

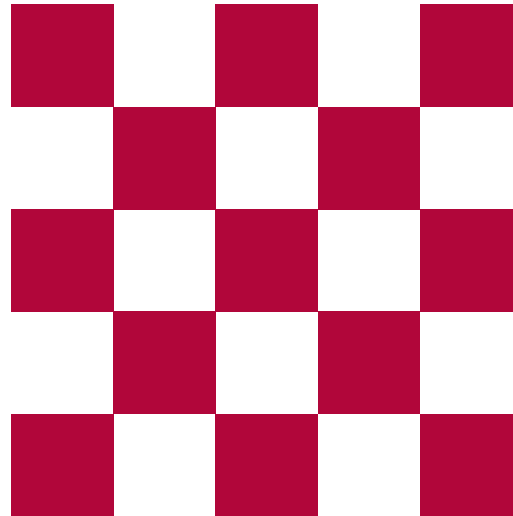
Covid Regulations

**Please, register to this event
with the Corona Warn-App!**



Large-Scale Time Series Analytics
F-E.06, Campus II, HPI

**Please, always wear your mask
and sit in a checkboard fashion.**



When you sit, you can take the
mask off.

Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 1

A large underwater photograph showing a dense school of fish swimming in a blue ocean. In the lower right foreground, the back and part of the gear of a scuba diver are visible, looking towards the school of fish.

Course: Large-Scale Time Series Analytics

Deep Dive into Time Series Anomaly Detection

Sebastian Schmidl
Phillip Wenig

Agenda

1. Organization
 - Introduction Round
 - Team Building
 - Fix Dates
2. Definitions
 - Time Series
 - Anomaly
3. Algorithms
 - Unsupervised
 - Semi-Supervised
 - Supervised
 - Preliminary Evaluation Results
4. Tools
 - GutenTAG
 - TimeEval
 - UltraMine
5. Research Example: S2G++
6. Homework

Organization

Introduction round

- Name
- Interests
 - Technologies
 - Research fields
- Hobbies?
- ...



Team Building

If you haven't found a team partner, do so now.



Fix Dates

1. Bi-weekly
 - Starting this week
 - All together
 - Every team presents current progress
 - Meeting room F-2.10 (most likely)
2. On-demand
 - Every team sends us date proposals
 - Minimum once every 2 weeks
 - More if needed

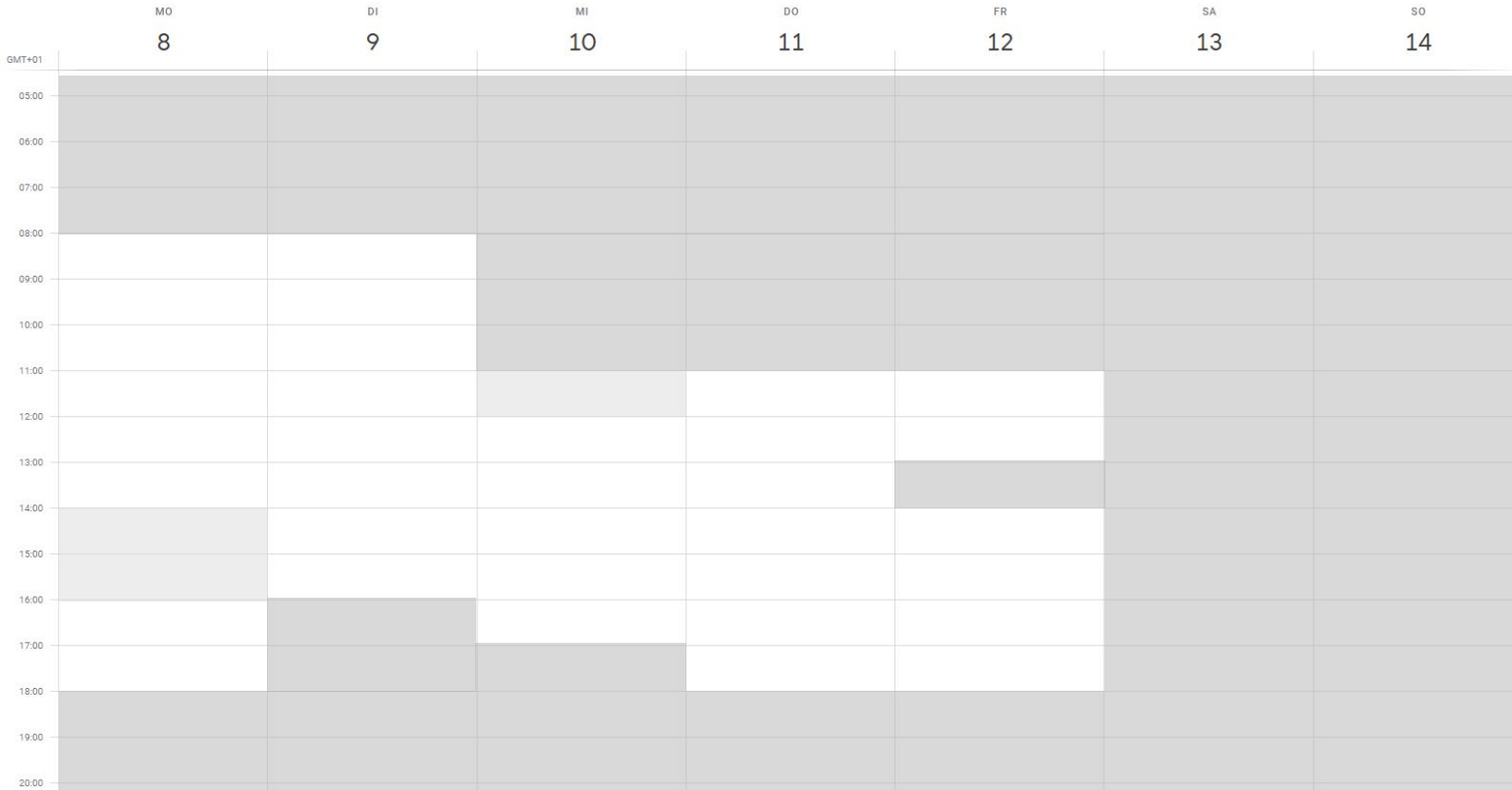


We prefer



When?

Fix Dates



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 8

Definitions

Time Series

Ordered set

**Time series =
data series**

In this study, we investigate algorithms for the detection of anomalous subsequences in time series data. A *time series* (or *data series* in general) is an ordered set $T = \{T_1, T_2, \dots, T_m\}$ of m real-valued, potentially multi-dimensional data points $T_i \in \mathbb{R}^n$. A *subsequence* $T_{i,j} = \{T_i, \dots, T_j\} \subseteq T$ is a contiguous segment of T with length $|T_{i,j}| = j - i$ and $|T_{i,j}| \geq 1$. Our evaluation assumes that the data points are equidistant, which is true for most real-world time series and relieves the algorithms from interpreting the concrete continuous measures (time, mass, angle, etc.); data series not following this assumption need to be discretized.

**Real-valued multi-
dimensional data
points**

equidistant

**Subsequence =
contiguous segment of
time series with length
 ≥ 1**

Time Series

or timestamps

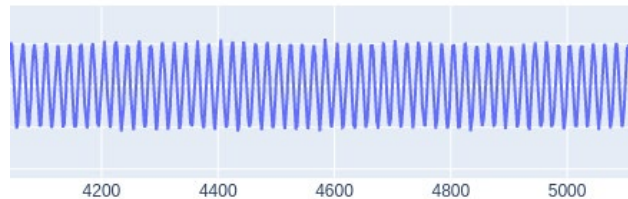
Index	value_0
0	0.199
1	0.323
2	0.339
3	0.401
4	0.356
5	0.299
6	0.200
7	0.154
...	...



```
import pandas as pd
```

```
time_series = pd.read_csv(PATH_TO_CSV)
```

```
time_series.plot()
```



Time Series

Index	value_0	value_1
0	0.199	0.199
1	0.323	0.323
2	0.339	0.339
3	0.401	0.401
4	0.356	0.356
5	0.299	0.299
6	0.200	0.200
7	0.154	0.154
...

Multivariate

- Same time
- Same subject
- But multiple (different) channels

Time Series

Index	value_0	value_1	is_anomaly
0	0.199	0.199	0
1	0.323	0.323	0
2	0.339	0.339	0
3	0.401	0.401	0
4	0.356	0.356	0
5	0.299	0.299	1
6	0.200	0.200	1
7	0.154	0.154	1
...

Labeling

- Single point labels

Or potentially overlapping subsequences

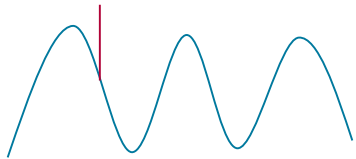
Anomaly

Definitions

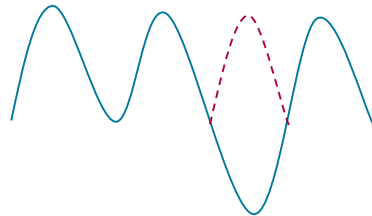
- Largest distance from rest of data points/subsequences
- Highest error in reconstructing learned behavior
- Most isolated part(s) of time series
- Part of the data set that cannot be explained by the underlying function that the data set follows.
- ...

