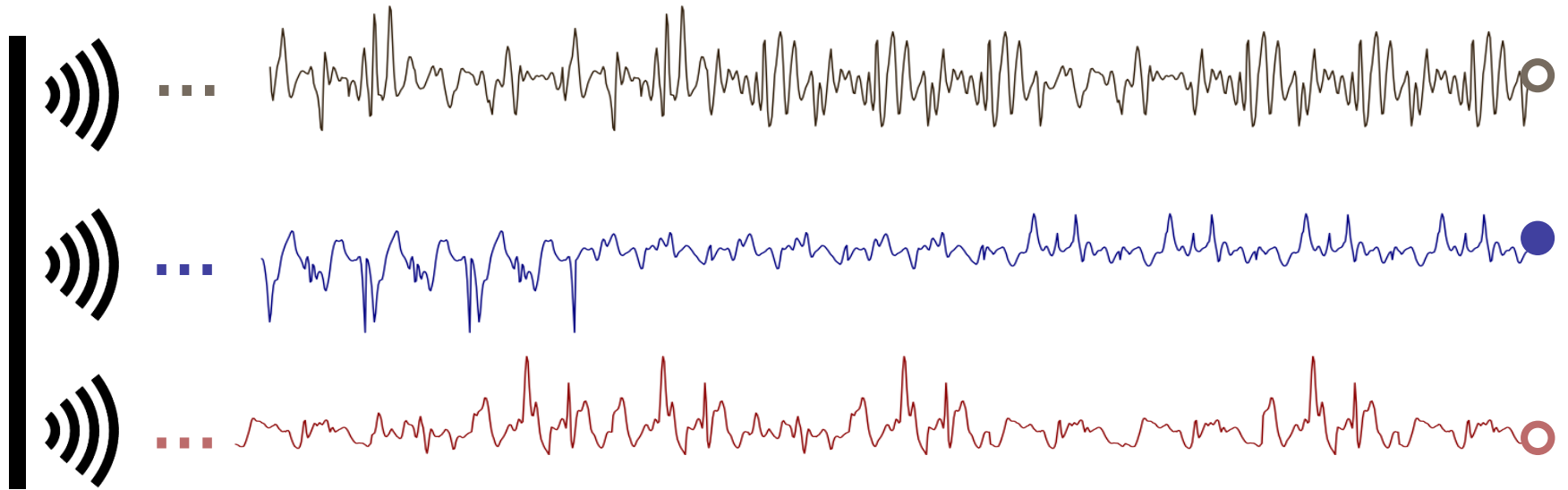
**Personal
healthcare**



**Weather forecasting,
climate modelling**

**Elements of a time series
arrive one after another
in several dimensions in an online mode**

Imputation of missing values



How to plausibly and quickly
synthesize missing values?

SANNI: Snippet and ANN-based Imputation

1. Labeling

- Automatically label a representative fragment of the time series undergo imputation through the **snippet concept**

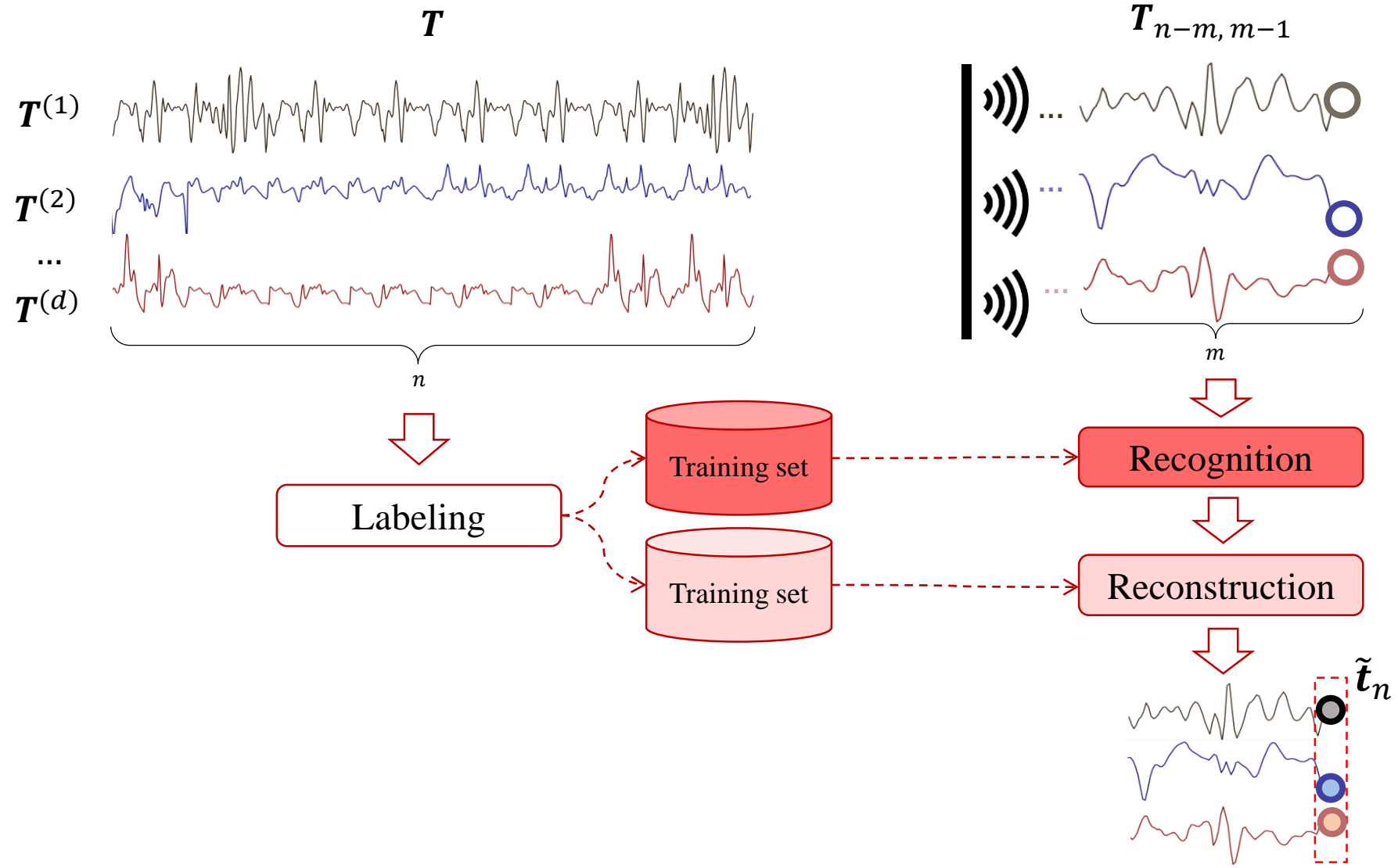
2. Recognition

- Recognize a snippet in a subsequence before the missing value through the **convolutional neural networks**

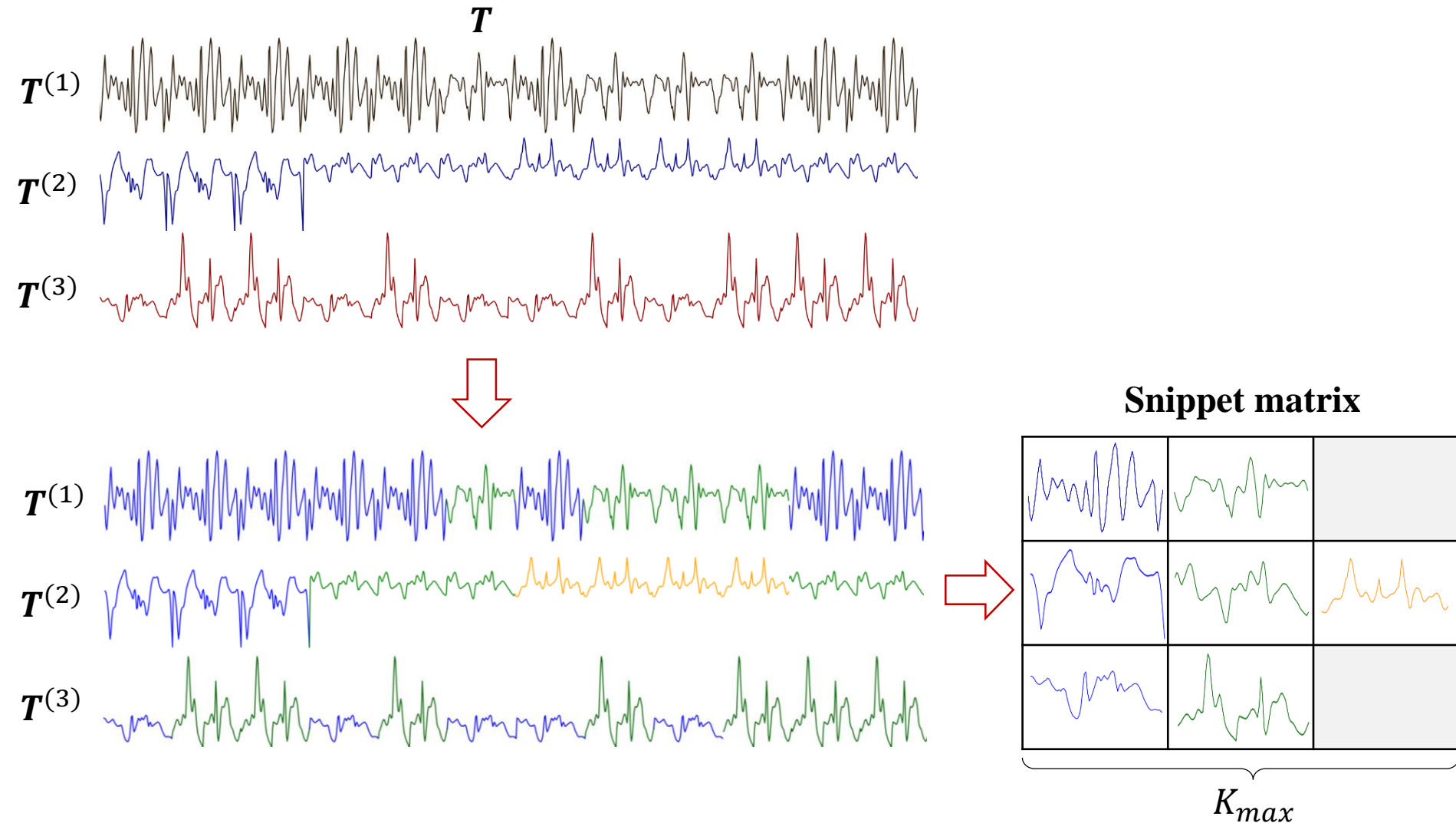
3. Reconstruction

- Generate a missing value based on the recognized snippet and **recurrent neural networks**

