

## APPENDIX

The Appendix supplies additional materials regarding the evaluation of TSEvo and the implementation. Appendix A show the datasetwise results for the mutation ablation study, appendix B shows a sample visualization of the multivariate dataset and appendix C describe the implementation and code basis.

### A. Mutation Types

This section contains the remaining results. Table VII, VI, and VIII show the results for sparsity ( $R_2$ ), proximity ( $R_1$ ), and plausibility ( $R_3$ ) on dataset level. Figure 5 visualizes the counterfactuals achieved with the different mutation types on the first test image.

Table VII shows the sparsity per dataset and mutation type. Except for GunPoint, Authentic Opposing Information produces the counterfactuals with the lowest number of changed time steps. The generally much larger sparsity for frequency mutation results from the exchange of frequency bands, leading to a changed frequency and, therefore, more changed time steps. TSEvo achieves the best proximity scores on authentic opposing information for most datasets, followed by combination mutation. The best plausibility values for each dataset are distributed between the three basic mutation types Opposing, Frequency, and Gaussian. Note that on  $R_3$  frequency mutation performs exceptionally well on sensor-related datasets, indicating that authentic opposing information might not provide a good solution for all types of datasets. Combination Mutation is ranked second for most datasets and metrics, indicating that different mutations might be necessary for different dataset types to achieve plausible results.

	Opposing	Frequency	Gaussian	Combination
CBF	<b>0.22</b> $\pm$ 0.13	0.28 $\pm$ 0.14	0.59 $\pm$ 0.17	0.29 $\pm$ 0.14
CharacterTrajectories	<b>0.21</b> $\pm$ 0.07	1.0 $\pm$ 0.01	0.66 $\pm$ 0.14	0.26 $\pm$ 0.13
Coffee	<b>0.03</b> $\pm$ 0.02	0.04 $\pm$ 0.01	0.05 $\pm$ 0.02	<b>0.03</b> $\pm$ 0.02
ECG5000	0.34 $\pm$ 0.12	0.39 $\pm$ 0.16	0.45 $\pm$ 0.14	<b>0.33</b> $\pm$ 0.11
ElectricDevices	0.12 $\pm$ 0.09	<b>0.2</b> $\pm$ <b>0.17</b>	0.58 $\pm$ 0.16	<b>0.2</b> $\pm$ <b>0.14</b>
FordA	<b>0.12</b> $\pm$ <b>0.11</b>	0.15 $\pm$ 0.13	0.24 $\pm$ 0.27	0.14 $\pm$ 0.12
GunPoint	0.25 $\pm$ 0.3	0.22 $\pm$ 0.28	<b>0.16</b> $\pm$ <b>0.12</b>	0.25 $\pm$ 0.3
Heartbeat	<b>0.0</b> $\pm$ <b>0.01</b>	0.01 $\pm$ 0.03	0.02 $\pm$ 0.03	0.01 $\pm$ 0.03
NATOPS	<b>0.15</b> $\pm$ <b>0.09</b>	0.17 $\pm$ 0.07	0.19 $\pm$ 0.1	0.17 $\pm$ 0.1
UWaveGestureLibrary	<b>0.35</b> $\pm$ <b>0.15</b>	0.43 $\pm$ 0.19	0.83 $\pm$ 0.12	0.57 $\pm$ 0.24

TABLE VI: Proximity ( $R_1$ ) for each dataset and mutation type. The lower, the smaller the changes made to the original instance.