Examples

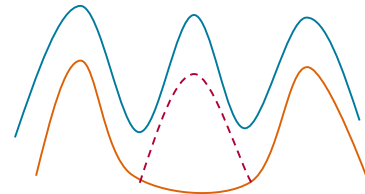
Single Point Extrema



Pattern Change



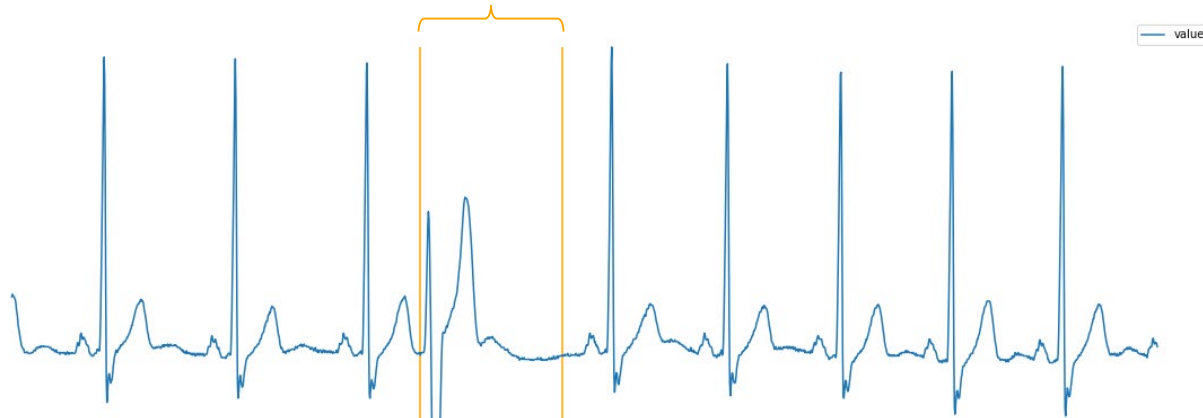
Correlation Anomaly



Anomaly

Example Dataset

Anomalous subsequence

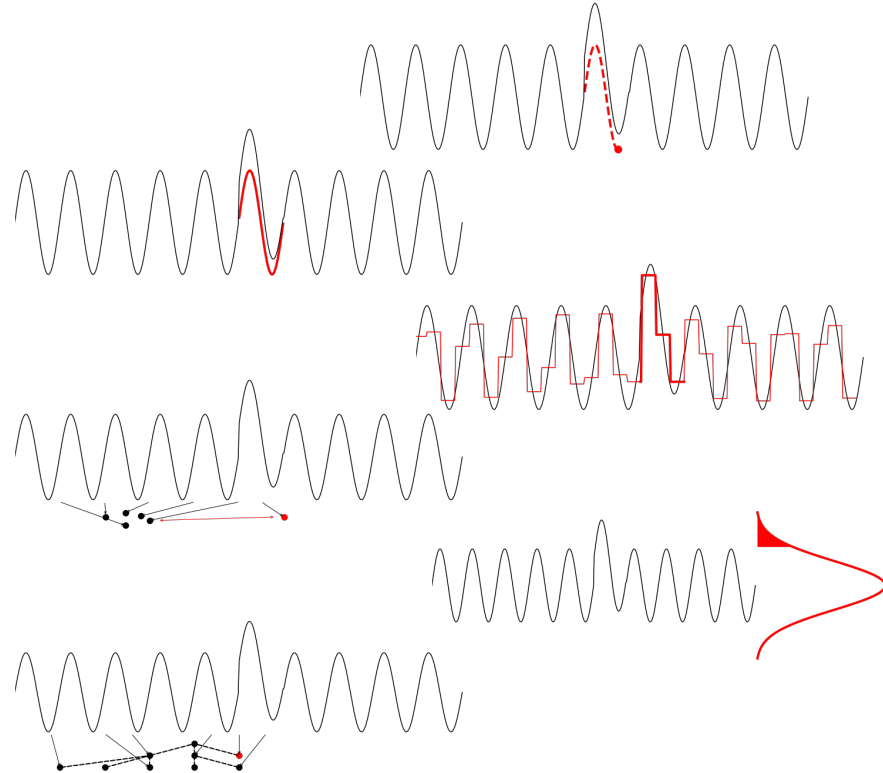


**Premature ventricular
contraction**

Algorithms

Algorithm Method Families

- Forecasting methods
- Reconstruction methods
- Encoding methods
- Distance methods
- Distribution methods
- Isolation tree methods



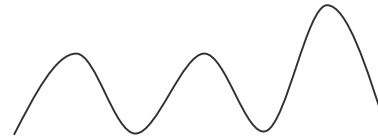
Algorithm Learning Types

Train

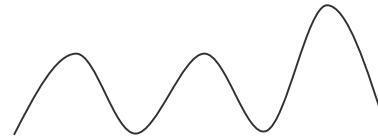
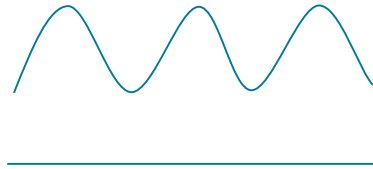
Test

- Unsupervised

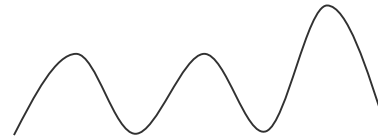
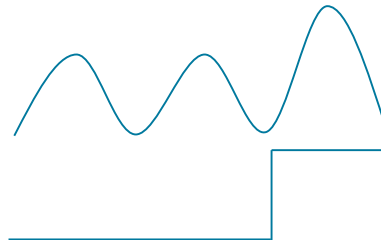
\emptyset



- Semi-supervised



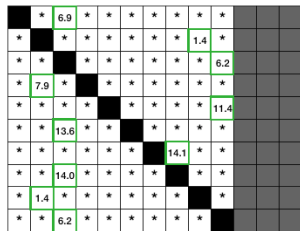
- Supervised



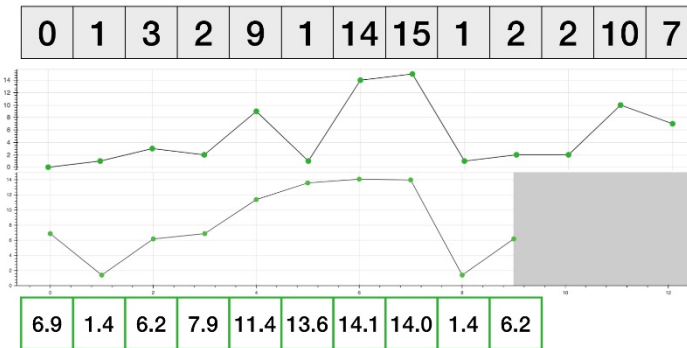
Unsupervised Algorithm Example

Matrix Profile

Matrix Profile



#DistanceProfiles

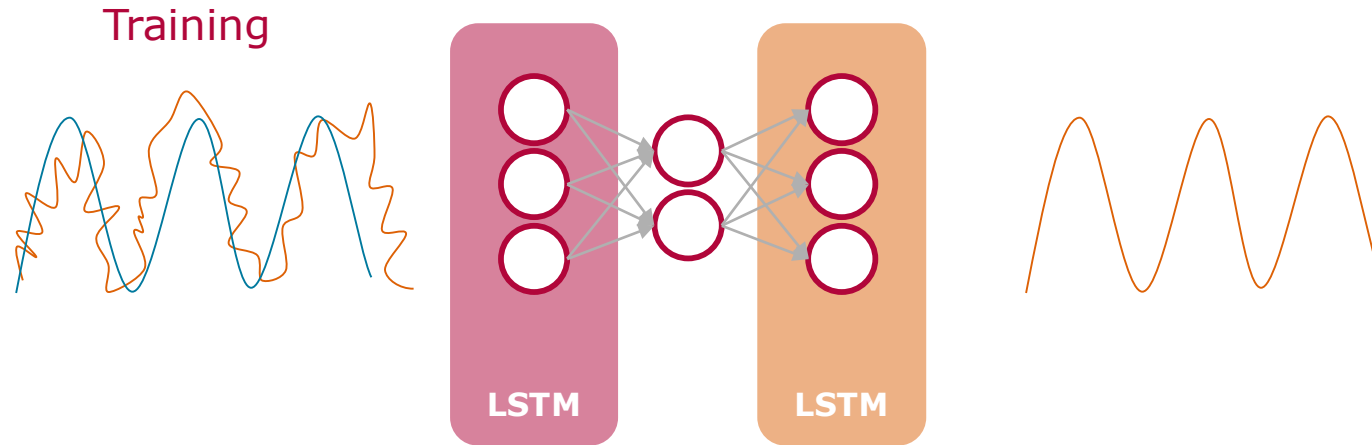


#MatrixProfileAnnotation

Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 20

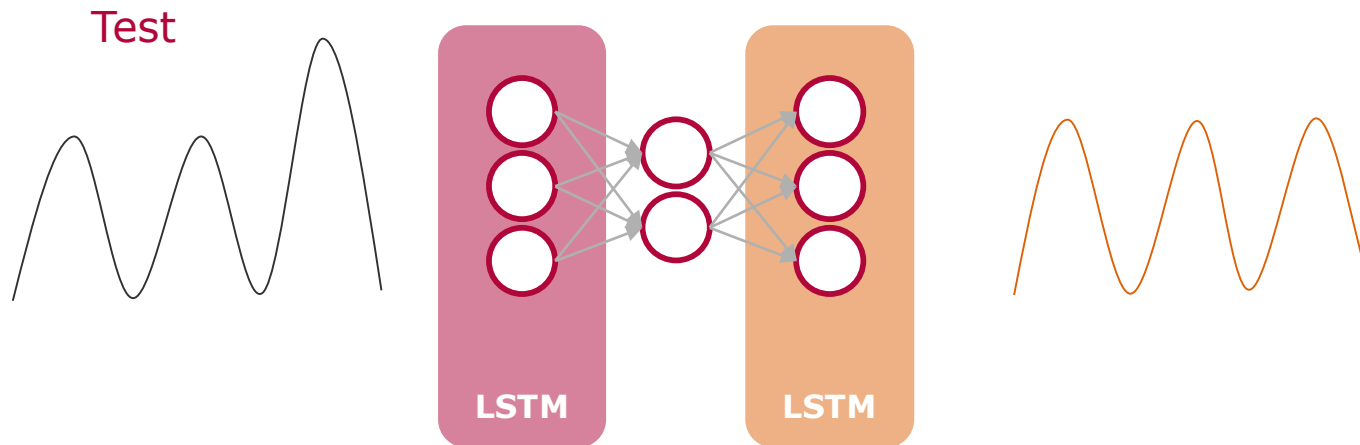
Semi-Supervised Algorithm Example

- Encoder Decoder Anomaly Detector (EncDec-AD)



Semi-Supervised Algorithm Example

- Encoder Decoder Anomaly Detector (EncDec-AD)

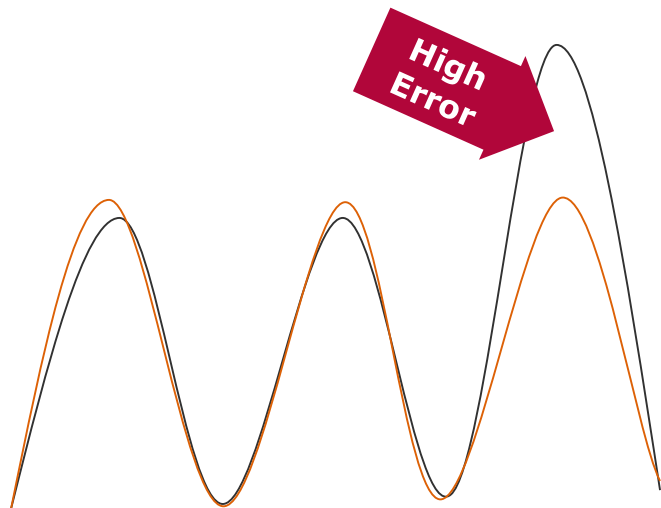


Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 22

Semi-Supervised Algorithm Example

- Encoder Decoder Anomaly Detector (EncDec-AD)

