
Supplementary Material: Nominality Score Conditioned Time Series Anomaly Detection by Point/Sequential Reconstruction

Anonymous Author(s)

Affiliation

Address

email

1 A The Best F1 Score

2 This section describes the method for calculating the best F1 score ($F1^*$) from a set of anomaly scores
3 $\mathbf{a} = \{a_1, \dots, a_T\}$ and a set of labels $\mathbf{y} = \{y_1, \dots, y_T\}$. Firstly, given \mathbf{a} and some arbitrary threshold
4 θ_a , we can calculate $\hat{\mathbf{y}} = \{\hat{y}_1, \dots, \hat{y}_T\}$, where $\hat{y}_t \triangleq \mathbb{1}_{a_t \geq \theta_a}$. Secondly, $\hat{\mathbf{y}}$ is used to calculate TP, FP,
5 and FN, which corresponds to the sets of time points for true positives, false positives, and false
6 negatives.

$$\text{TP} \triangleq \{t | \hat{y}_t = 1, y_t = 1\}, \quad \text{FP} \triangleq \{t | \hat{y}_t = 1, y_t = 0\}, \quad \text{FN} \triangleq \{t | \hat{y}_t = 0, y_t = 1\} \quad (1)$$

7 Thirdly, we calculate the precision (P) and recall (R), and then calculate the F1 score, which is the
8 harmonic mean between R and P.

$$\text{F1} \triangleq \frac{2\text{PR}}{\text{P} + \text{R}}, \quad \text{P} \triangleq \frac{\text{n}(\text{TP})}{\text{n}(\text{TP}) + \text{n}(\text{FP})}, \quad \text{R} \triangleq \frac{\text{n}(\text{TP})}{\text{n}(\text{TP}) + \text{n}(\text{FN})} \quad (2)$$

9 Finally, we calculate $F1^*$ by using the threshold that yields the highest F1.

$$F1^*(\mathbf{a}; \mathbf{y}) \triangleq \max_{\theta_a} F1(\hat{\mathbf{y}}(\mathbf{a}, \theta_a); \mathbf{y}) \quad (3)$$

10 B Model Architecture

11 Performers (an improved variant of Transformers) are competitive in terms of execution speed
12 compared with other Transformer variants [1, 2], hence we use them as the basic building block of
13 our models. For both \mathcal{M}_{pt} and \mathcal{M}_{seq} , we use a linear layer with input and output dimensions equal
14 to D as the token embedding layer, a fixed positional embedding layer at the beginning, a feature
15 redraw interval of 1, and a tanh activation function immediately before the output. GELUs [3] are
16 used as the activation layer for all linear layers. We do not change any other predefined activation
17 layer inside Performers.

18 B.1 Performer-based autoencoder

19 For the Performer-based autoencoder \mathcal{M}_{pt} , the input with the shape $(batch, W, D)$ is passed through
20 a plain Performer with N_{perf} layers after the positional embedding step, where W represents the
21 input window size. This is followed by a linear encoding layer that transforms the dimensionality
22 from D to D_{lat} , resulting in the shape $(batch, W, D_{lat})$ for the latent variables. In the case of \mathcal{M}_{pt} ,
23 there is no compression along the time domain. The latent variables then pass through another
24 linear decoding layer that transforms the dimensionality from D_{lat} back to D , followed by another
25 Performer with N_{perf} layers.

26 B.2 Performer-based stacked encoder

27 For the Performer-based stacked encoder \mathcal{M}_{seq} , the input with shape $(batch, W_0, D)$ is passed
 28 through a plain Performer with one layer after the positional embedding step, followed by a linear
 29 encoding layer that transforms the window size from W_0 to W_1 , where $W_0 = 2\gamma$ (cf. section 2.6).
 30 It should be noted that unlike \mathcal{M}_{pt} , compression is done along the time domain for \mathcal{M}_{seq} . The
 31 one-layer Performer and linear layer are stacked N_{enc} times, where in the i -th stack, the window
 32 size is compressed from W_{i-1} to W_i . $W_{N_{enc}}$ is equal to the target output window size δ . For both
 33 \mathcal{M}_{pt} and \mathcal{M}_{seq} , we optimize N_{perf} , D_{lat} , W , N_{enc} , W_i for $i \in \{0, \dots, N_{enc}\}$, and δ for the best
 34 performance. Note that \mathcal{M}_{seq} isn't capable of reconstructing the first and last γ time points due to its
 35 architecture, hence we discard the first and last γ points reconstructed by \mathcal{M}_{pt} so that rest of the time
 36 points have exactly two reconstructed values corresponding to using \mathcal{M}_{seq} and \mathcal{M}_{pt} , respectively.