	Opposing	Frequency	Gaussian	Combination
CBF	<b>0.26</b> $\pm$ <b>0.15</b>	0.88 $\pm$ 0.14	0.62 $\pm$ 0.16	0.36 $\pm$ 0.15
CharacterTrajectories	<b>0.21</b> $\pm$ <b>0.07</b>	1.0 $\pm$ 0.01	0.66 $\pm$ 0.14	0.26 $\pm$ 0.13
Coffee	<b>0.12</b> $\pm$ <b>0.08</b>	0.78 $\pm$ 0.21	0.37 $\pm$ 0.1	0.15 $\pm$ 0.1
ECG5000	<b>0.38</b> $\pm$ <b>0.13</b>	0.98 $\pm$ 0.05	0.56 $\pm$ 0.14	0.45 $\pm$ 0.19
ElectricDevices	<b>0.32</b> $\pm$ <b>0.18</b>	0.8 $\pm$ 0.21	0.64 $\pm$ 0.17	0.4 $\pm$ 0.24
FordA	<b>0.21</b> $\pm$ <b>0.19</b>	0.82 $\pm$ 0.25	0.22 $\pm$ 0.24	0.31 $\pm$ 0.29
Heartbeat	<b>0.03</b> $\pm$ <b>0.03</b>	0.88 $\pm$ 0.18	0.06 $\pm$ 0.08	0.11 $\pm$ 0.21
GunPoint	0.35 $\pm$ 0.38	0.66 $\pm$ 0.39	<b>0.33</b> $\pm$ <b>0.2</b>	0.37 $\pm$ 0.38
NATOPS	<b>0.25</b> $\pm$ <b>0.15</b>	0.96 $\pm$ 0.04	0.32 $\pm$ 0.15	0.57 $\pm$ 0.33
UWaveGestureLibrary	<b>0.29</b> $\pm$ <b>0.12</b>	0.96 $\pm$ 0.1	0.72 $\pm$ 0.1	0.85 $\pm$ 0.24

TABLE VII: Sparsity ( $R_2$ ) for each dataset and mutation type. The lower, the smaller the number of changes made to the original instance.

	Opposing	Frequency	Gaussian	Combination
CBF	0.9688	<b>0.9719</b>	0.9688	0.9703
CharacterTrajectories	<b>0.9813</b>	0.9791	0.9791	0.9791
Coffee	<b>0.9902</b>	0.9881	0.9881	0.9874
ECG5000	0.9714	<b>0.9743</b>	0.9729	0.9729
ElectricDevices	0.9625	<b>0.9708</b>	0.9646	0.9667
FordA	0.9924	0.9932	<b>0.996</b>	0.9924
GunPoint	0.9813	0.984	0.976	0.9813
Heartbeat	0.9906	<b>0.9916</b>	0.9911	0.9906
NATOPS	<b>0.9333</b>	0.9294	0.9294	0.9294
UWaveGestureLibrary	0.9873	0.9873	<b>0.9905</b>	0.9873

TABLE VIII: Plausibility ( $R_3$ ) of the generated counterfactuals for each dataset and mutation type. The higher, the better - i.e., more data support from the training data.

### B. Benchmarking

Figure 6 shows the visualization of the counterfactual for the first time series in the test set of NATOPS, generated with TSEvo and by Ates et al. [1]. The counterfactuals (yellow) are plotted onto the original time series (blue). Each line in the plot denotes a feature time series. We did not include W-CF in the visualizations as only a few W-CF were valid.

### C. Implementation

The implementation of TSEvo is based on DEAP<sup>5</sup>, a highly flexible python framework for evolutionary computation. The classification models for the evaluation of the counterfactual methods are written in PyTorch<sup>6</sup>. Although not used in this paper, TSEvo also supports tensorflow. For usage with tensorflow please refer to our github repository.

For Ates et al. [1], we used their python implementation on github: <https://github.com/peacelab/CoMTE>. We only adapted the prediction function to enable the usage with pytorch. We reimplemented the approach of Wachter et al. [32] in PyTorch, as most library implementations use tensorflow as Basis (e.g. Alibi<sup>7</sup>) or only cope with univariate classification (e.g. CARLA [23]). NUN-CF and NUN-CF GradCam were adapted from the authors github implementation: [https://github.com/e-delaney/Instance-Based\\_CFE\\_TSC](https://github.com/e-delaney/Instance-Based_CFE_TSC) (commit: 19f87a1) to enable the usage in a multiclass classification on PyTorch.

TSEvo as well as the implemented or adapted benchmarks can be found in our github repository: <https://anonymous.4open.science/r/TSEvo-869B/Readme.md>.