Test



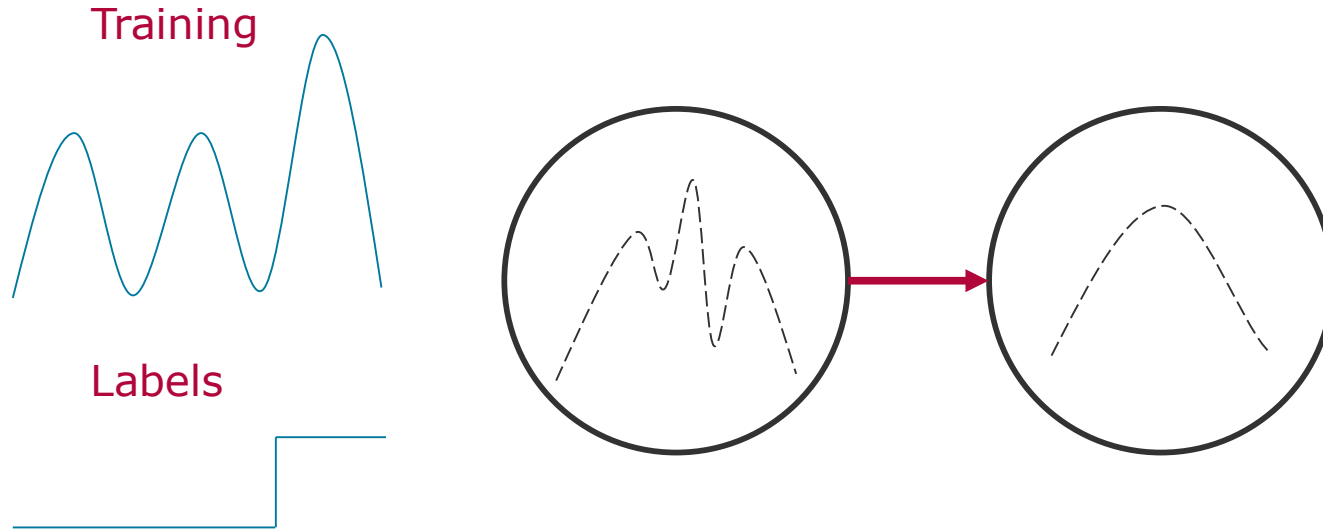
Anomaly Score



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 23

Supervised Algorithm Example

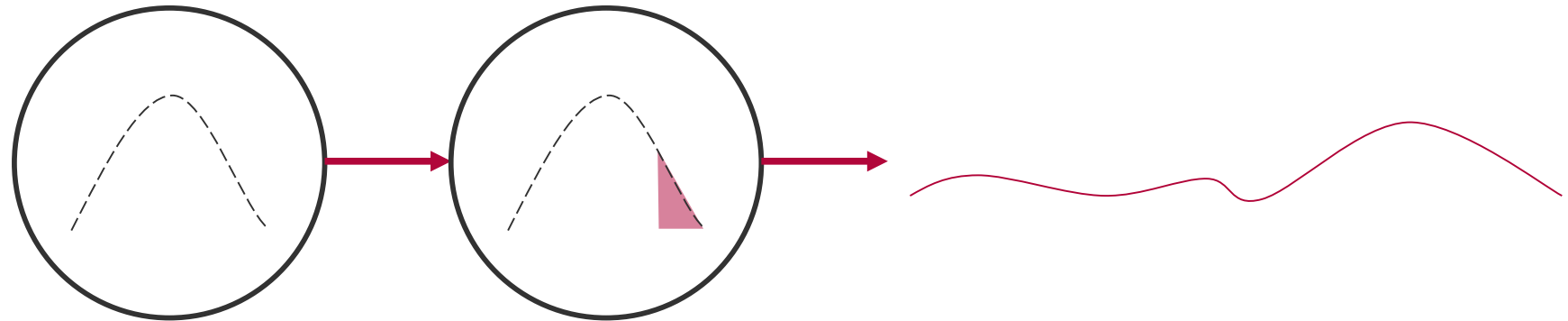
- Normalizing Flows



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 24

Supervised Algorithm Example

- Normalizing Flows



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 25

Preliminary Evaluation Results



Tools

GutenTAG

A Good Timeseries Anomaly Generator



timeseries:

- name: **demo**
- length: **10000**
- channels: **1**

- semi-supervised: **true**
- supervised: **true**

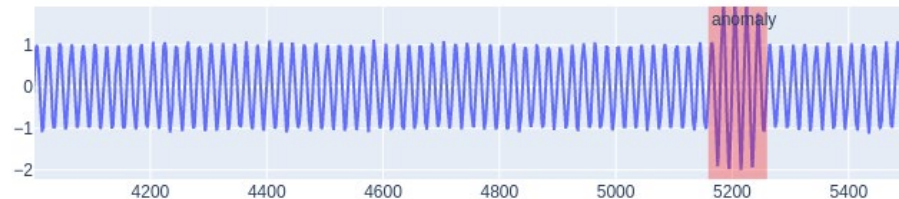
- base-oscillation:
 - kind: **sinus**
 - frequency: **10.0**
 - amplitude: **1.0**
 - variance: **0.05**

anomalies:

- position: **end**
- length: **100**
- kinds:

- kind: **amplitude**
 - parameters:
 - amplitude_factor: **2.0**

- position: *beginning*



Output format

Time Series Properties

Anomalies Properties

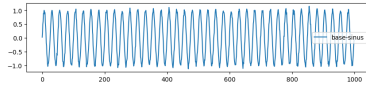
Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart **28**

GutenTAG

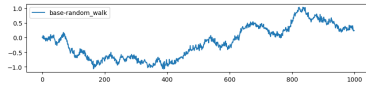
A Good Timeseries Anomaly Generator



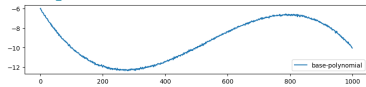
Sine wave



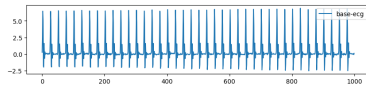
Random walk



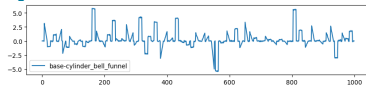
Polynomial



ECG

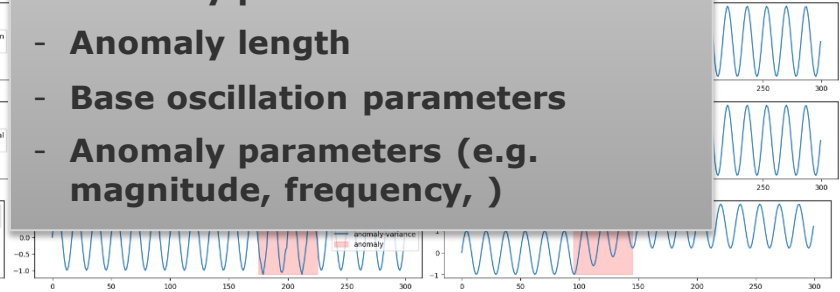
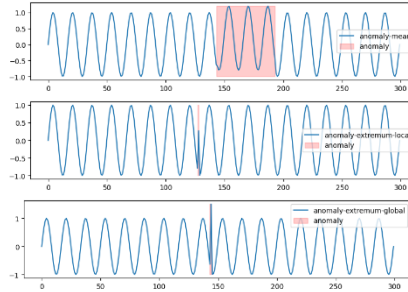


Cylinder-bell-funnel



Variation of

- Noise
- Trend
- #dimensions
- #anomalies (similar or different)
- Anomaly position
- Anomaly length
- Base oscillation parameters
- Anomaly parameters (e.g. magnitude, frequency,)



Anomaly types

Chart 29

GutenTAG

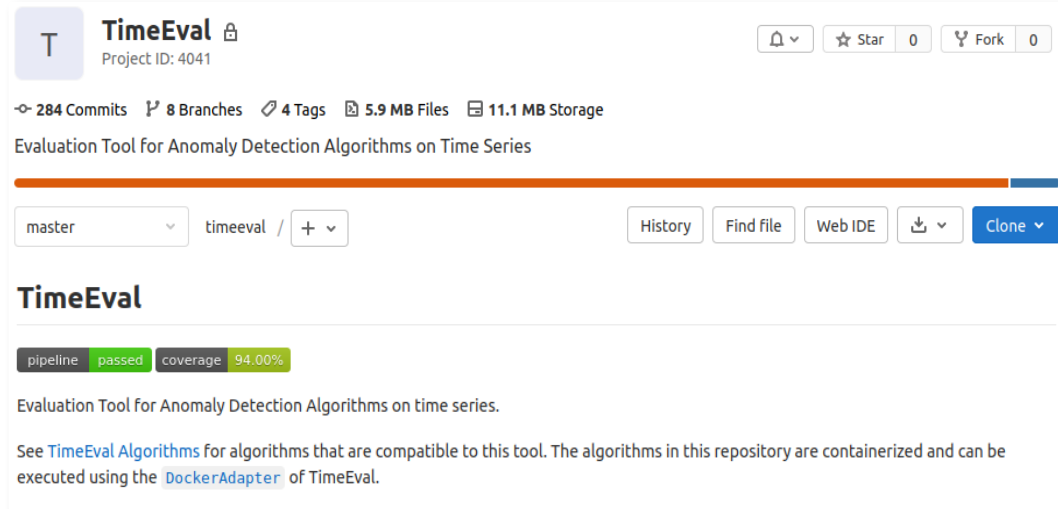
A Good Timeseries Anomaly Generator



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart **30**

TimeEval

- Integrated benchmark datasets
- **Adapter** architecture for algorithm integration (FunctionAdapter, DockerAdapter, ...)
- Automatic **quality scoring using different metrics**
- **Automatic timing** of the algorithm execution
- **Distributed** experiment execution
- Output and logfile tracking for subsequent inspection / debugging