37 C Data Preprocessing and Training Details

38 Table 1 shows the hyperparameters used for implementing NPSR on the experimented datasets. To
 39 ensure fair comparison, the same preprocessing method is applied to all algorithms for the same
 40 dataset. The search for hyperparameters is done manually, starting from some reasonable value (e.g.
 41 a learning rate of 10^{-4}). The authors believe that there is still room for improvement by fine-tuning
 42 these hyperparameters. To speed up training, we load all training inputs and outputs, and testing
 43 inputs onto the GPU before training. We use a local GPU, which can be either GeForce RTX 3070
 44 (8GB), 3080 (12GB) or 3090 (24GB). For an individual experiment using a single dataset and training
 45 method, the training time ranges from approximately 2 minutes to 12 hours. For single-entity datasets
 46 and multi-entity datasets that use the combined training method, we run the experiments for at least 3
 47 times and confirm that the results are stable given different random seeds. For multi-entity datasets
 48 with entities trained individually, the results are averaged across all entities. Generally, datasets
 49 with single entities train faster than those with multiple entities. There are some additional remarks
 50 regarding the preprocessing of the datasets. For SWaT, we use `SWaT_Dataset_Attack_v0.csv` and
 51 `SWaT_Dataset_Normal_v1.csv` from the folder `SWaT.A1 & A2_Dec 2015` (manually converted
 52 from `*.xlsx`). We corrected some original flaws in the dataset (e.g. redundant blank spaces in some
 53 labels), and set the 5th and 10th columns to all 0. For WADI, we use the 2017 year dataset. Columns
 54 with excessive NaNs (more than half of the entire length) are deleted. Other NaNs are forward-filled.
 55 After deleting all the columns with excessive NaNs, the 86th column is further set to all 0. For PSM,
 56 we forward-fill all NaNs. For MSL and SMD, two additional blank channels are added to make the
 57 number of channels divisible by the number of heads.

Table 1: Implementation details. (c) stands for the combined training method (cf. section 3.2).

Parameter \ Dataset	SWaT	WADI	PSM	MSL	MSL (c)	SMAP	SMAP (c)	SMD	SMD (c)
	Preprocess								
Downsample	10	10	10	1	1	1	1	2	2
Normalization	Minmax	Minmax	Minmax	Minmax	Minmax	Minmax	Minmax	Minmax	Minmax
Stride	10	10	10	10	10	10	10	10	10
W (for \mathcal{M}_{pt})	100	100	100	100	100	50	50	50	50
W_0 (for \mathcal{M}_{seq})	100	100	100	50	50	50	50	50	50
δ	20	20	20	6	6	6	6	6	6
clamping (test data)	[-4,4]	[-4,4]	[-4,4]	[-4,4]	[-4,4]	[-4,4]	[-4,4]	[-4,4]	[-4,4]
	Model architecture								
# of heads	9	14	5	11	12	5	10	8	11
D_{lat}	10	10	10	10	10	10	10	10	10
ff_mult	4	4	4	4	4	4	4	4	4
N_{perf}	4	4	4	4	4	4	4	4	4
N_{enc}	8	8	8	8	8	8	8	8	8
	Induced anomaly score								
Gate function	soft	soft	soft	soft	hard	soft	hard	soft	soft
d	16	16	64	128	128	64	64	16	1
Ratio of N_{trn} for θ_N	99.85%	99.85%	99.85%	99.85%	97.5%	99.85%	∞	99.85%	99.85%
	Training								
Learn rate	10^{-4}	10^{-4}	10^{-4}	10^{-4}	10^{-4}	10^{-4}	10^{-4}	10^{-4}	10^{-4}
Optimizer	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam	Adam
Batch size	64	64	64	64	64	64	64	64	64
Training epochs	100	100	80	40	20	50	20	40	10

58 D The Point-adjusted Best F1 Score

59 Analogous to $F1^*$, the point-adjusted best F1 score ($F1_{PA}^*$) corresponds to $F1^*$ calculated after
 60 point-adjustment. Table 2 shows the $F1^*$ and $F1_{PA}^*$ of different algorithms, including NPSR, applied
 61 to several datasets.