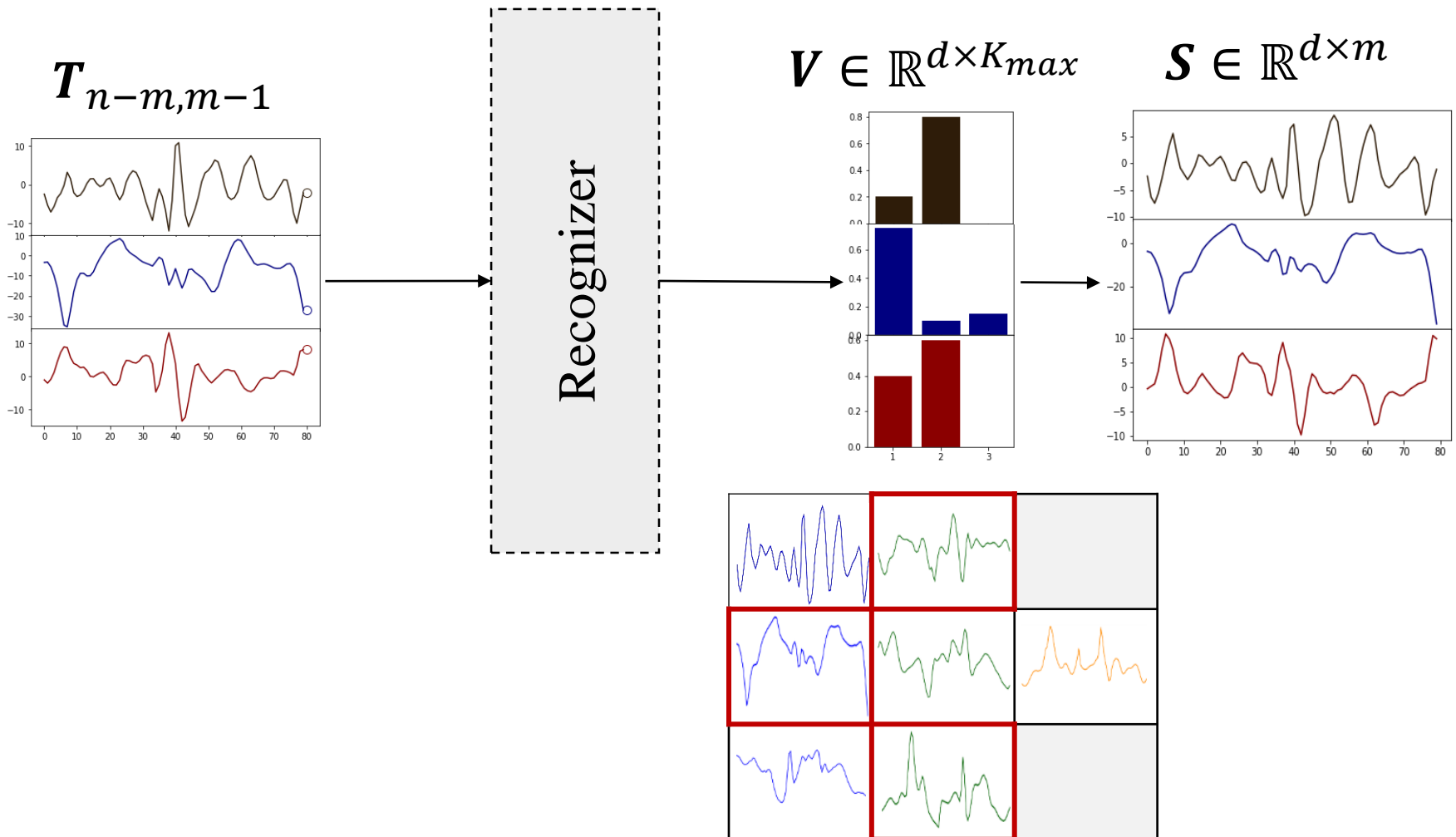
SANNI: Architecture



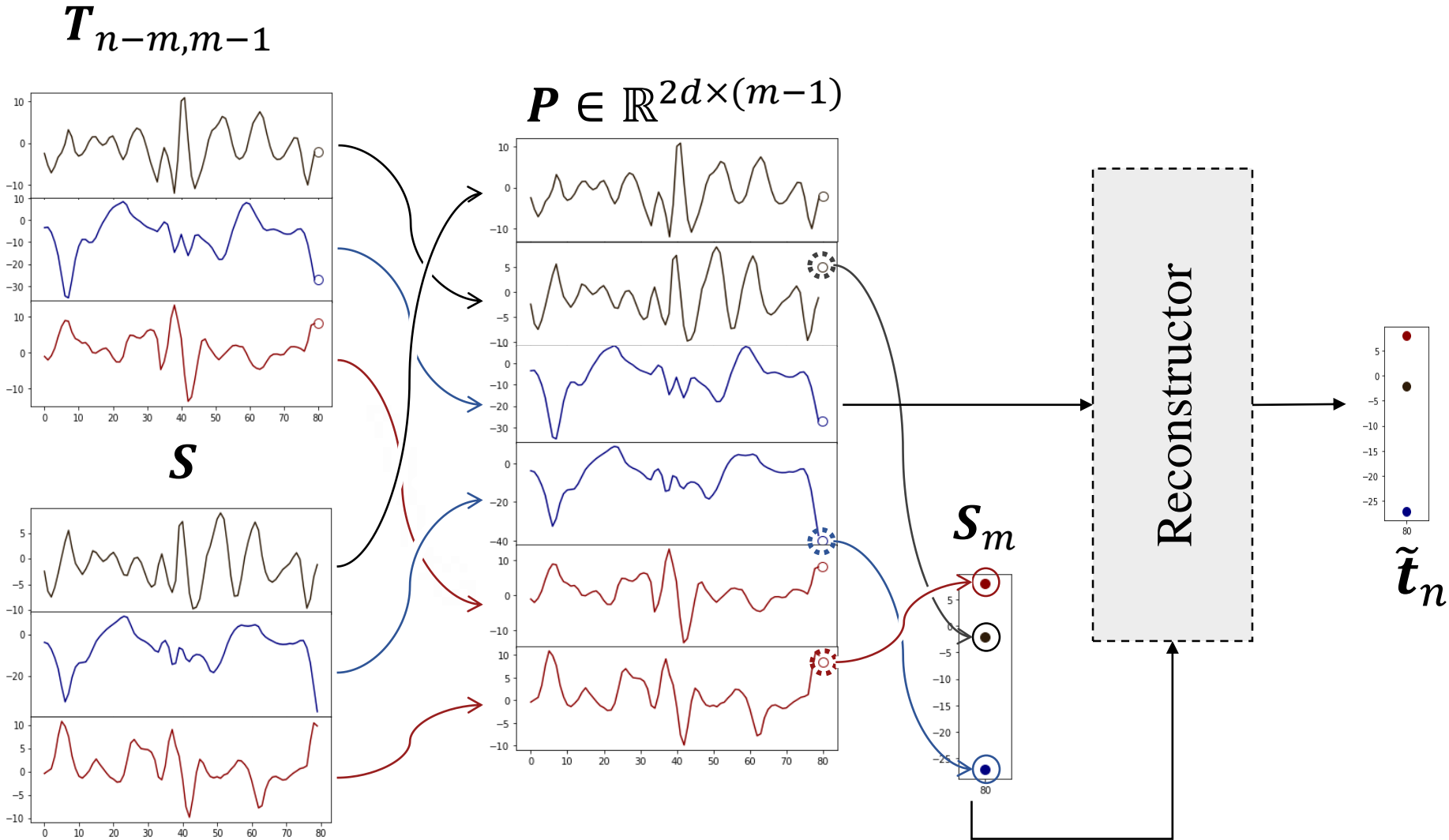
SANNI: Labeling



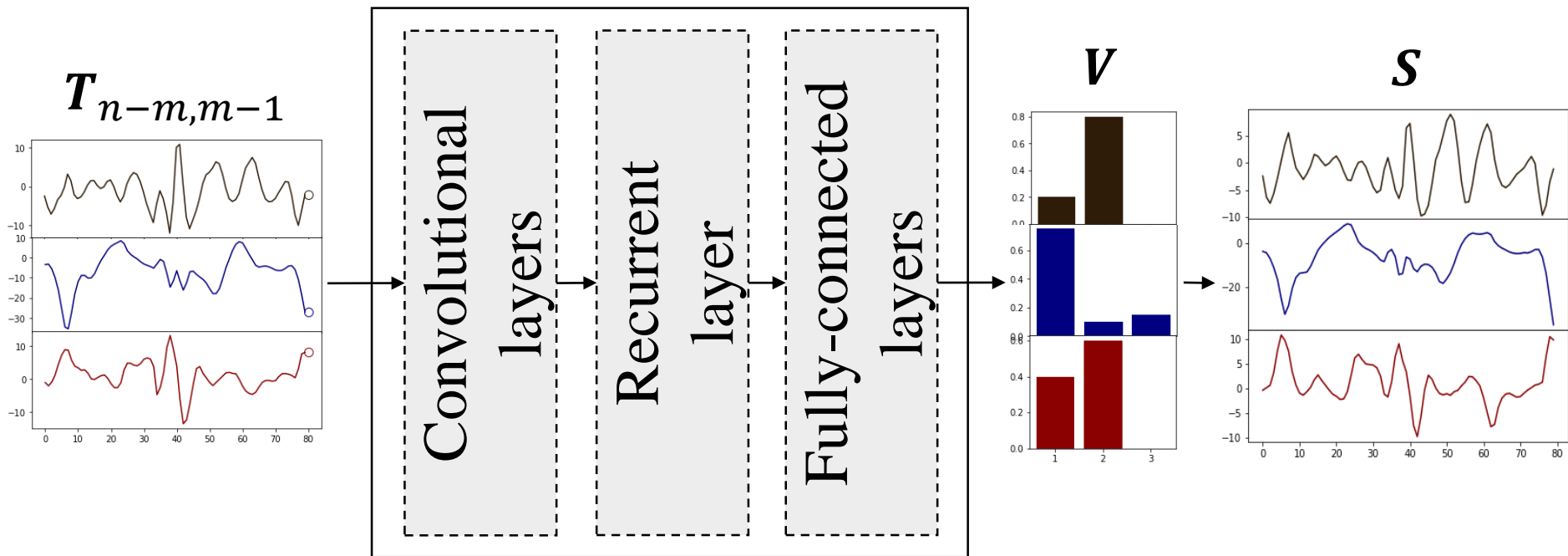
SANNI: Recognition



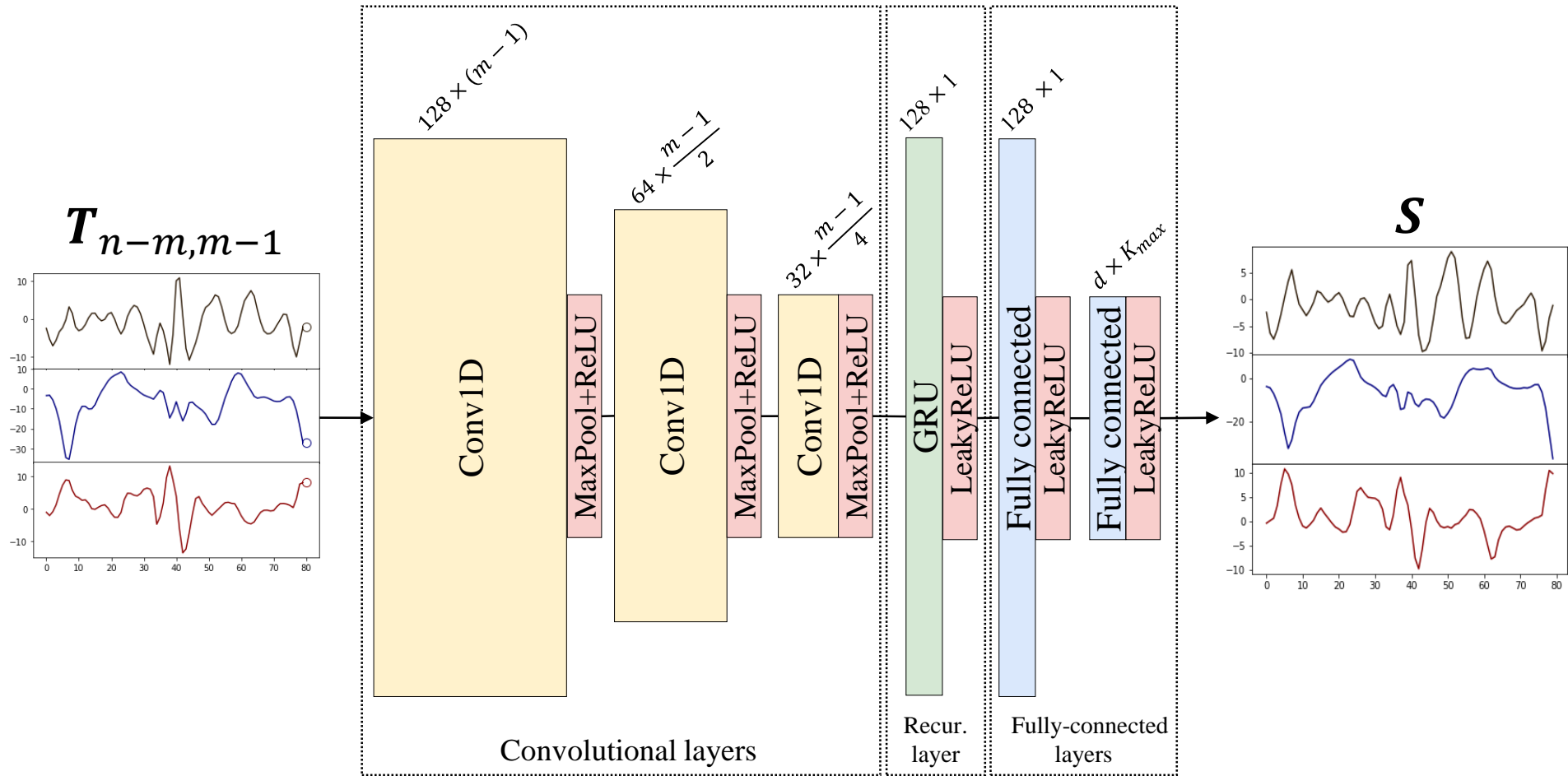
SANNI: Reconstruction



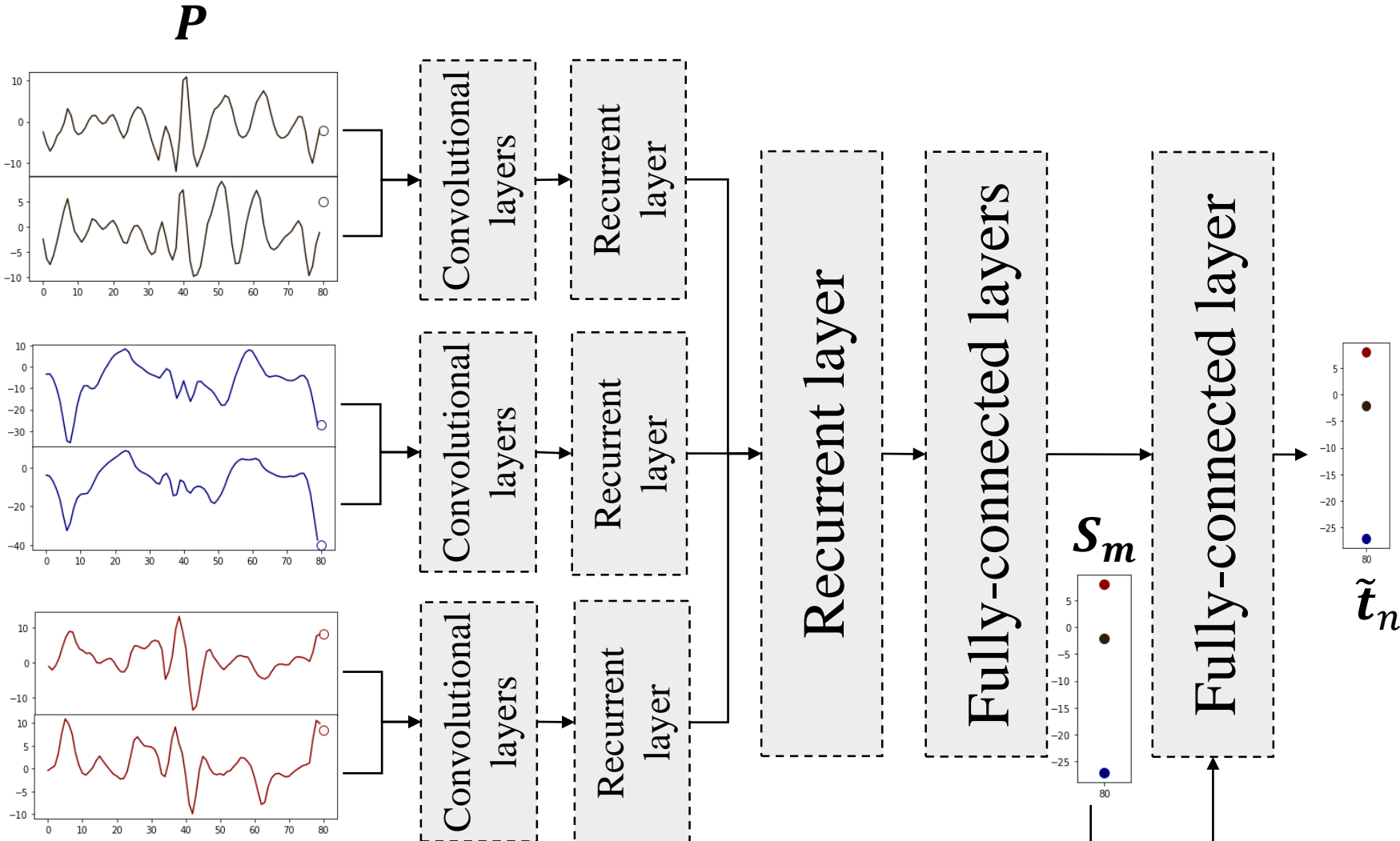
Recognizer



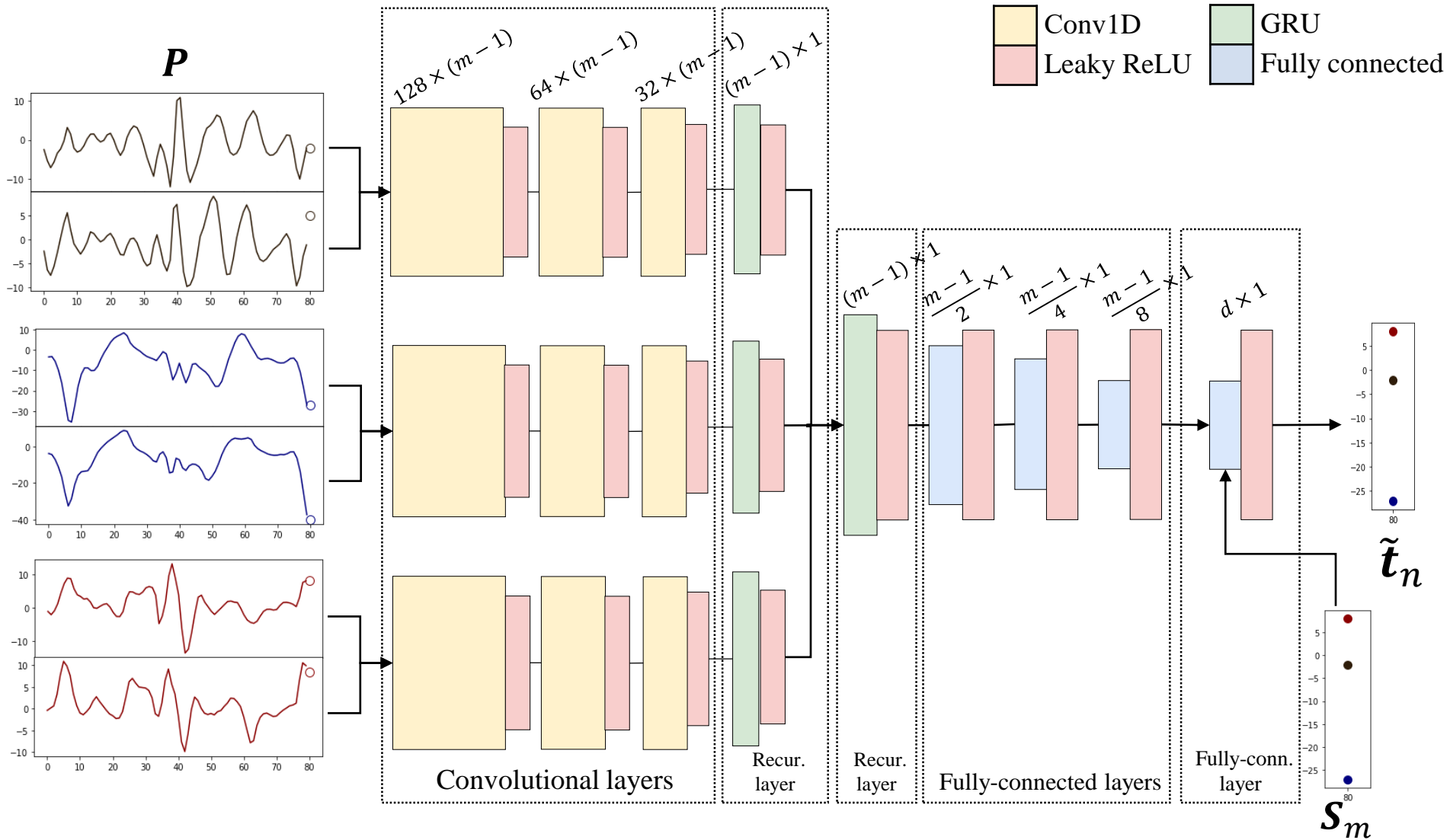
Recognizer



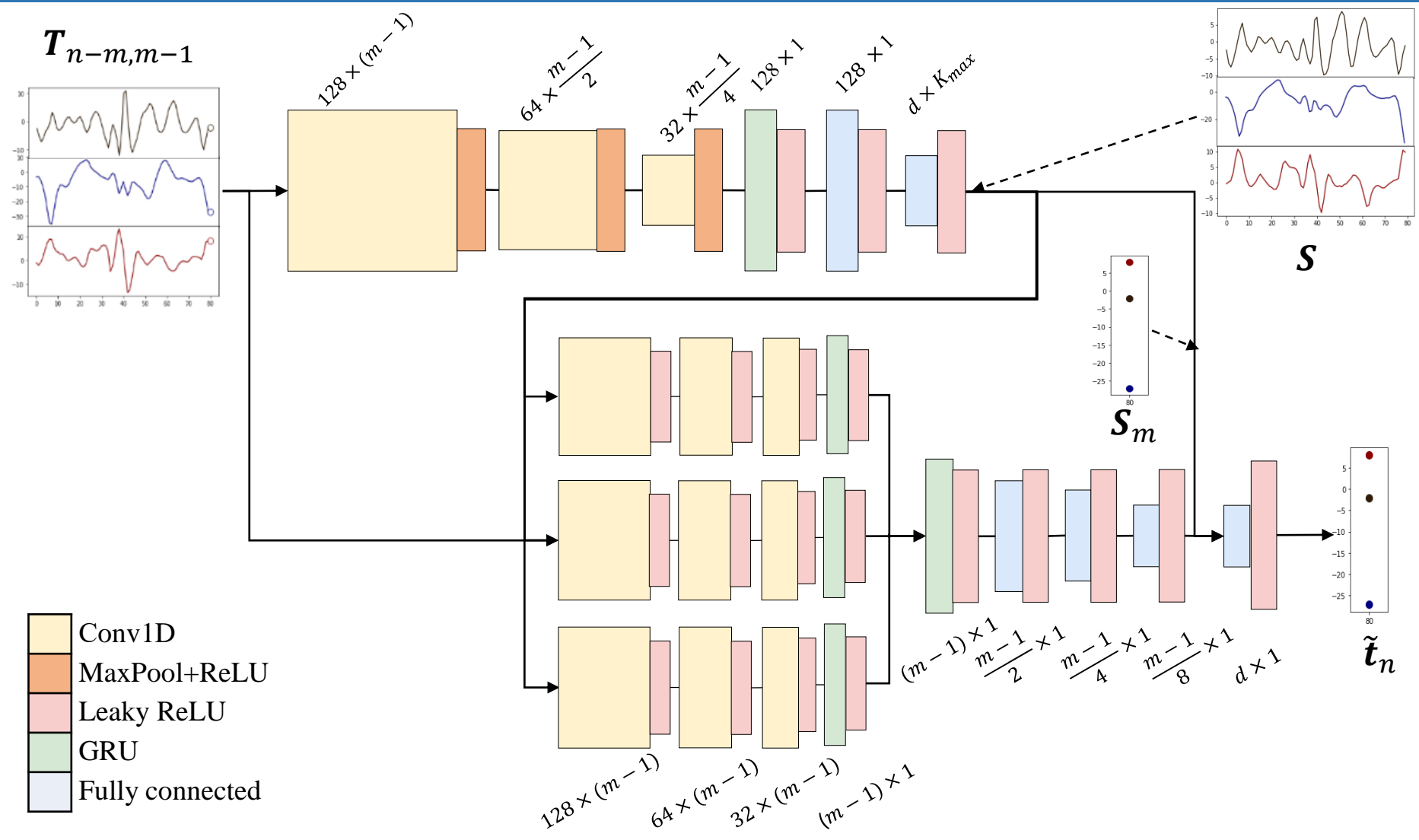
Reconstructor



Reconstructor

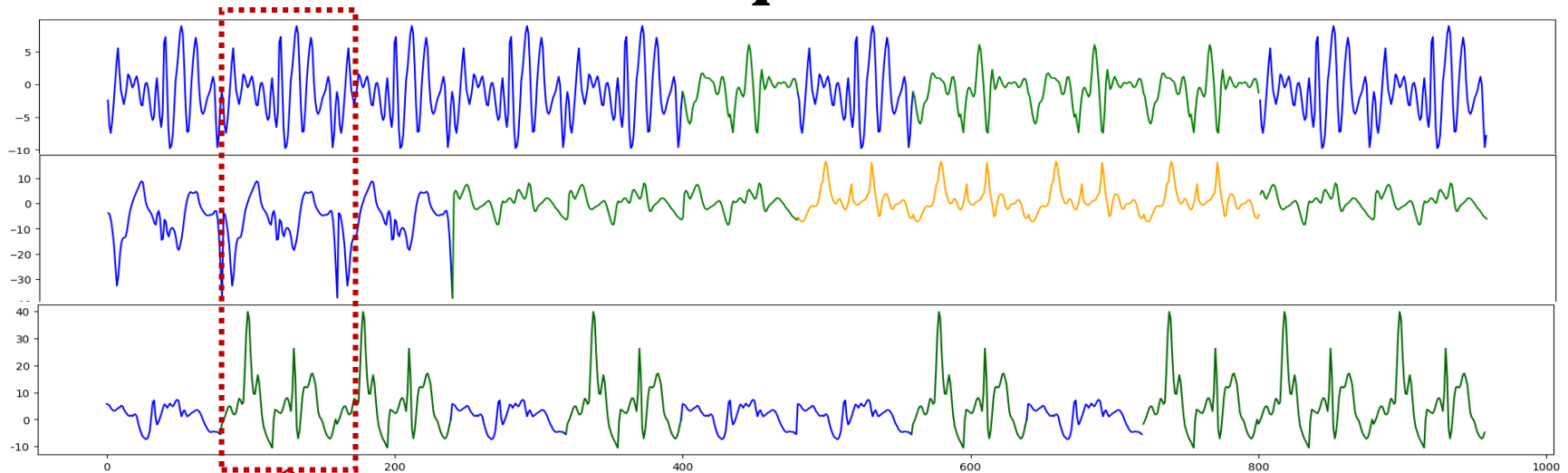


SANNI

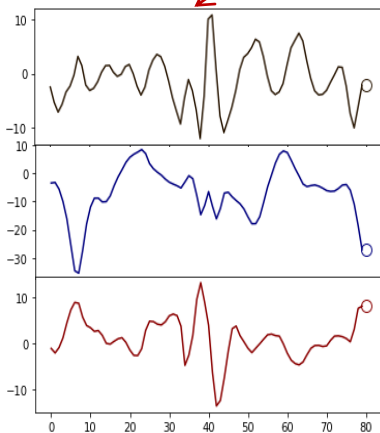