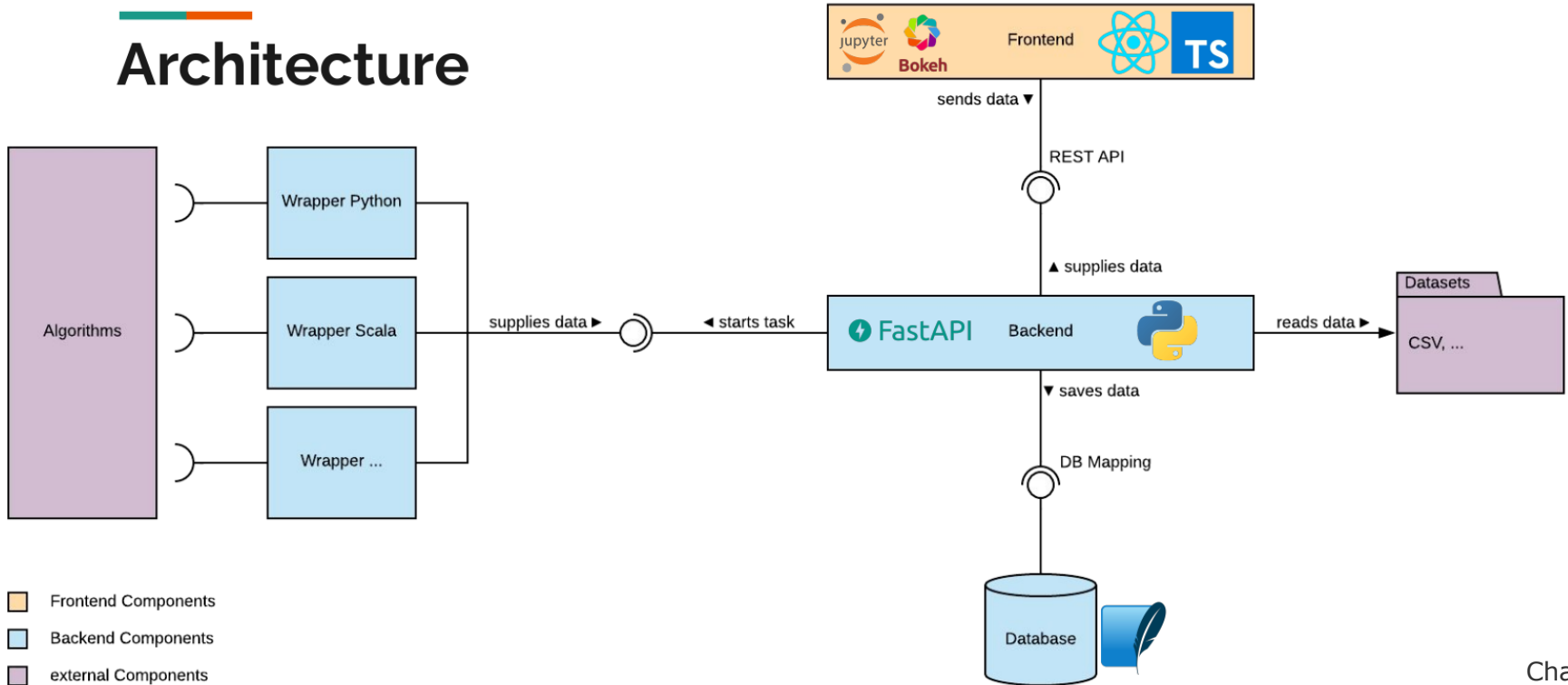


The screenshot shows the GitLab repository page for TimeEval. At the top, it displays the repository name 'TimeEval' with a lock icon and 'Project ID: 4041'. To the right are icons for notifications, stars (0), and forks (0). Below this, it shows repository statistics: 284 Commits, 8 Branches, 4 Tags, 5.9 MB Files, and 11.1 MB Storage. The description reads 'Evaluation Tool for Anomaly Detection Algorithms on Time Series'. There is a branch selector set to 'master' and a '+', along with buttons for 'History', 'Find file', 'Web IDE', a download icon, and 'Clone'. The main content area shows a 'TimeEval' heading, a status bar with 'pipeline passed' and 'coverage 94.00%', and a description: 'Evaluation Tool for Anomaly Detection Algorithms on time series.' It also includes a link to 'TimeEval Algorithms' and mentions that algorithms are containerized and can be executed using the 'DockerAdapter'.





Architecture





Research Example: DADS

(more efficient)

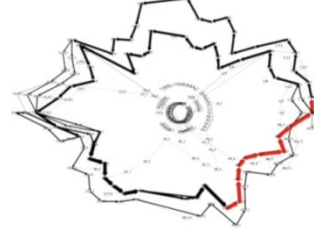
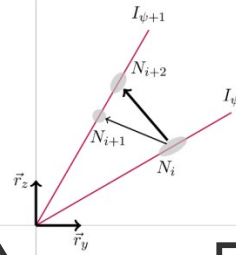
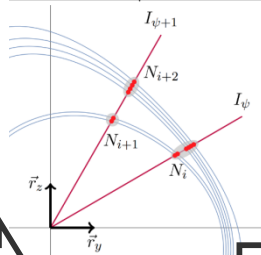
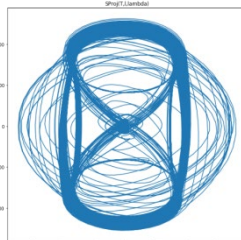
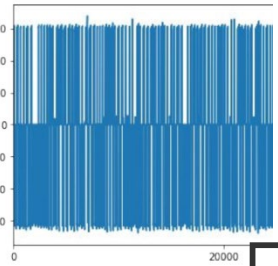
Series2Graph

Subsequence
Embedding

Intersection
Calculation

Edge
Extraction

Subsequence
Scoring



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 36

Distributed Detection of Sequential Anomalies in Univariate Time Series

Johannes Schneider · Phillip Wenig · Thorsten Papenbrock

Received: date / Accepted: date

Abstract The automated detection of sequential anomalies in time series is an essential task for many applications, such as the monitoring of technical systems, fraud detection in high-frequency trading, or the early detection of disease symptoms. All these applications require the detection to find all sequential anomalies possibly first on potentially very large time series. In other words, the detection needs to be effective, efficient and scalable w.r.t. the input size. Series2Graph is an effective solution based on Graph-embeddings that is robust against re-occurring anomalies, can discover sequential anomalies of arbitrary length and works without training data. Yet, Series2Graph is neither efficient nor scalable due to its single-threaded approach: it can, in particular, not process arbitrarily large sequences due to the memory constraints of a single machine.

In this paper, we propose our Distributed Anomaly Detection System, short DADS, which is an efficient and scalable adaptation of Series2Graph. Based on the actor programming model, DADS distributes the input time sequence, intermediate state and the computation to all processes of a cluster in a way that minimizes communication costs and synchronization barriers. Our evaluation shows that DADS is orders of magnitude faster than S2G, scales almost linearly with the number of processes in the cluster and can process much larger input sequences due to its scale-out property.

J. Schneider
Hasso Plattner Institute, University of Potsdam
E-mail: johannes.schneider@hpi.de

P. Wenig
Hasso Plattner Institute, University of Potsdam
E-mail: philipp.wenig@hpi.de

T. Papenbrock
Hasso Plattner Institute, University of Potsdam
E-mail: thorsten.papenbrock@hpi.de

Keywords distributed programming, sequential anomaly, actor model, data mining, time series

1 Sequential Anomaly Detection

Univariate time series are ordered sequences of one-dimensional, real-valued records [1, 18, 31, 38, 41]. The ordering is usually time-related – hence the name – but sometimes also follows other continuous measures, such as size, distance, or speed. Anomalies in time series denote subsequences with significant, i.e., particularly subsequent values or patterns. In general, we distinguish between point anomalies s.a.k.a. outliers and sequential anomalies s.a.k.a. collective anomalies [8]. Because point anomalies are sequential anomalies of size one, sequential anomalies cover point anomalies. For example, a temperature reading of 50°C in an ordinary room is a point anomaly; a sudden temperature increase of 10°C directly followed by a drastic drop of 15°C is a sequential anomaly, if common temperature changes are slow.

Anomalies often indicate meaningful events in time series and are, therefore, important for various applications. Examples for such applications range from intrusion detection in digital networks [17, 28], over finding unexpected phenomena in anthropological measurements [27], to medical studies of disease cases [15]. Figure 1 shows another example from aircraft engineering: We see a time series T of engine measurements with various regular spikes and one anomalously keen spike. None of the records in the anomaly is an outlier, but its shape is clearly different from other patterns that we find in the sequence. An anomaly score, which is depicted in red in Figure 1, is a proper indicator of this anomaly. In this paper, we propose a scalable algorithm that efficiently calculates such anomaly scores.

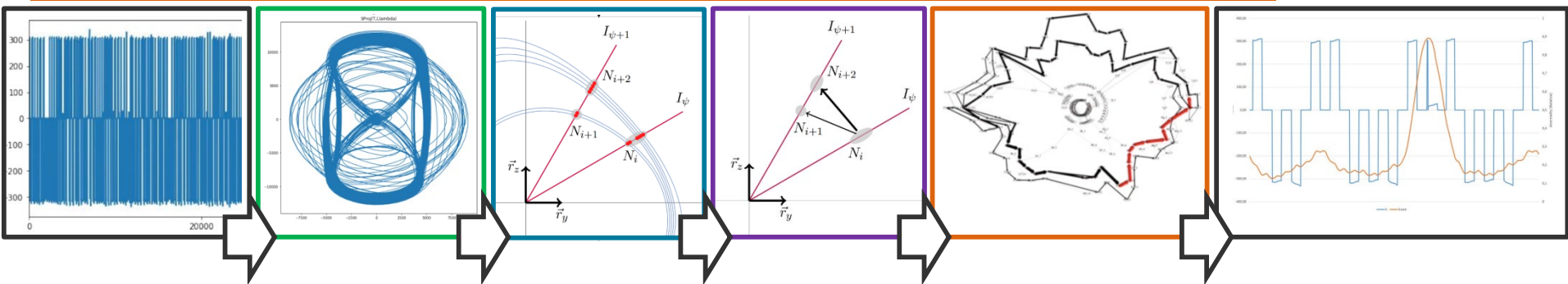
Sequential anomalies are hard to detect, because they often differ in length, vary in degree of anomaly and might

- **Distributed Anomaly Detection System (DADS)**
- **Distributed implementation of the Series2Graph (S2G) algorithm**
- **Distributes computation as well as temporary data**
- **Reactive / minimizes communication barriers**
- **Implemented using Akka**



Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 37

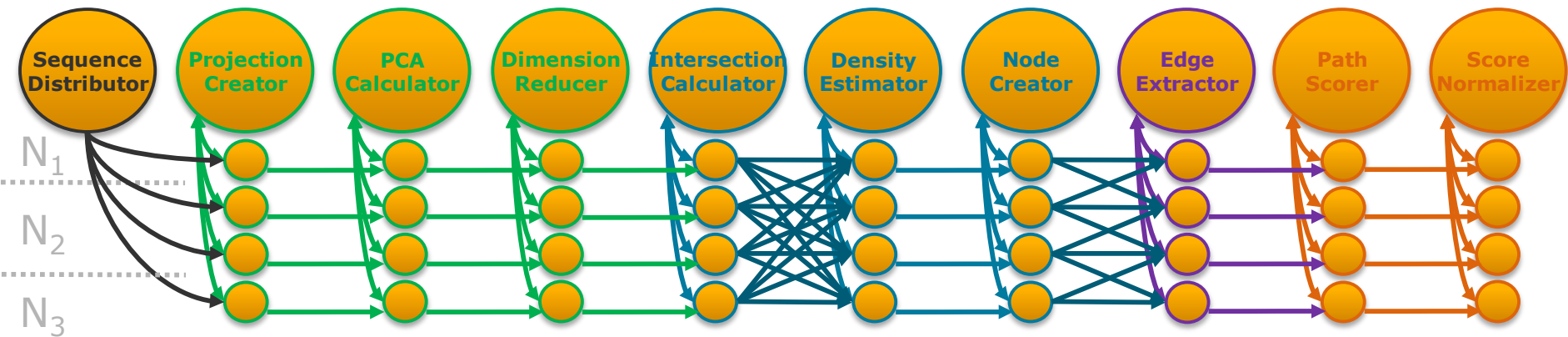
DADS



DADS: Subsequence Embedding

Node Extraction

Edge Extraction Subsequence Scoring



DADS



 T ($\times 10^6$)	S2G	S2G⁺	DADS (1P,1T)	DADS (1P,20T)	DADS (12P,20T)
0.01	4	4	6	5	7
0.1	35	29	12	7	7
1	388	297	145	19	9
10	5,897	3,077	904	145	26
50	TL	15,947	4,575	727	106
100	ML	ML	ML	ML	206
200	ML	ML	ML	ML	401
300	ML	ML	ML	ML	586
400	ML	ML	ML	ML	802
500	ML	ML	ML	ML	986