<sup>5</sup><https://deap.readthedocs.io/en/master/>

<sup>6</sup><https://pytorch.org/>

<sup>7</sup><https://docs.seldon.io/projects/alibi/en/latest/>

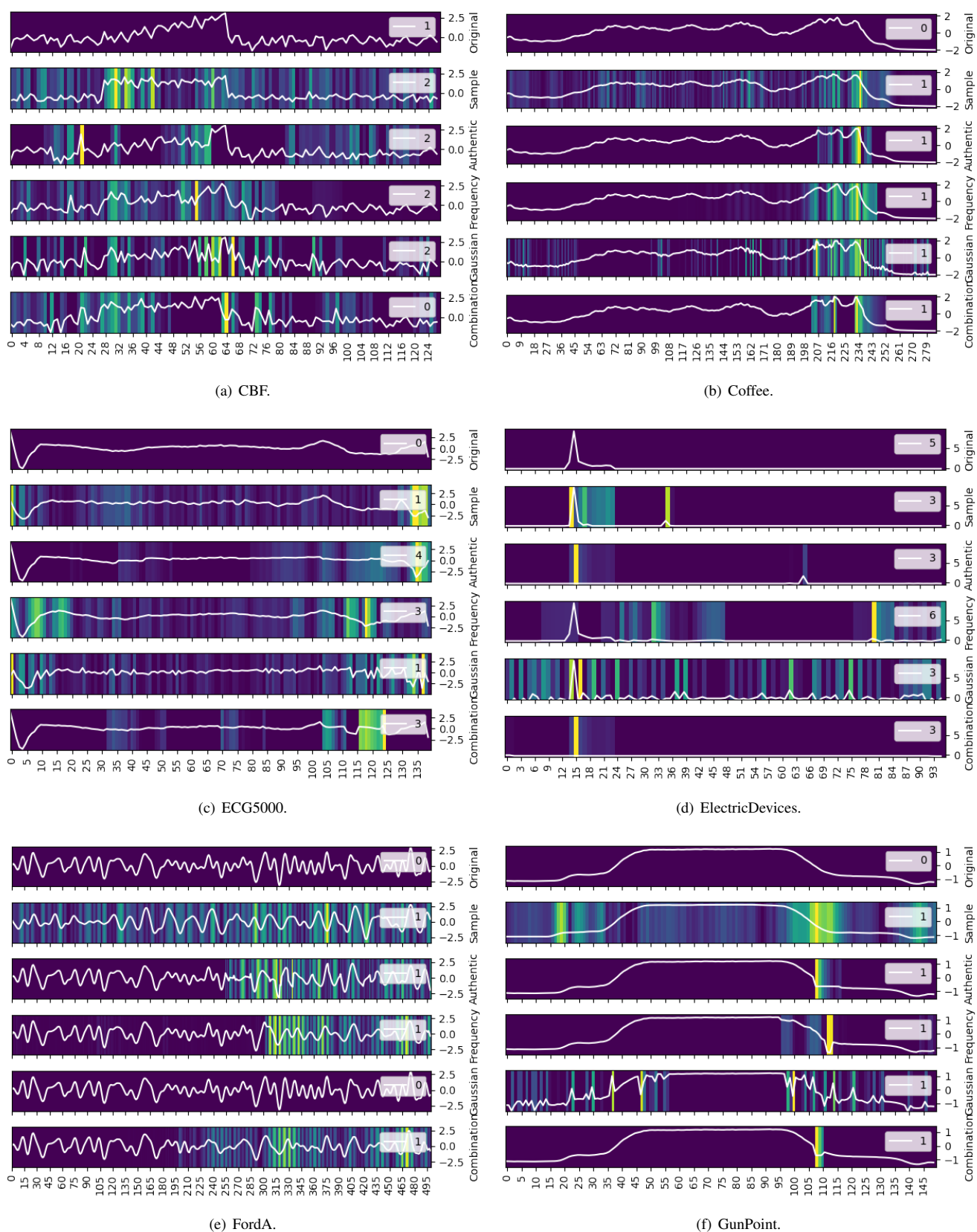
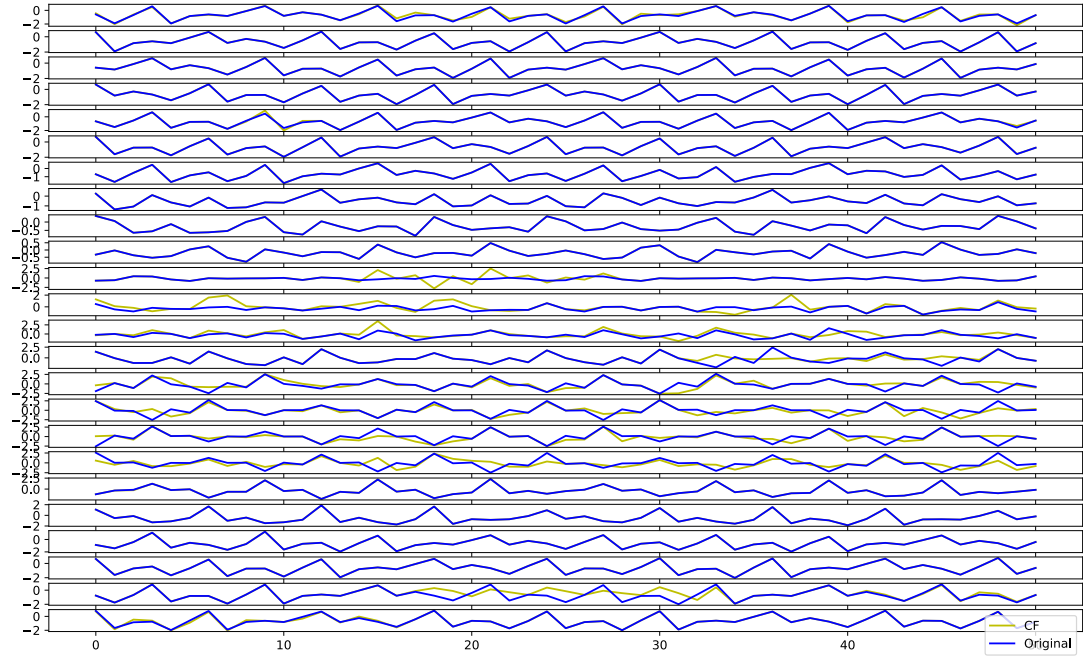
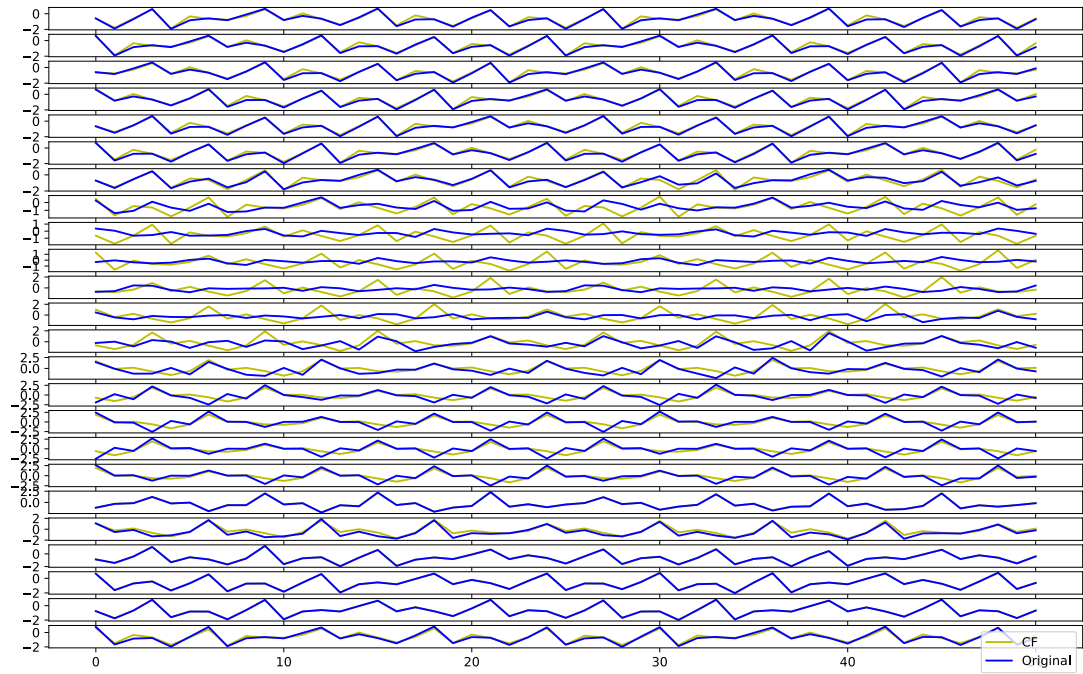


Fig. 5: Counterfactual for the first timeseries of the test set from CBR, Coffe, ECG5000, Electric Devices, GunPoint and FordA obtained with the different mutation types. If the labels are consistent with the original classification, the method failed to generate a true counterfactual.



(a) TSEvo.



(b) Ates.

Fig. 6: Visualization of the original and counterfactual of the first test time series on the multivariate dataset NATOPS. Blue is the original time series feature and yellow is the CF time series feature. If no yellow is used, the time-series feature stays unchanged.