Learning of Recognizer

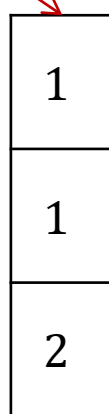
T



X



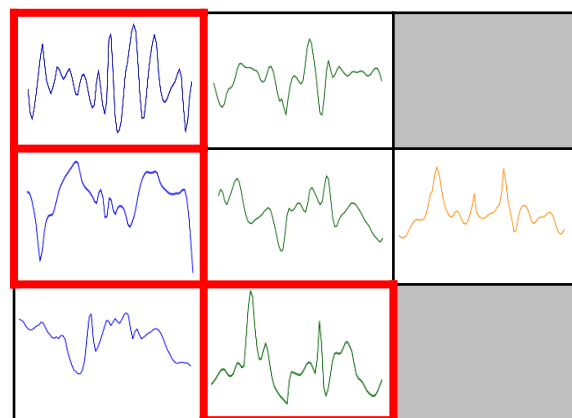
Y



1

2

3



Loss function: **cross entropy**

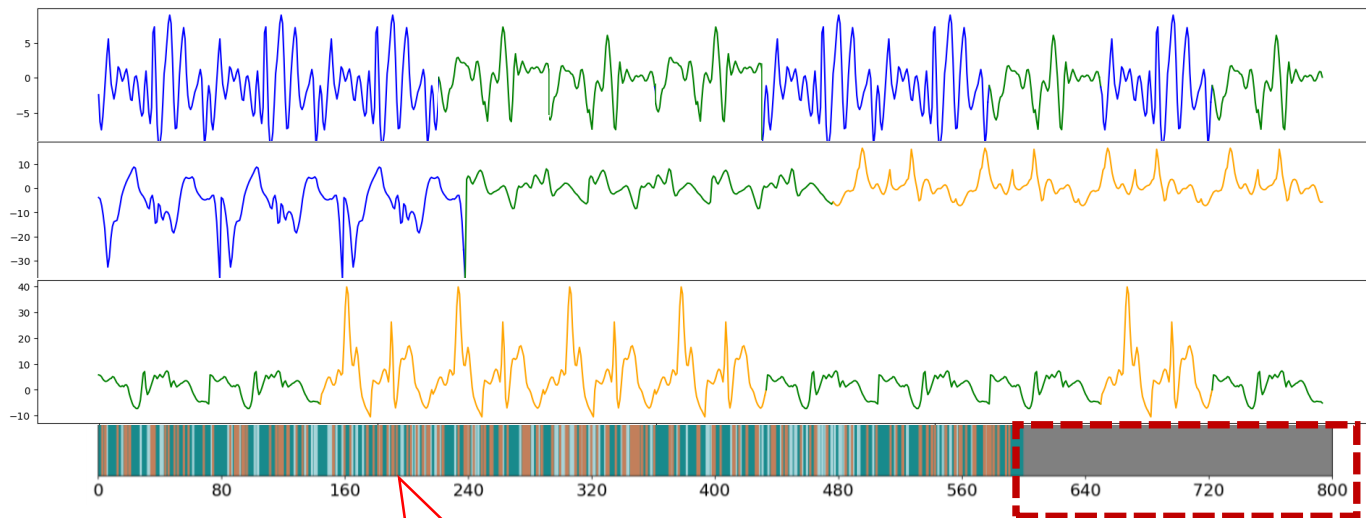
Optimizer: **Adam**

Epochs: **1000**

Learning rate: **10^{-3}**

Learning of Recognizer

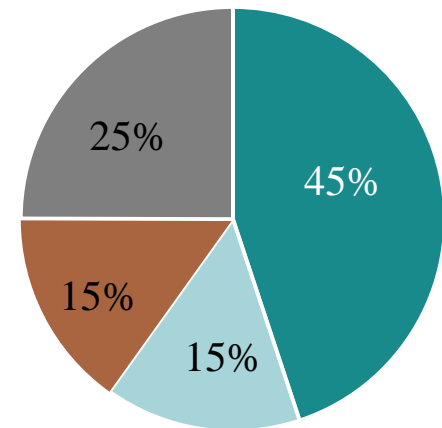
T



Shuffle to evenly
distribute snippets
across the training set

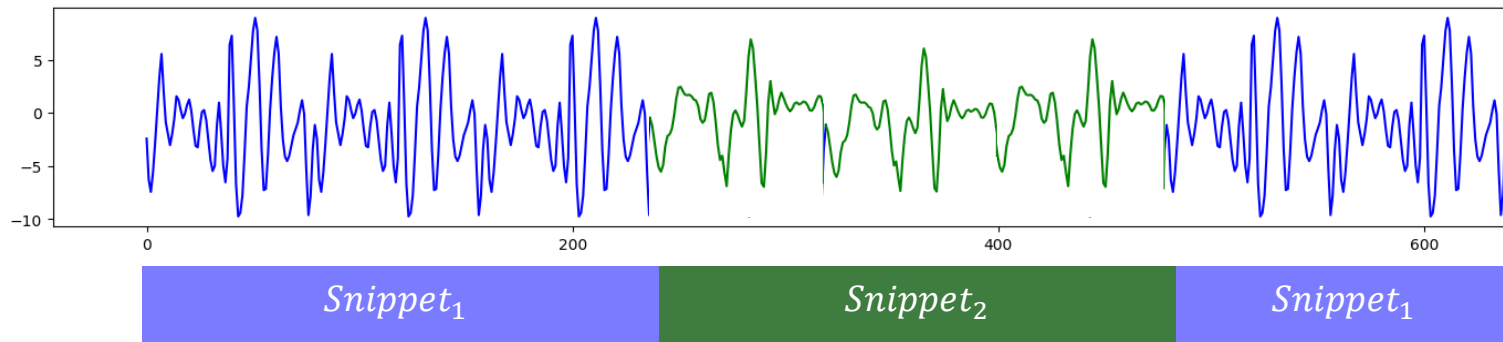
Leave
for the SANNI
testing

Data splitting

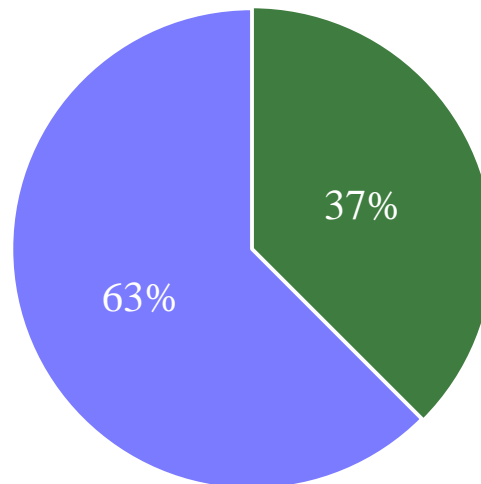