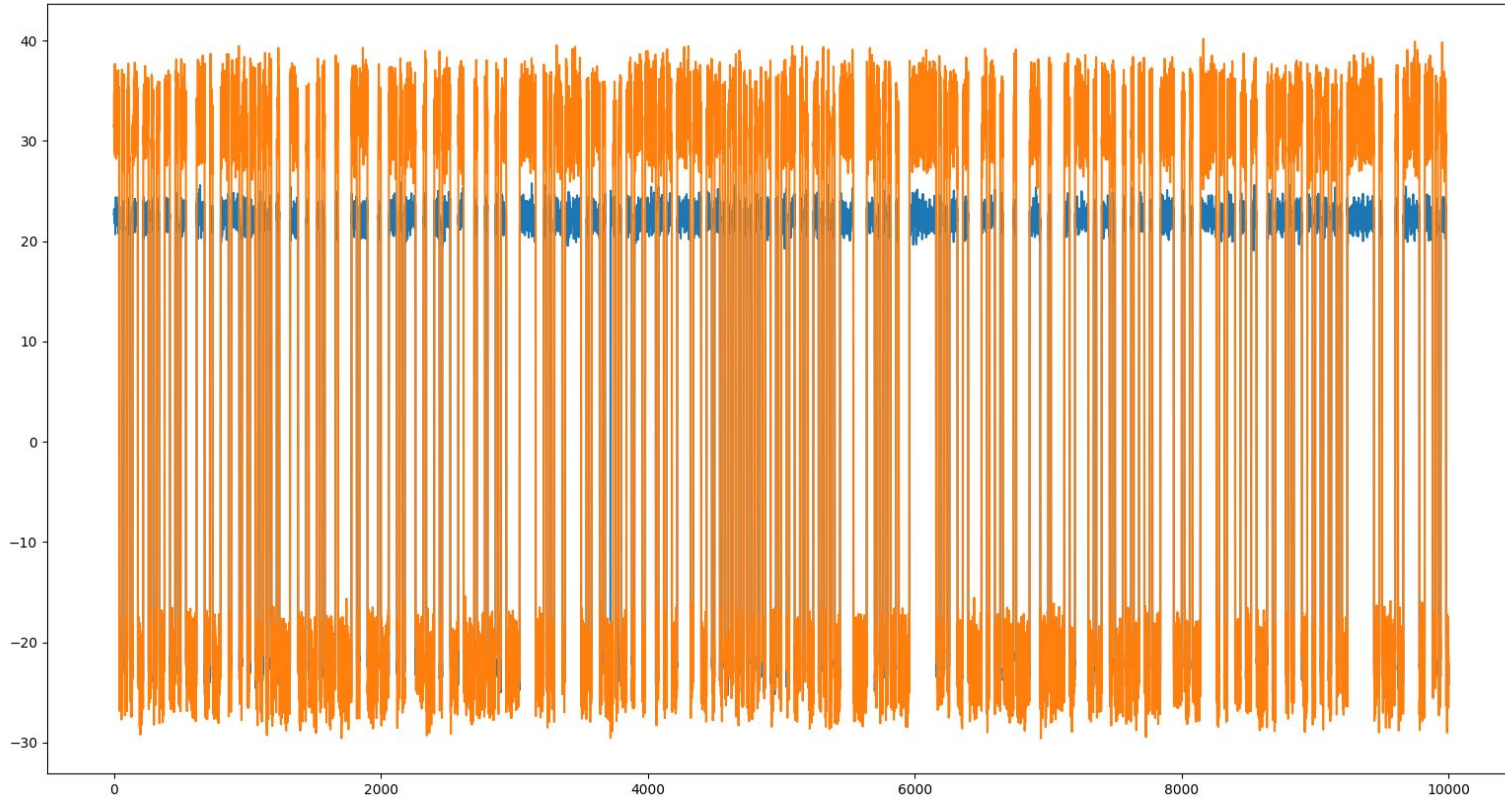
runtime in seconds

Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart **39**

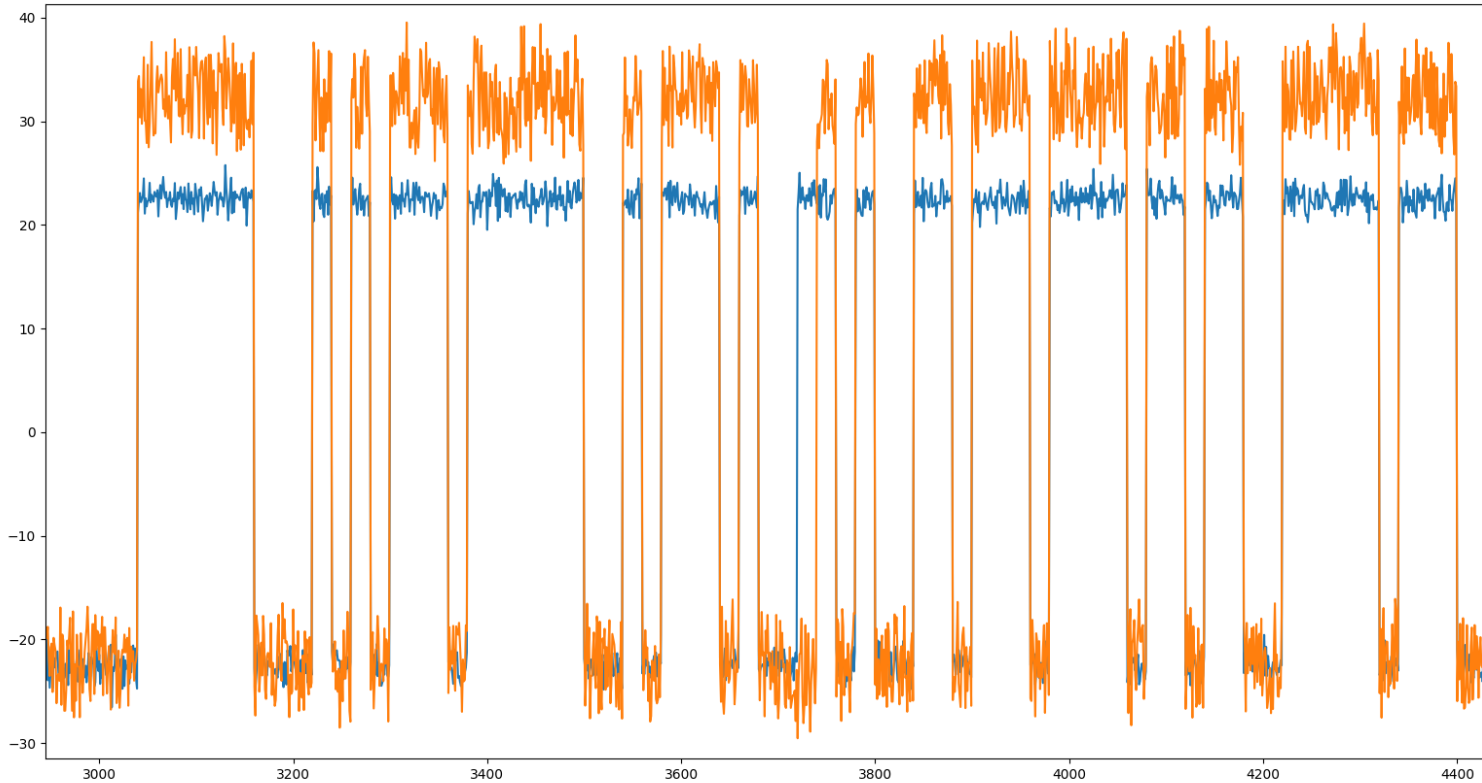
Research Example: S2G++

(more capable; work in progress)

Series2Graph

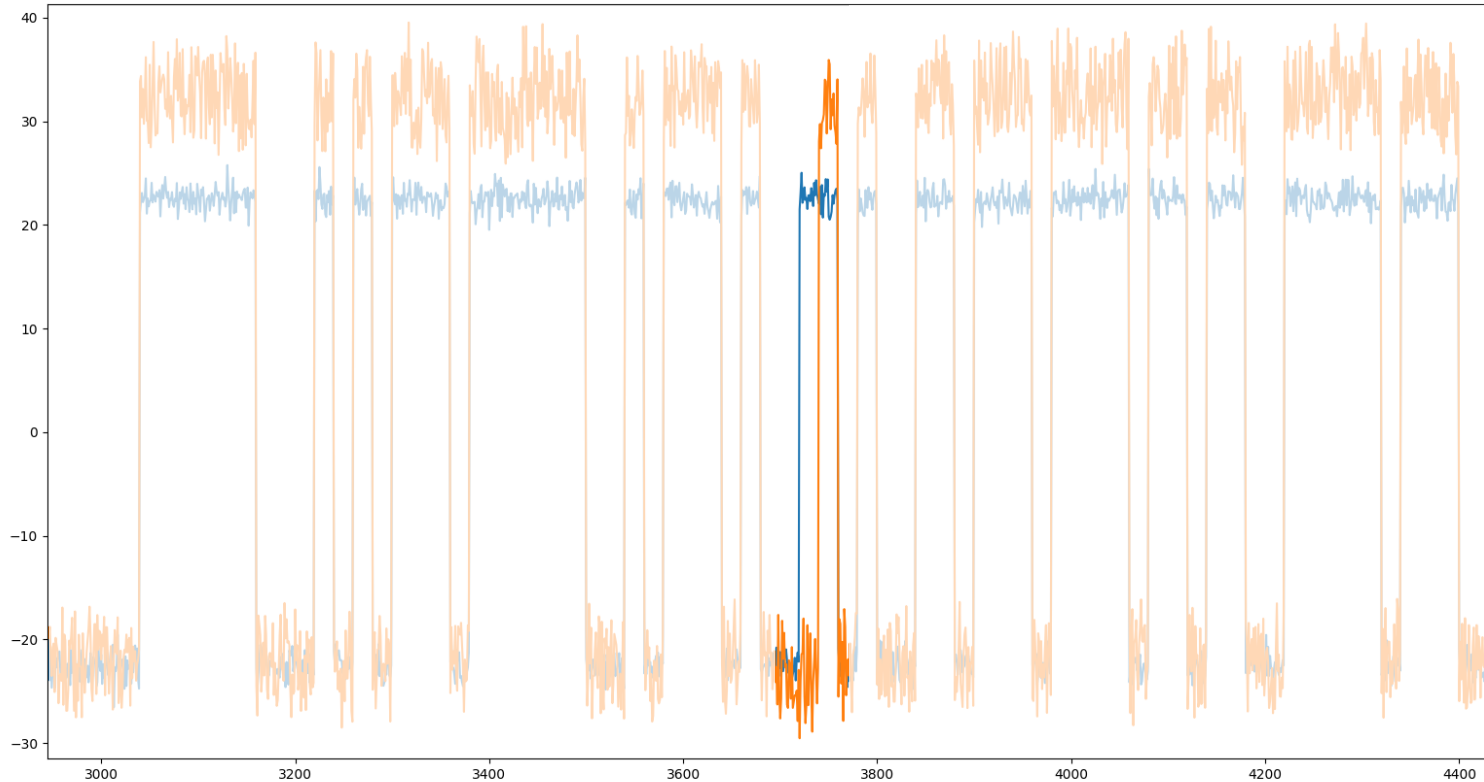


Series2Graph



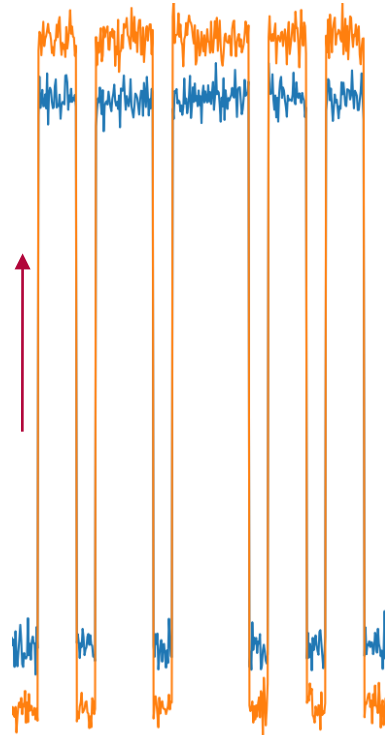
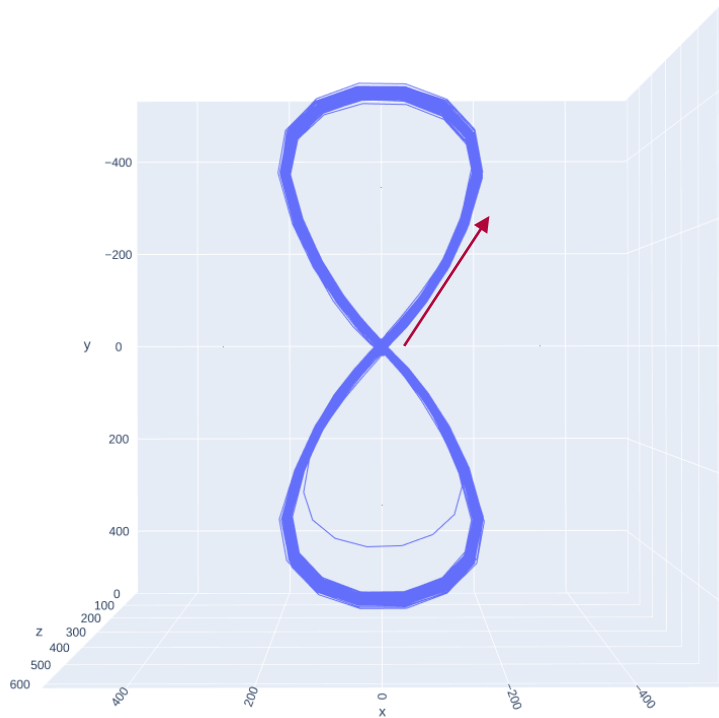
Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart **42**

Series2Graph

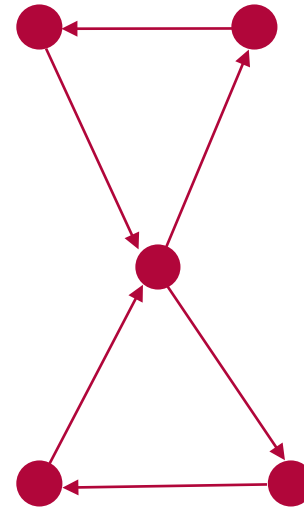
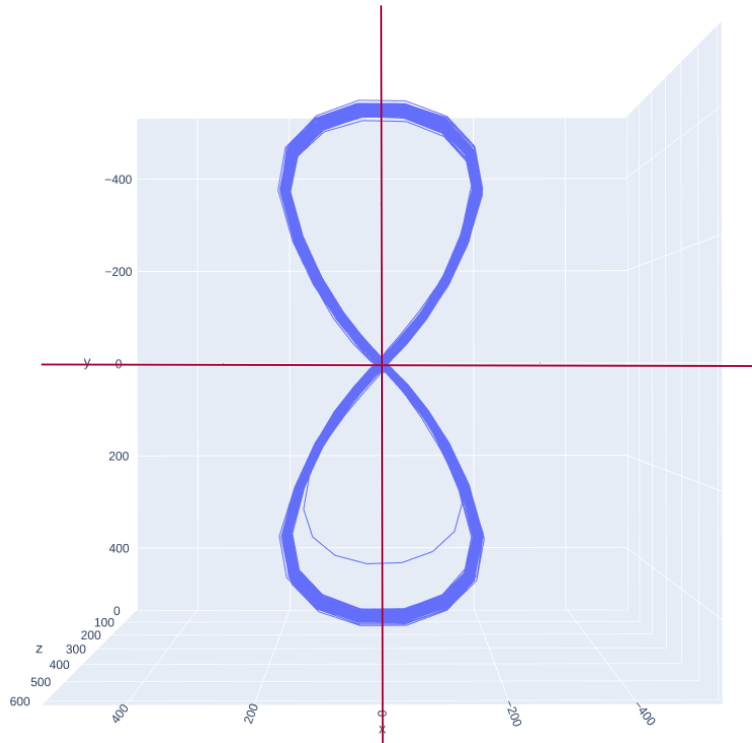


Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart **43**

Series2Graph

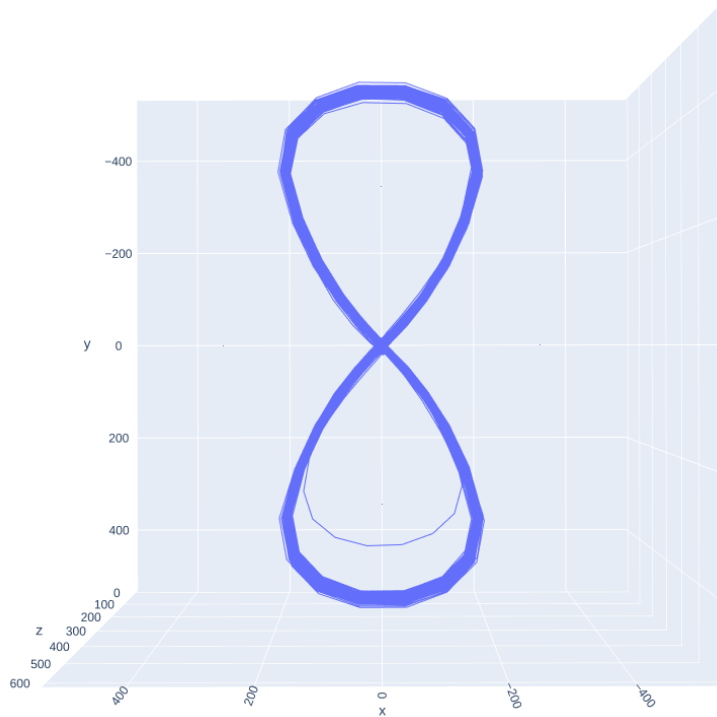


Series2Graph

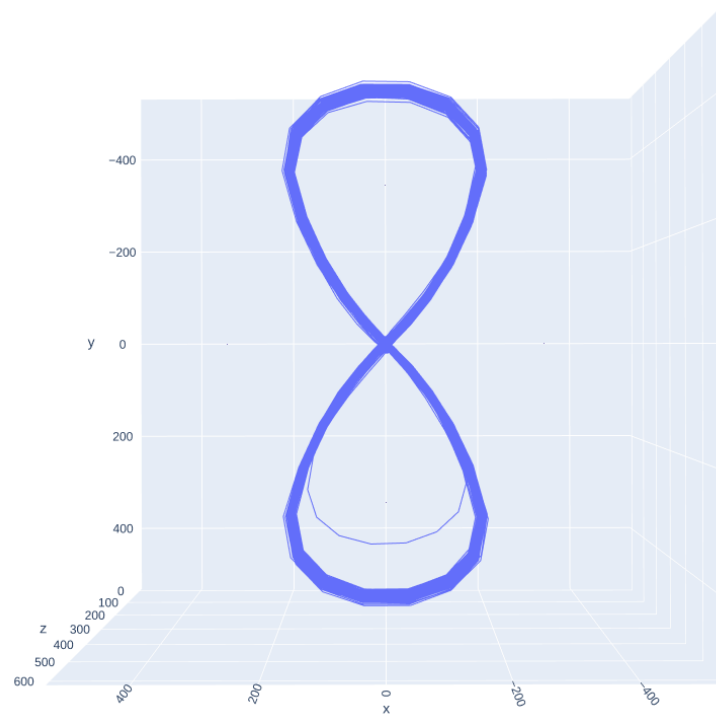


Series2Graph++

Time Series 1

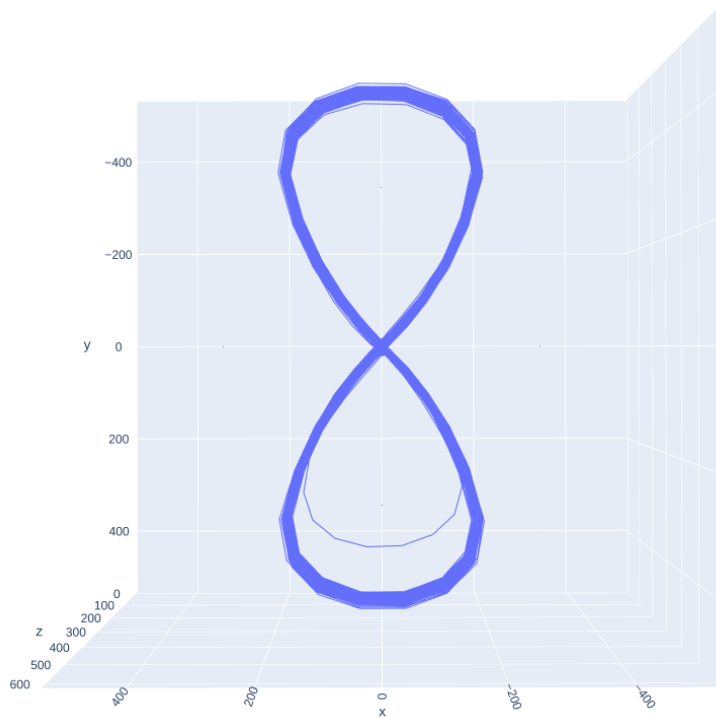


Time Series 2

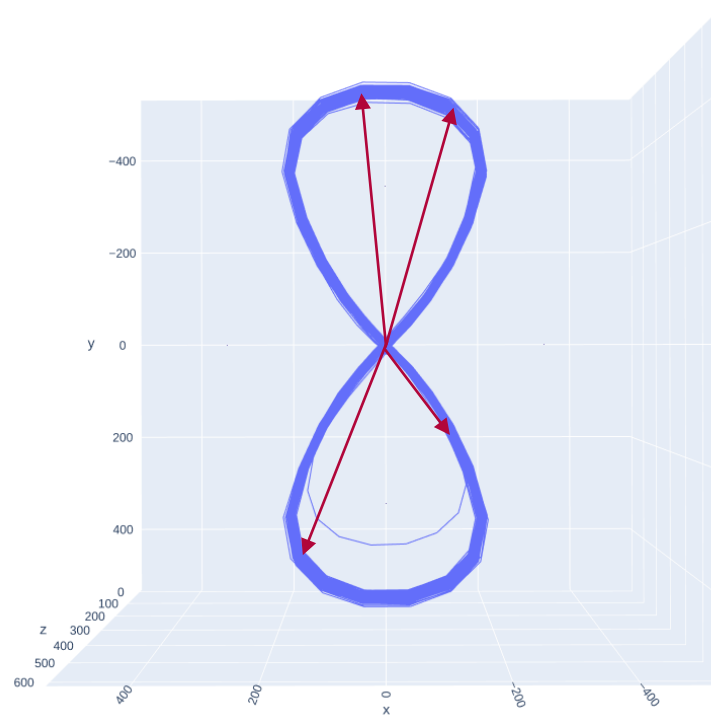


Series2Graph++

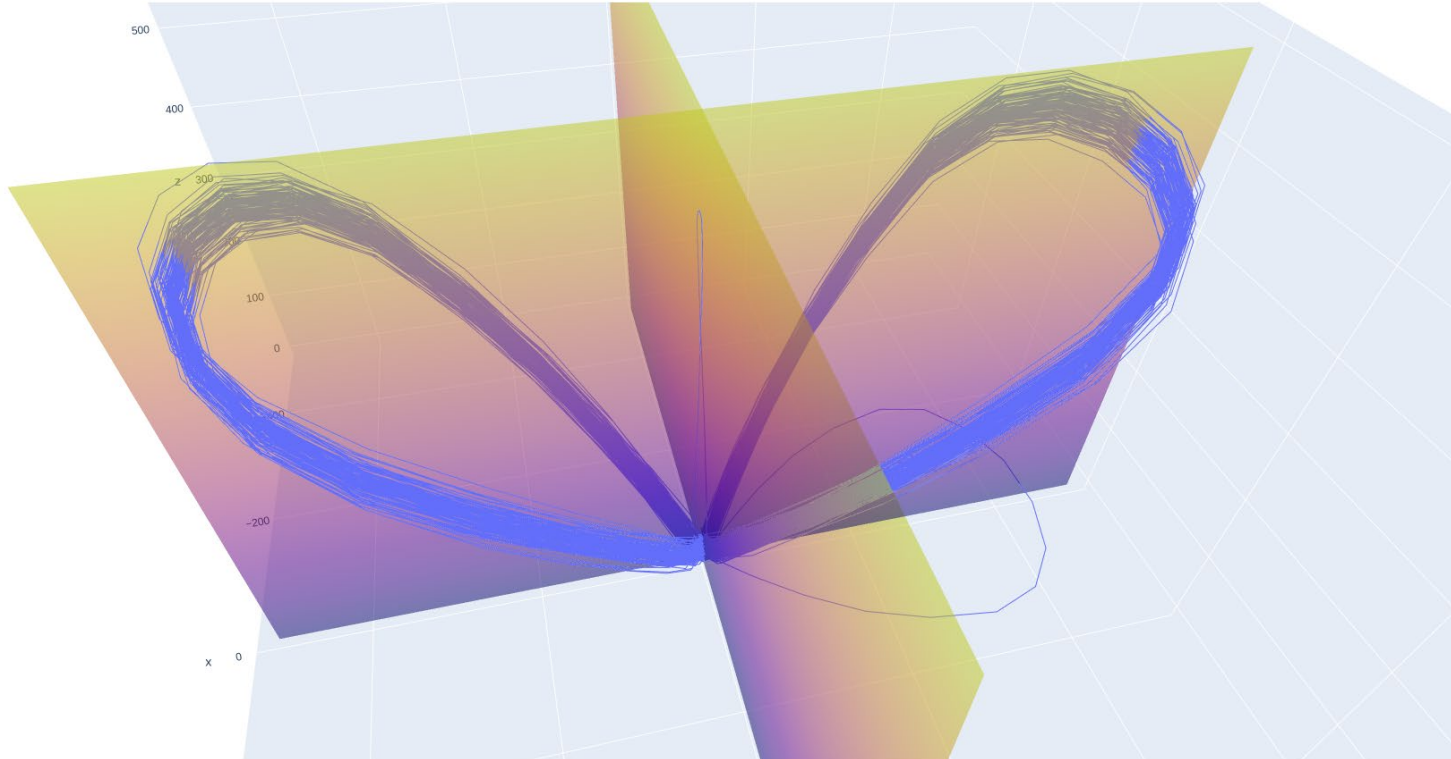
Time Series 1



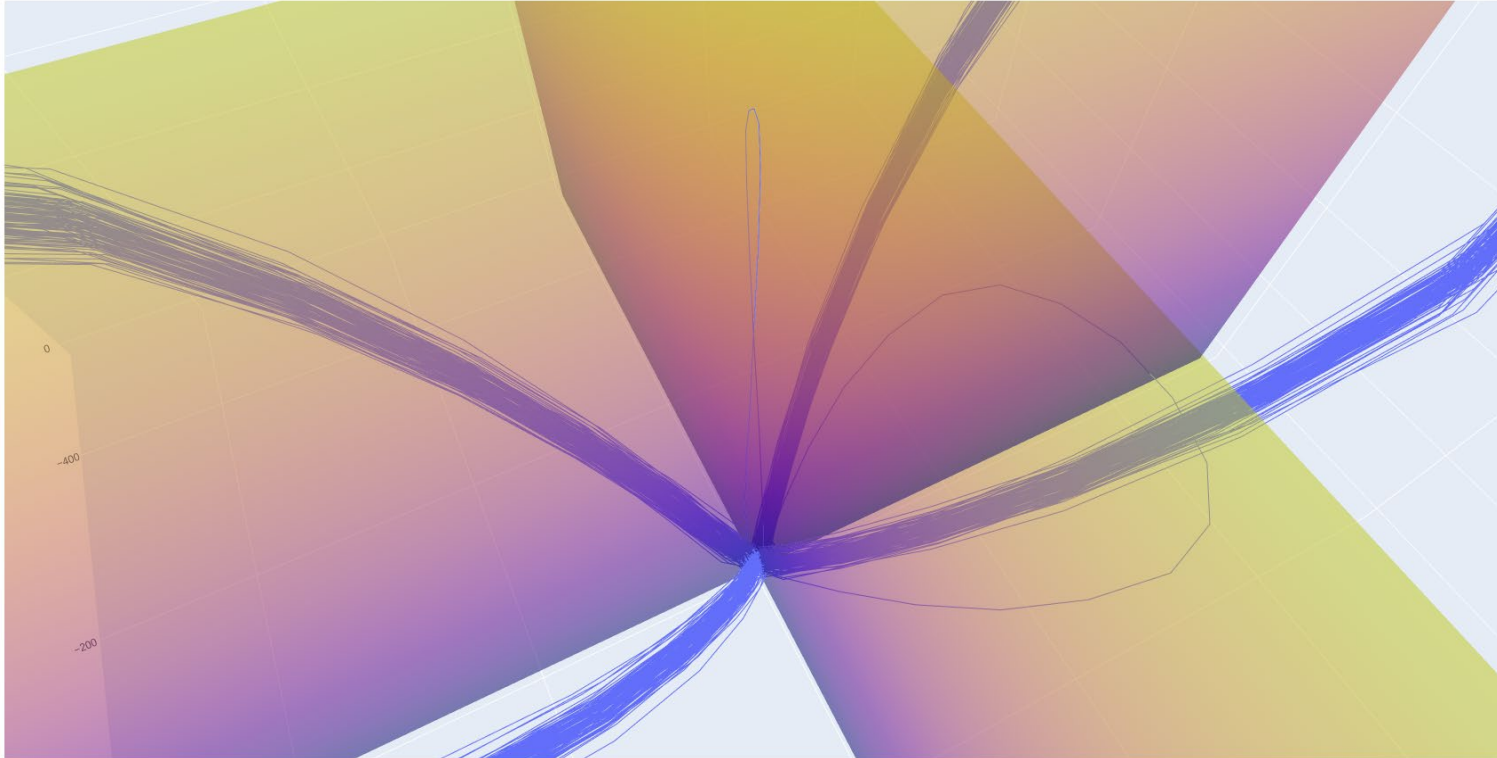
Time Series 2



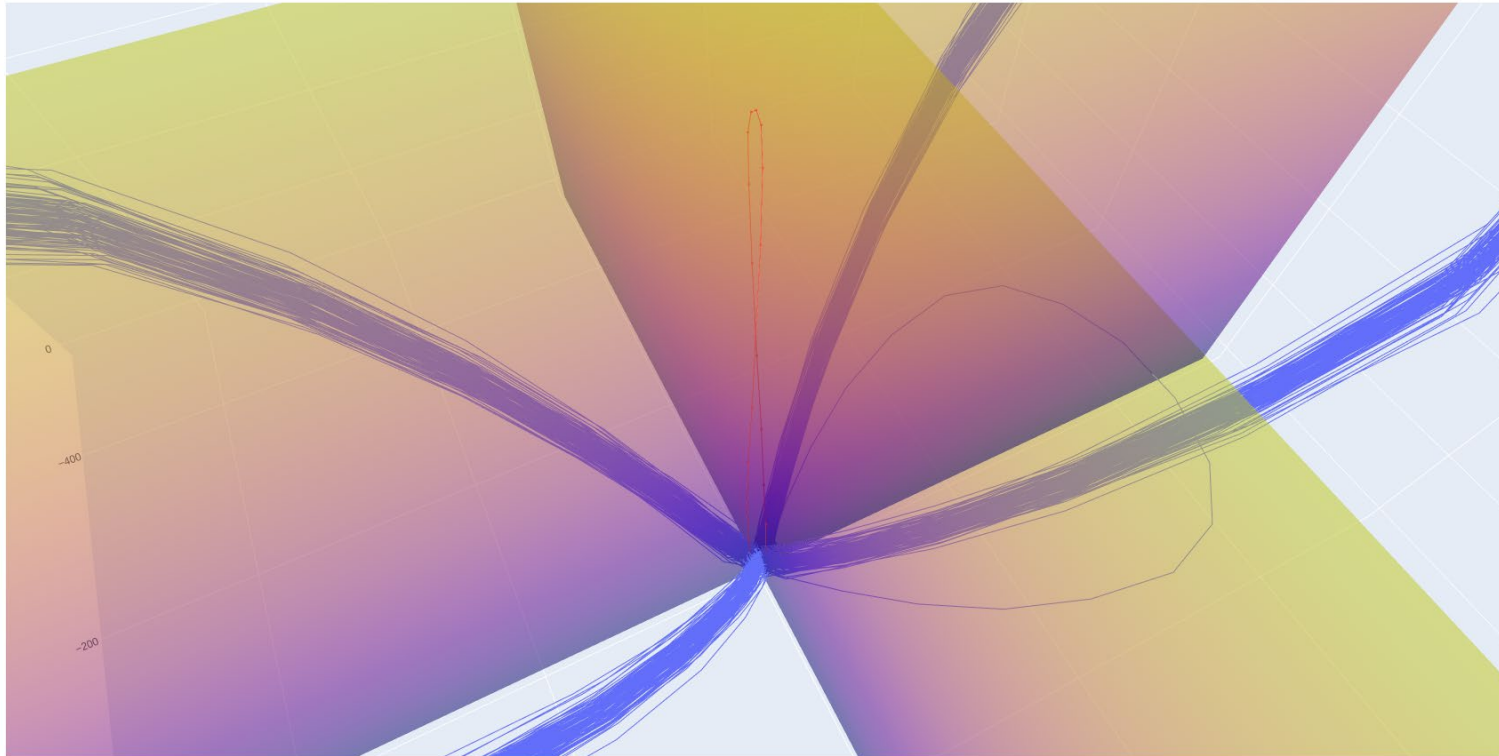
Series2Graph++



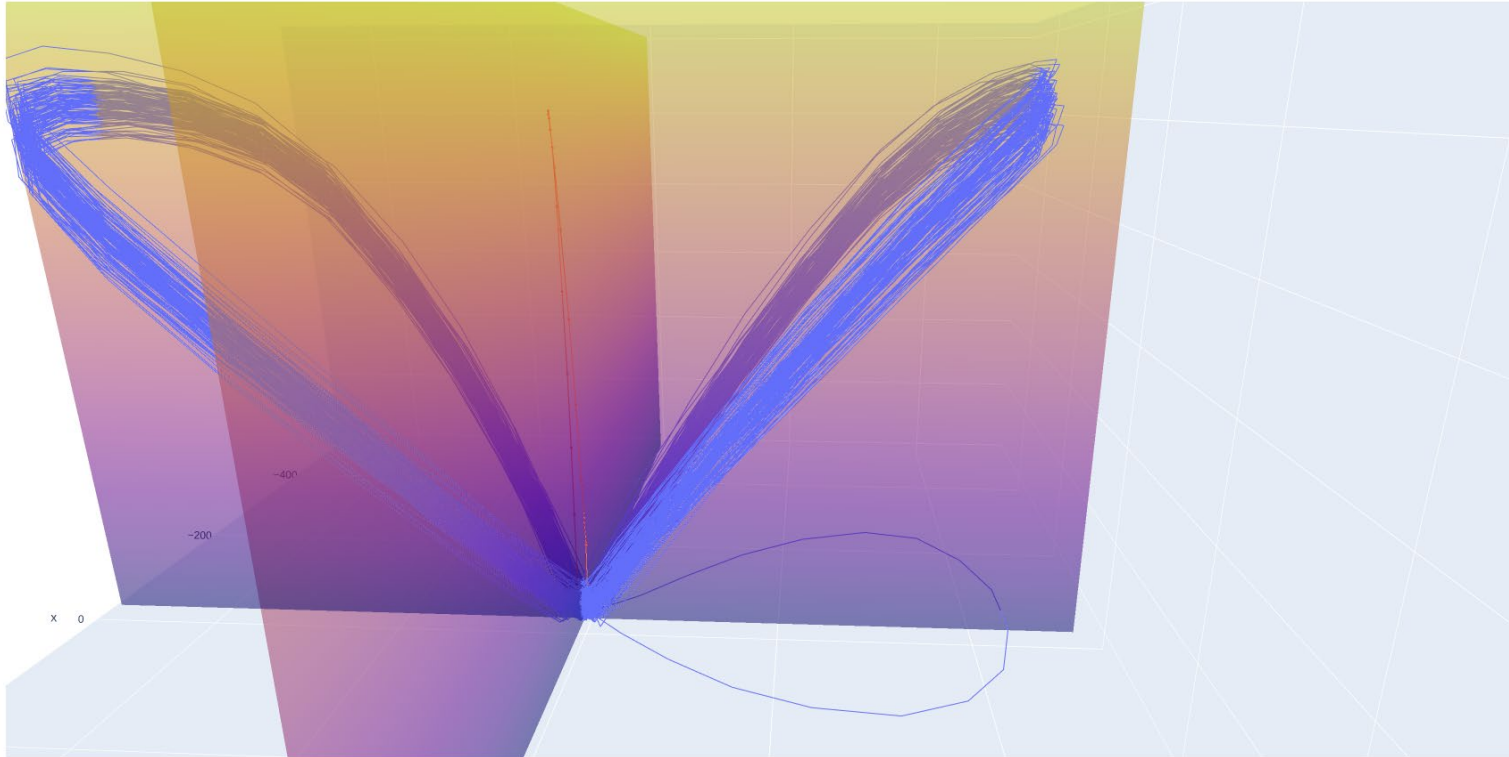
Series2Graph++



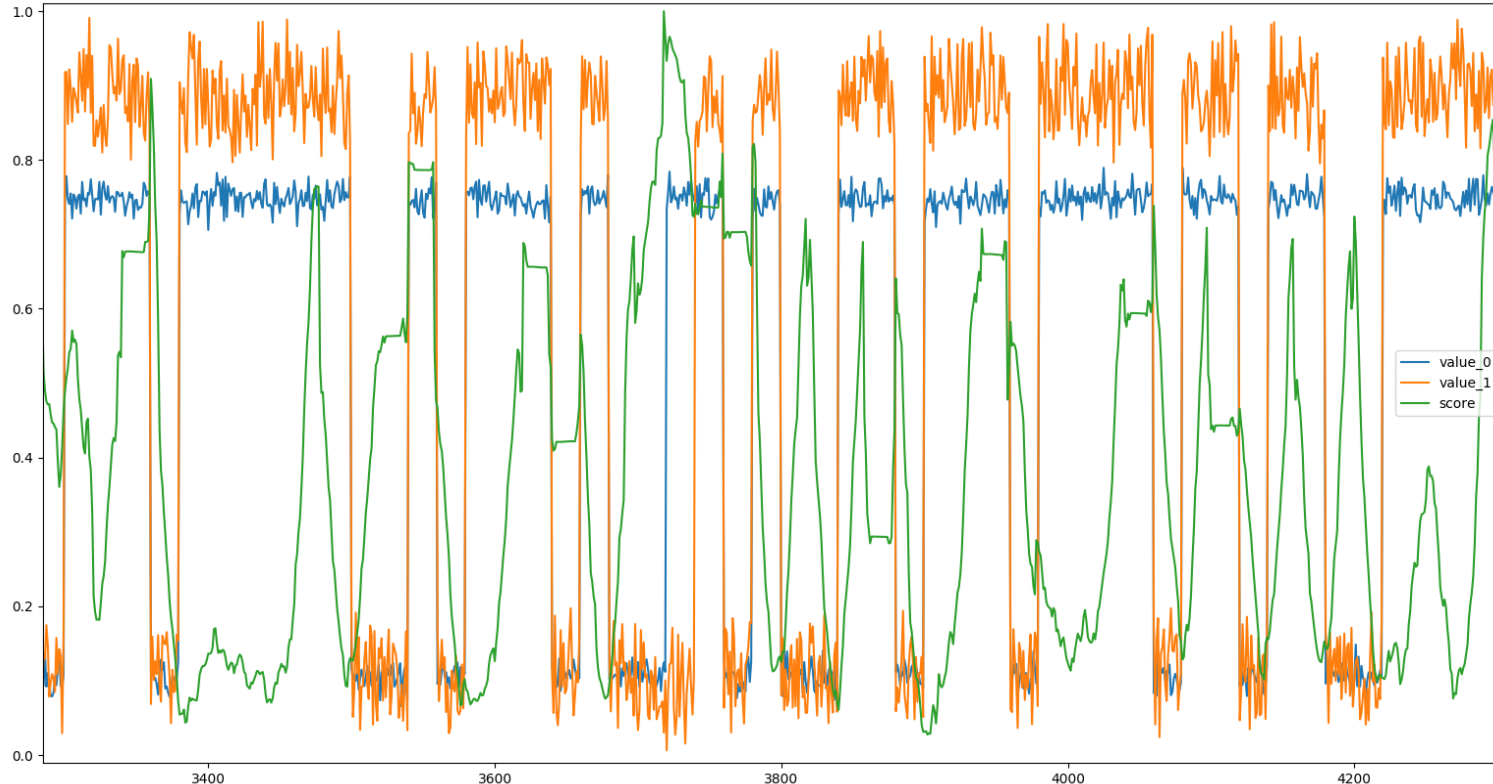
Series2Graph++



Series2Graph++



Series2Graph++



Homework

Homework

- Literature: Draft of our experimental evaluation paper

- Introduction to subsequence anomaly detection in time series
- Overview about current State-of-the-Art methods
- Categorization of methods
- Experiment setup and assumptions



Don't share or publish!

Find in the PDF in our seminar OwnCloud share.

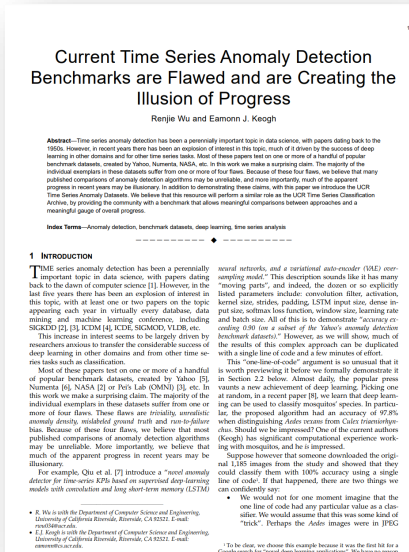
Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 54

Homework

- Literature: Ongoing discussion on how to assess time series anomaly detection methods, e.g.

- Current datasets

- are too **trivial**
- have unrealistically **dense** anomalies
- have **mislabeled** data
- have more anomalies in the **last** points



Source:

<https://doi.org/10.1109/TKDE.2021.3112126>

Schmidl & Wenig
Large-Scale TSA
Winter 2021/22
Chart 55

Homework