- Train
- Validate
- Test
- Unattended

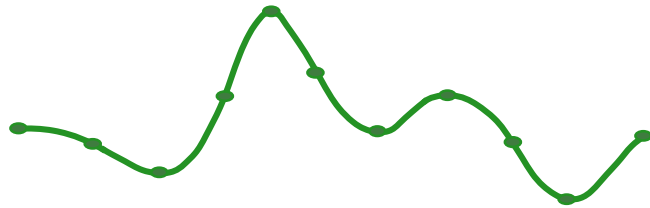
Snippets with small coverage



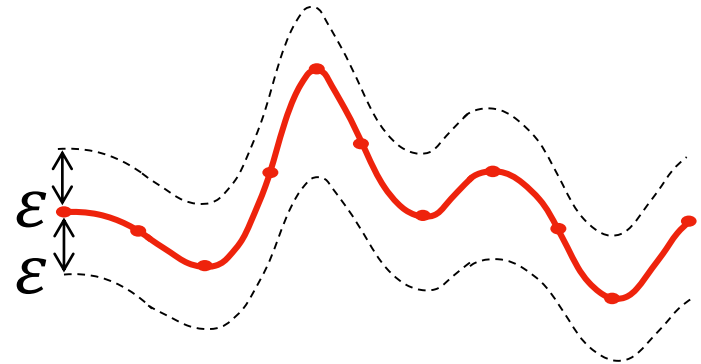
Coverage



Augmentation of snippets with small coverage



A snippet



A nearest neighbor of the snippet

$$\varepsilon = \text{ED}(\text{snippet}, \text{neighbor})$$

snippet.NN is a set of all the nearest neighbors of the snippet

$$\forall \text{neighbor} \in \text{snippet.NN} \exists c \in R^m: \sum_{k=1}^m c_k = \varepsilon,$$