- Check out the material in our seminar share (E-Mail with link follows)
 - Datasets
 - Algorithms
 - Literature (Papers)
 - Metadata and evaluation results
 - ...

- Decide on your team direction (goal):
 - More reliable
 - More accurate
 - More efficient
 - More capable

Further Survey/Review papers

Chandola, Varun, Arindam Banerjee, and Vipin Kumar. **Anomaly Detection: A Survey.** ACM Computing Surveys 41, no. 3 (2009): 1–58.
<https://doi.org/10.1145/1541880.1541882>.

Gupta, Manish, Jing Gao, Charu C. Aggarwal, and Jiawei Han. **Outlier Detection for Temporal Data: A Survey.** IEEE Transactions on Knowledge and Data Engineering (TKDE) 26, no. 9 (2014): 2250–67. <https://doi.org/10.1109/TKDE.2013.184>.

Choudhary, Dhruv, Arun Kejariwal, and Francois Orsini. **On the Runtime-Efficacy Trade-off of Anomaly Detection Techniques for Real-Time Streaming Data.** ArXiv:1710.04735, October 12, 2017. <http://arxiv.org/abs/1710.04735>.

Chalapathy, Raghavendra, and Sanjay Chawla. **Deep Learning for Anomaly Detection: A Survey.** ArXiv:1901.03407, 2019. <http://arxiv.org/abs/1901.03407>.

Blázquez-García, Ane, Angel Conde, Usue Mori, and Jose A. Lozano. **A Review on Outlier/Anomaly Detection in Time Series Data.** ArXiv:2002.04236, February 11, 2020. <http://arxiv.org/abs/2002.04236>.

Braei, Mohammad, and Sebastian Wagner. **Anomaly Detection in Univariate Time-Series: A Survey on the State-of-the-Art.** ArXiv:2004.00433, April 1, 2020.
<http://arxiv.org/abs/2004.00433>.