$$\begin{aligned} \text{synthetic1} &= \text{neighbor} + c, \\ \text{ED}(\text{snippet}, \text{synthetic1}) &< \varepsilon, \end{aligned}$$

$$\begin{aligned} \text{synthetic2} &= \text{neighbor} - c \\ \text{ED}(\text{snippet}, \text{synthetic2}) &< \varepsilon \end{aligned}$$

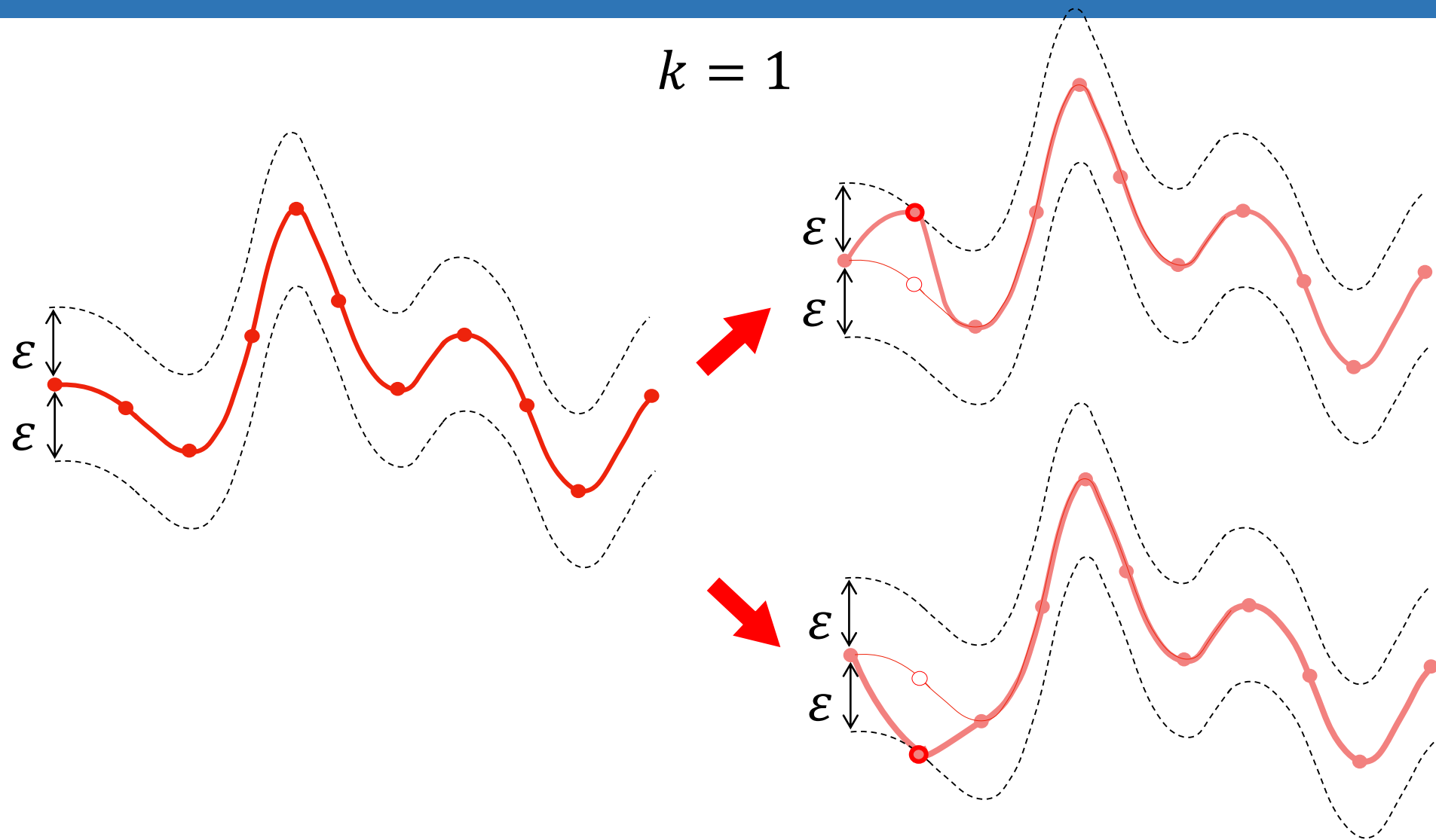
Number of synthetic neighbors is

$$C_m^{m+k-1} = \frac{(m+k-1)!}{m!(k-1)!} *$$

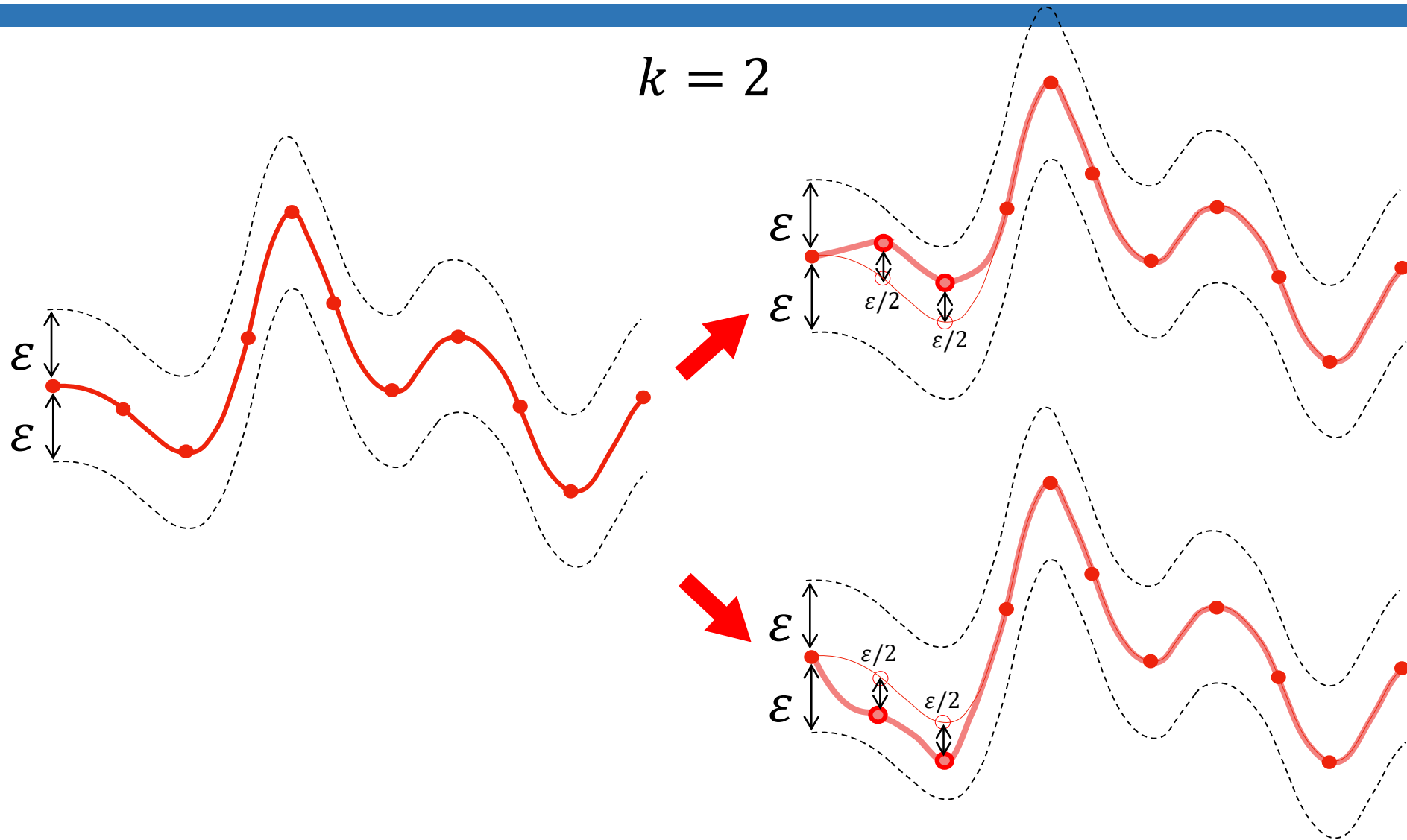
* Reingold E. *et al.* Combinatorial Algorithms: Theory and Practice. Prentice Hall, 1977. 433 p.

Synthetic neighbors of a small-coverage snippet

$k = 1$

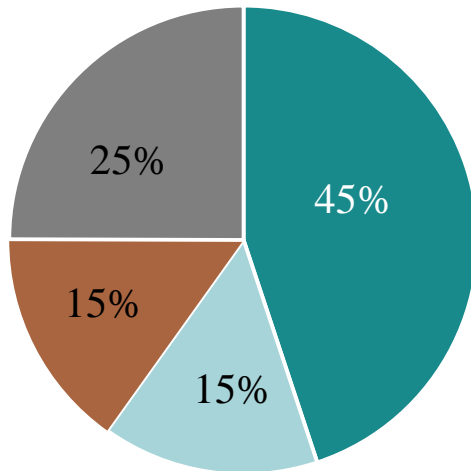


Synthetic neighbors of a small-coverage snippet



Learning Recognizer with augmentation data

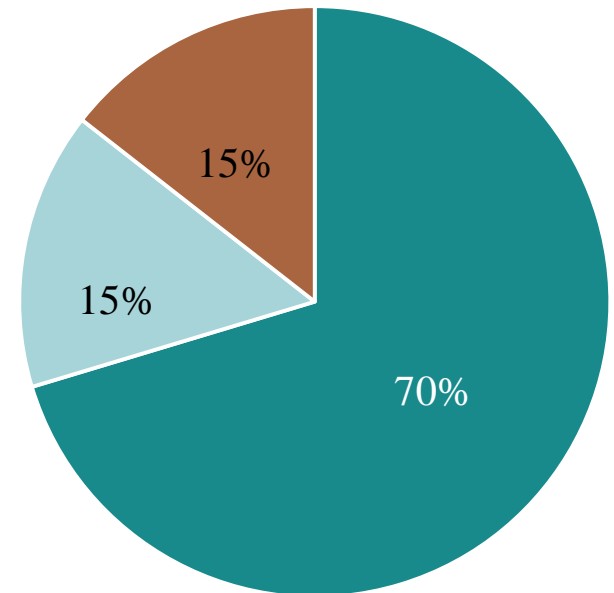
Data splitting



- Train
- Validate
- Test
- Unattended

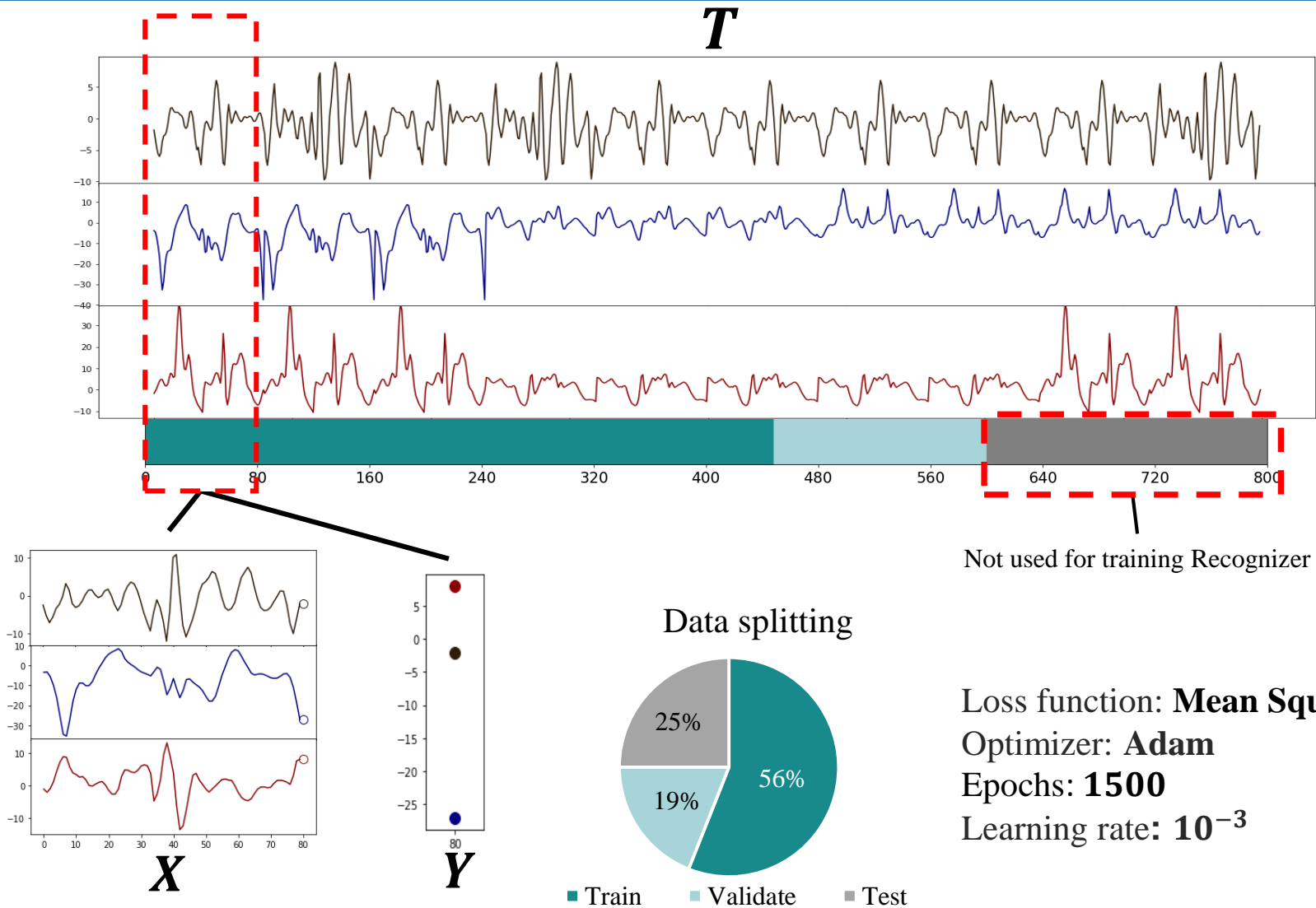
Augmentation

Data splitting



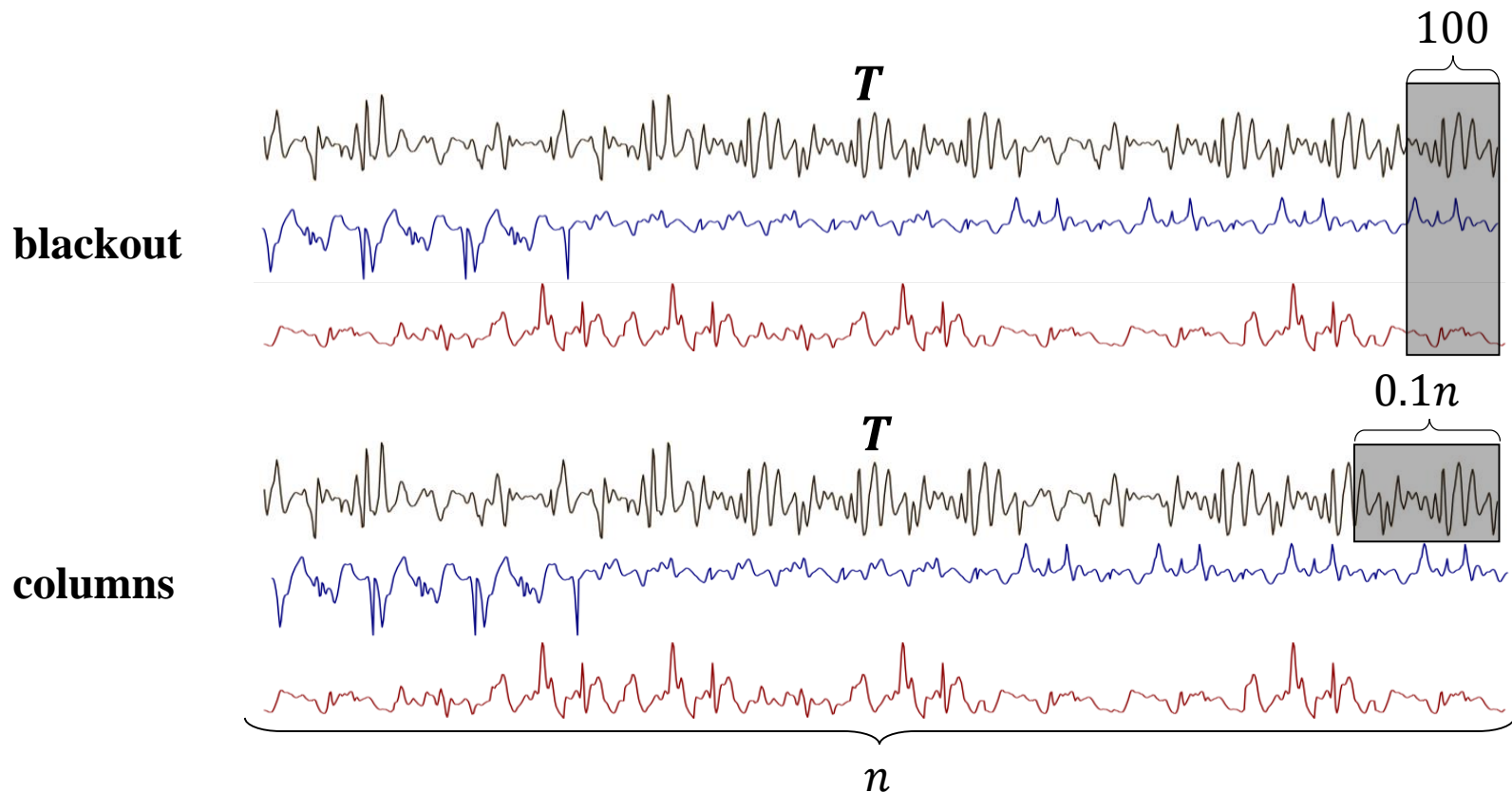
- Train
- Validate
- Test

Learning of Reconstructor





Experiments: Setup

- GPU: NVIDIA Tesla V100 SXM2 (5120 cores @1.312 GHz, memory 32 Gb).
- Scenario under ORBITS¹ framework:



¹ Khayati M., *et al.* ORBITS: Online Recovery of Missing Values in Multiple Time Series Streams. Proc. VLDB Endow. 2020. Vol. 14, no. 3. P. 294-306.

Experiments: Datasets

 With activities
 No activities

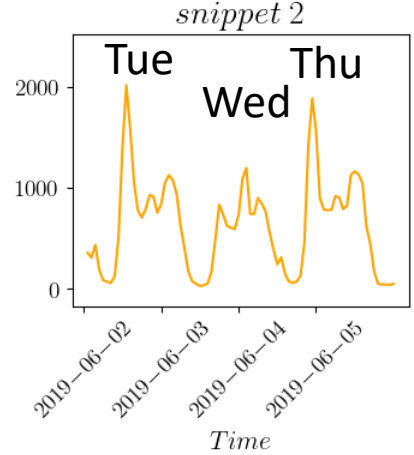
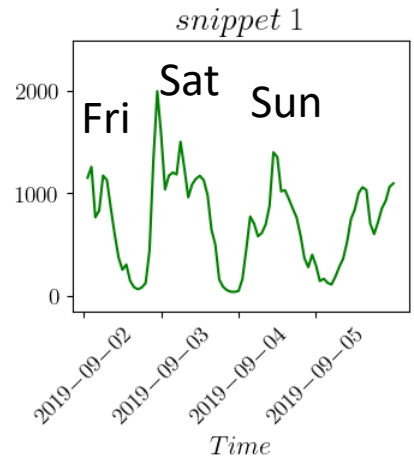
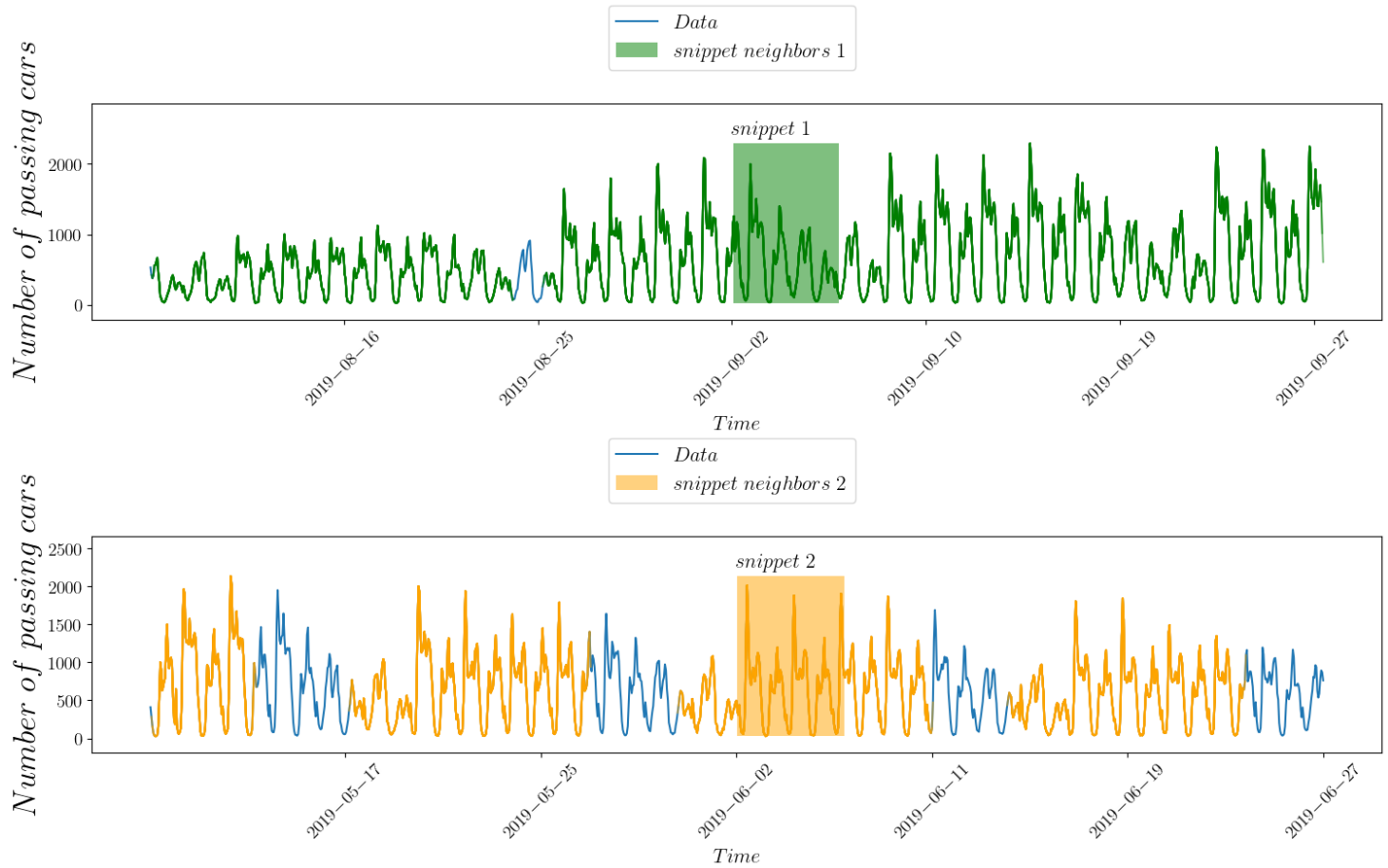
Dataset	Length n , $\times 10^3$	Dim. d	Snippet length m	Domain
Airq ¹	1	10	100	Air pollution in a city
BAFU ¹	50	12	48	Water discharge in Swiss rivers
Electricity ¹	5	20	120	Power demand of several consumers
Temp ¹	5	50	50	Daytime air temperature in different China regions
Gas ¹	3	100	100	Gas concentration measurements of a smart gas delivery platform at San Diego, US
Motion ²	10	20	120	Smartphone accelerometer data
MADRID ²	29	10	110	Road traffic (AVR statistics) in Madrid
WalkRun ³	220	19	100	Accelerometer, gyroscope, and magnetometer measurements during alternation of running and walking
WalkStop ³	140	19	250	Accelerometer, gyroscope, and magnetometer measurements during alternation of walking and still standing

¹ Khayati M., *et al.* ORBITS: Online Recovery of Missing Values in Multiple Time Series Streams. Proc. VLDB Endow. 2020.

² Zymbler M.L., Poluyanov A.N., Kraeva Ya.A. Parallel Algorithm for Real-time Sensor Data Recovery for a Many-core Processor. Bulletin of the South Ural State University. Series: Computational Mathematics and Software Engineering, 2022. Vol. 11, no. 3. P. 68–89. (in Russian) DOI: 10.14529/cmse220305

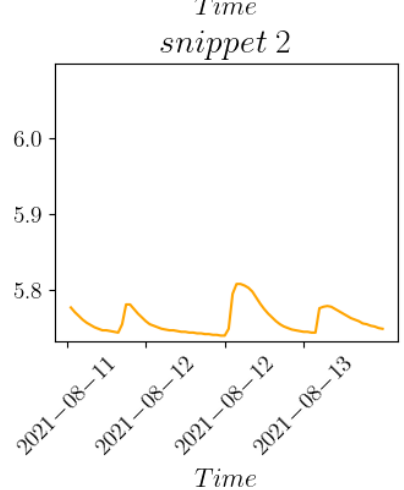
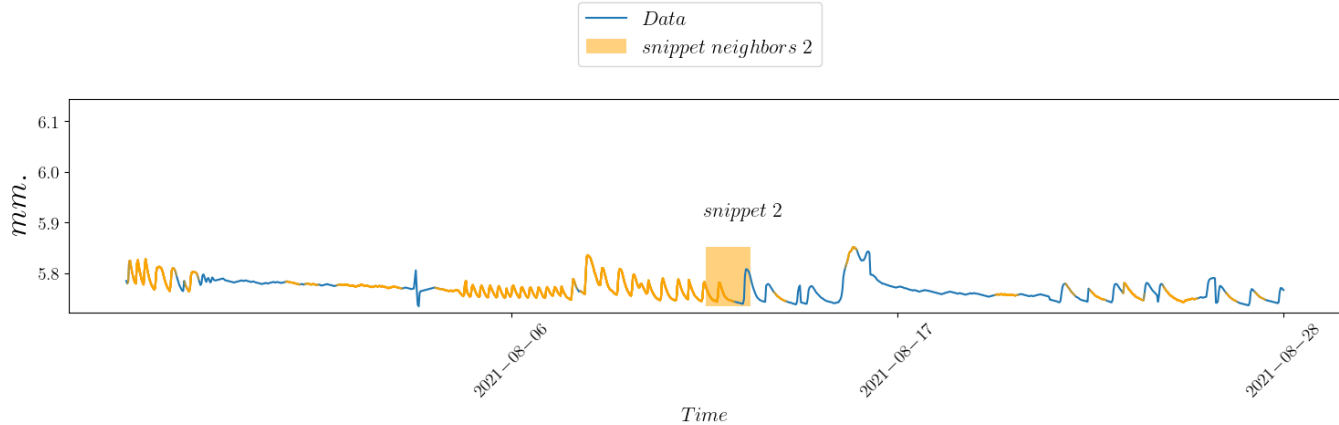
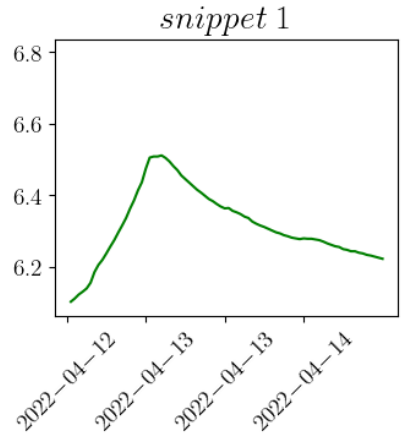
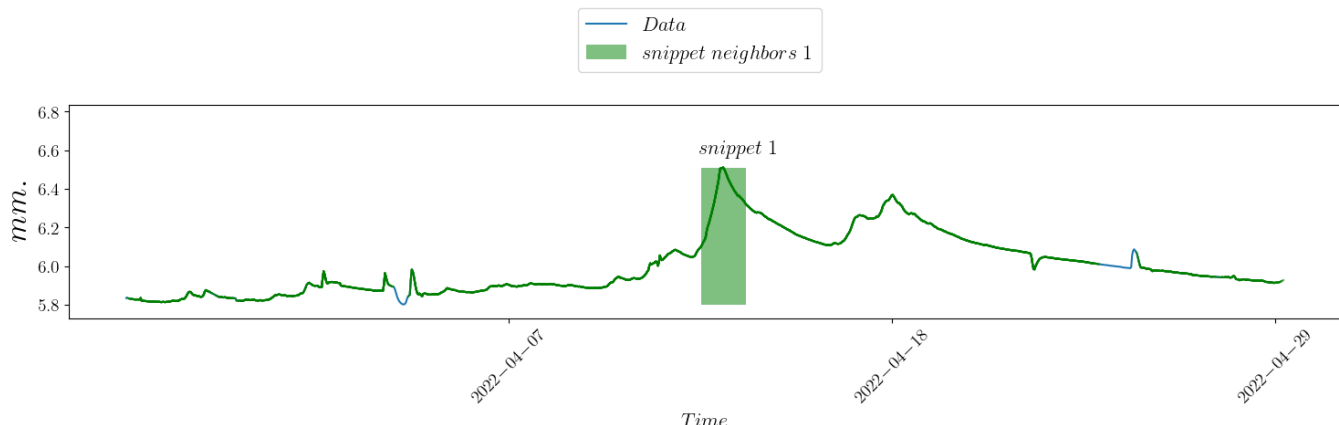
³ Our generated dataset

MADRID: dataset with activities



Snippet length is 110 points (4 days)

BAFU: dataset no activities



Snippet length is 48 points (3 days)

Experiments: Rivals

- Analytical approaches
 - CD-REC¹, OGDImpute², DynaMMo³, SoftImpute⁴, TKCM⁵, GROUSE⁶, SPIRIT⁷, ORBITS⁸
- ANN-based approaches
 - BRITS⁹, NAOMI¹⁰, MRNN¹¹

¹ Khayati M., et al. Scalable recovery of missing blocks in time series with high and low cross-correlations. Proc. Computer Science 2019.

² Anava O., et al. Online Time Series Prediction with Missing Data. Proc. ICML 2015.

³ Lei L., et al. DynaMMo: mining and summarization of coevolving sequences with missing values. KDD '09. 2009.

⁴ Mazumder R., et al. Spectral Regularization Algorithms for Learning Large Incomplete Matrices. Journal of Machine Learning Research. 2010.

⁵ Wellenzohn K., et al. Continuous Imputation of Missing Values in Streams of Pattern-Determining Time Series. EDBT 2017. P. 330-341.

⁶ Zhang D., et al. Global Convergence of a Grassmannian Gradient Descent Algorithm for Subspace Estimation. Electrical Engineering and Computer Science 2017.

⁷ Papadimitriou S., et al. Streaming Pattern Discovery in Multiple Time-Series. VLDB 2005.

⁸ Khayati M., et al. ORBITS: Online Recovery of Missing Values in Multiple Time Series Streams. Proc. VLDB Endow. 2020

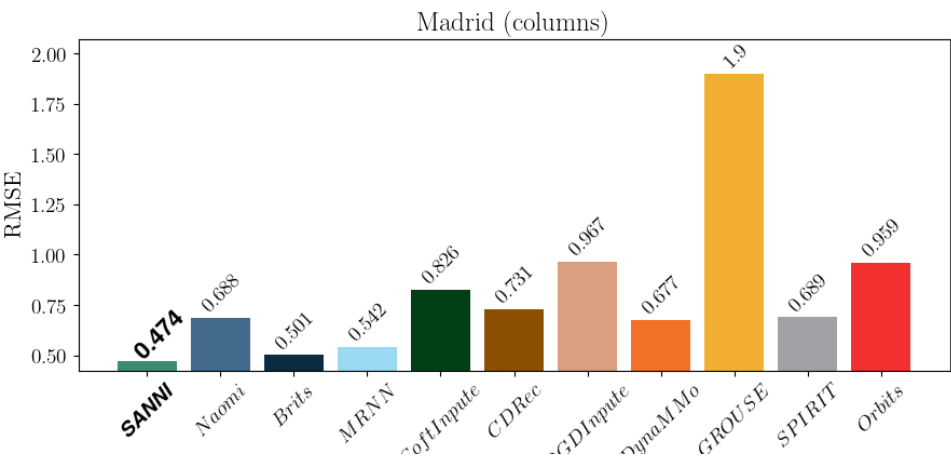
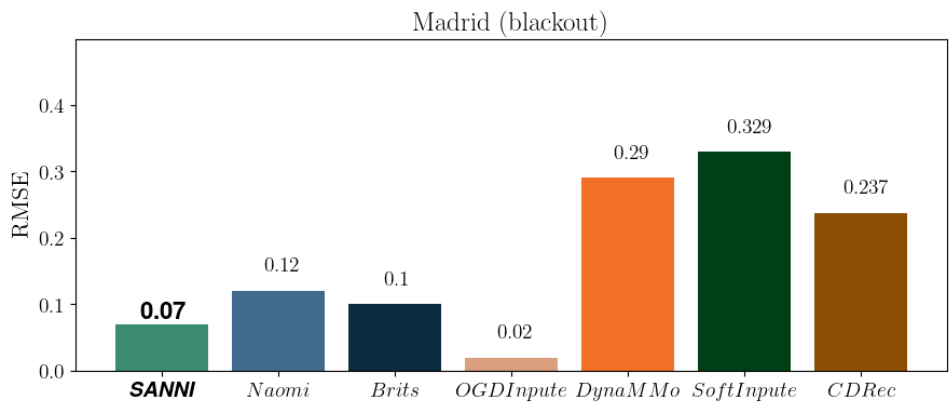
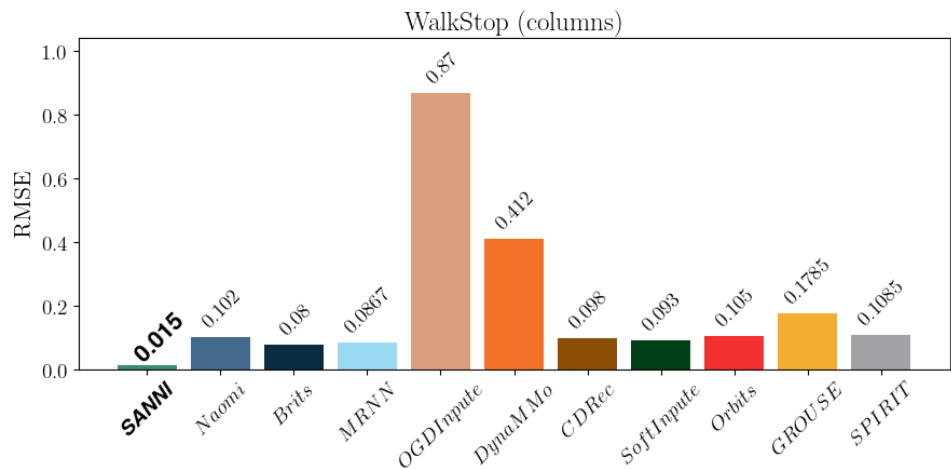
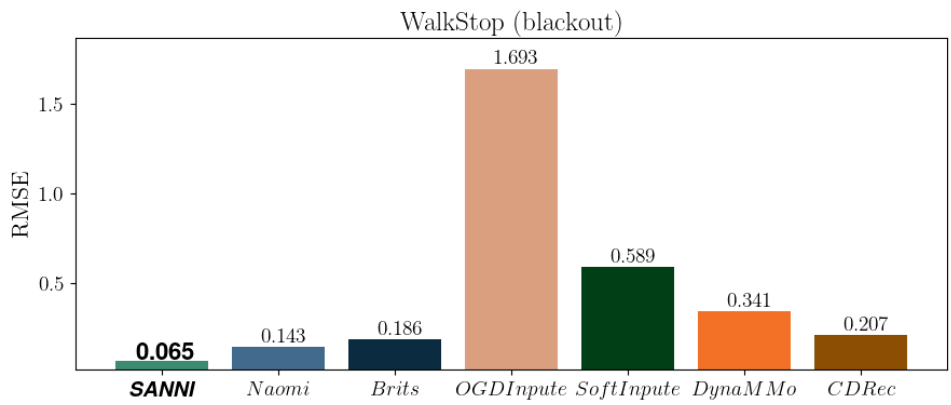
⁹ Wei C., et al. ORBITS: Online Recovery of Missing Values in Multiple Time Series Streams. 2018

¹⁰ Liu Y., et al. NAOMI: Non-Autoregressive Multiresolution Sequence Imputation. 2019

¹¹ Yoon J., et al. Estimating Missing Data in Temporal Data Streams Using Multi-Directional Recurrent Neural Networks. IEEE. 2019

Experiments: Accuracy (activities)

$$RMSE = \sqrt{\frac{1}{h} \sum_{i=1}^h (t_i - \tilde{t}_i)^2}$$

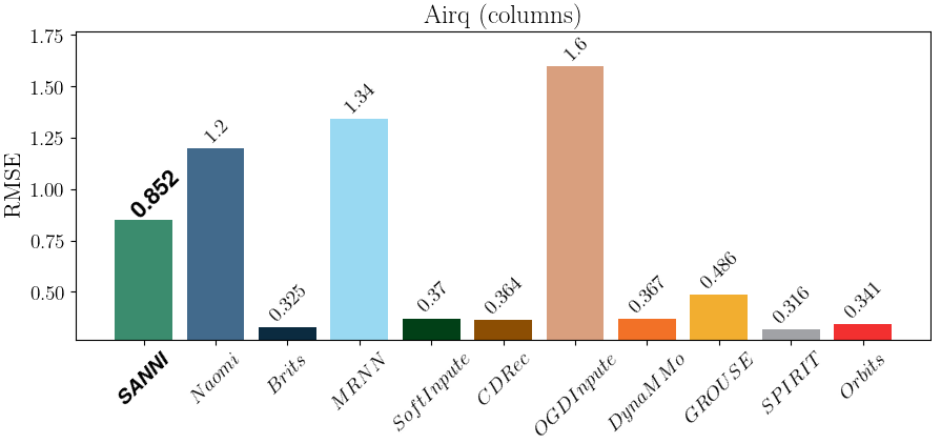
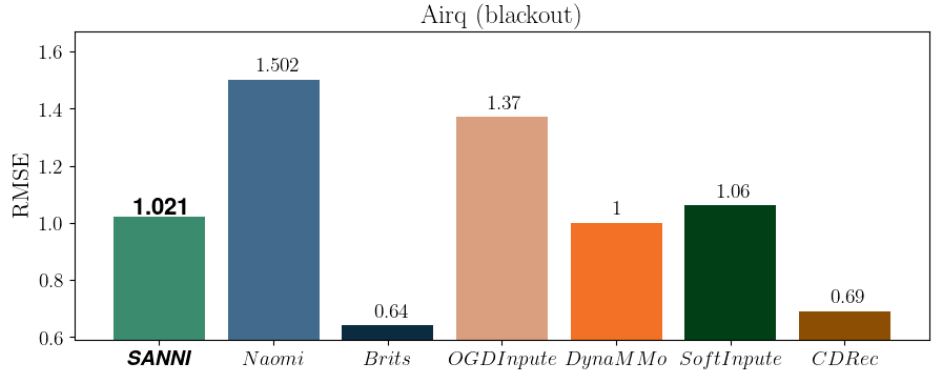
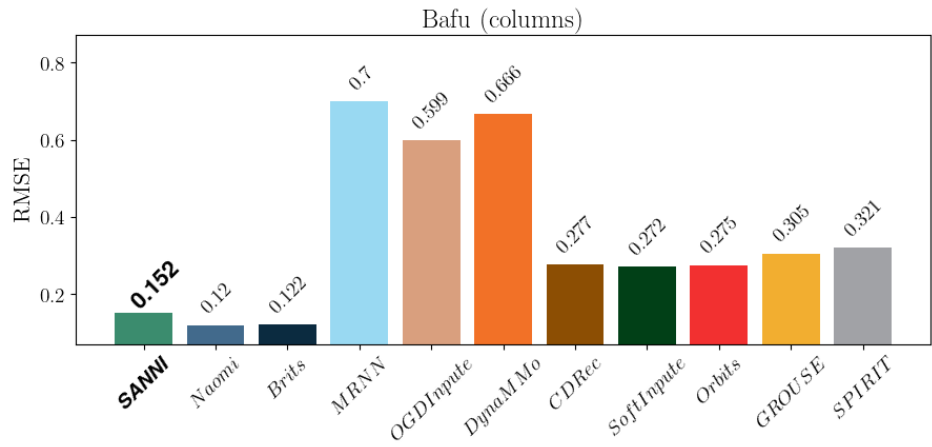
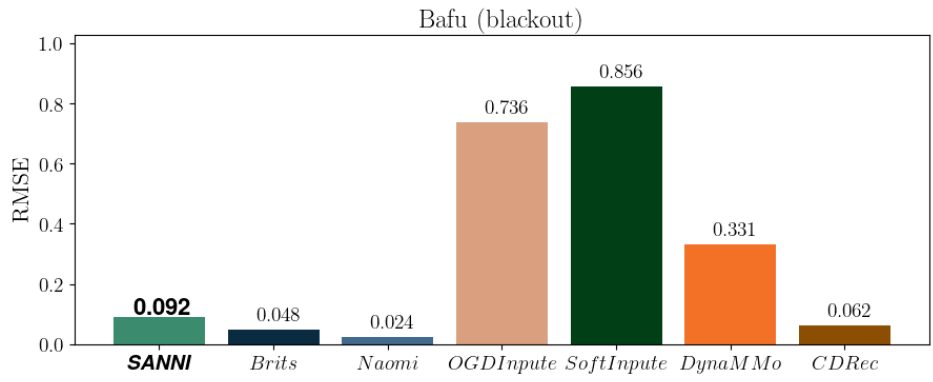


SANNI is ahead for datasets with activities

- SANNI
- Brits
- SoftInpute
- OGDInpute
- GROUSE
- Orbits
- Naomi
- MRNN
- CDRec
- DynaMMo
- SPIRIT

Experiments: Accuracy (no activities)

$$RMSE = \sqrt{\frac{1}{h} \sum_{i=1}^h (t_i - \tilde{t}_i)^2}$$



SANNI shows an **modest** result for datasets **no** activities

- SANNI
- Brits
- SoftInpute
- OGDInpute
- GROUSE
- Orbits
- Naomi
- MRNN
- CDRec
- DynaMMo
- SPIRIT

SANNI: *pro et contra*

- Pros
 - ahead of all analogs w.r.t. accuracy for the case when time series reflect activities
 - domain independence
- Cons
 - modest accuracy for the case when time series do not reflect activities

Conclusions and further research

- SANNI is novel snippet and ANN based method for imputation missing values in multivariate streaming time series that is ahead of many analogs w.r.t. accuracy for the case when measurements reflect activities
- Further research:
 - Extensive experimental evaluation of SANNI
 - Augmentation of small coverage snippet through GAN

Thank you for paying attention! Questions?
Alexey Yurtin, iurtinaa@susu.ru