62 The results suggest that $F1_{PA}^*$ may not be reliable - on the SWaT, WADI, PSM, and MSL datasets,
 63 simple heuristic approaches (e.g. using the mean squared value of an input time point as the anomaly
 64 score) outperform all deep learning methods when evaluated using $F1_{PA}^*$. Moreover, optimizing
 65 on $F1_{PA}^*$ does not necessarily guarantee a higher $F1^*$. NPSR is optimized on $F1^*$ by tuning the
 66 algorithm-specific parameters (e.g. d) and general parameters. To additionally optimize on $F1_{PA}^*$, we
 67 simply added *spikes* to the induced anomaly score ($\hat{A}_{spike}(\cdot)$) with value ∞ for some fixed interval s .
 68 The results show that $\hat{A}_{spike}(\cdot)$ can also achieve competitive $F1_{PA}^*$ values.

Table 2: Point-adjusted best F1 score ($F1_{PA}^*$) and best F1 score ($F1^*$) results on several datasets, with bold text denoting the highest and underlined text denoting the second highest value. The deep learning methods are sorted with older methods at the top and newer ones at the bottom.

Dataset	SWaT		WADI		PSM		MSL		SMAP		SMD	
Metric	$F1_{PA}^*$	$F1^*$										
Simple Heuristics [4, 5, 6]	0.969	0.789	0.965	0.353	0.985	0.509	0.965	0.239	0.961	0.229	0.934	0.494
DAGMM [7]	0.853	0.750	0.209	0.121	0.761	0.483	0.701	0.199	0.712	0.333	0.723	0.238
LSTM-VAE [8]	0.805	0.776	0.380	0.227	0.809	0.455	0.854	0.212	0.756	0.235	0.808	0.435
MSCRED [9]	0.807	0.757	0.374	0.046	0.626	0.556	0.936	0.250	0.866	0.170	0.841	0.382
OmniAnomaly [10]	0.866	0.782	0.417	0.223	0.664	0.452	0.901	0.207	0.854	0.227	<u>0.962</u>	0.474
MAD-GAN [11]	0.815	0.770	0.556	0.370	0.658	0.471	0.917	0.267	0.865	0.175	0.915	0.220
MTAD-GAT [12]	0.860	0.784	0.602	0.437	0.780	0.571	0.908	0.275	0.901	0.296	0.908	0.400
USAD [13]	0.846	0.792	0.430	0.233	0.725	0.479	0.911	0.211	0.819	0.228	0.946	0.426
THOC [14]	0.881	0.612	0.506	0.130	0.895	-	0.937	0.190	0.952	0.240	0.541	0.168
UAE [4]	0.869	0.453	<u>0.957</u>	0.354	0.936	0.427	0.920	<u>0.451</u>	0.896	0.390	0.972	0.435
GDN [15]	0.935	<u>0.810</u>	0.855	<u>0.570</u>	0.923	0.552	0.903	0.217	0.708	0.252	0.716	<u>0.529</u>
GTA [16]	0.910	0.761	0.84	0.531	0.855	0.542	0.911	0.218	0.904	0.231	0.919	0.351
Anomaly Transformer [17]	0.941	0.019	0.714	0.015	<u>0.979</u>	0.022	0.936	0.021	<u>0.967</u>	0.019	0.923	0.021
TranAD [18]	0.815	0.669	0.495	0.415	0.882	0.649	0.949	0.251	0.892	0.247	0.961	0.310
NPSR (combined)	-	-	-	-	-	-	<u>0.960</u>	0.261	0.978	0.511	0.850	0.252
NPSR	<u>0.953</u>	0.839	0.938	0.642	0.957	<u>0.648</u>	-	0.551	-	<u>0.437</u>	-	0.535

69 E Source of Data

70 Table 3 shows the data sources used to produce Table 2, as well as the sources for section 3.3. The
 71 reference number is followed by a number between 1 and 3, where 1 indicates that the data comes
 72 from the original work, 2 indicates that the data comes from reproduced values from another literature,
 73 and 3 indicates that we have reproduced the values using public repositories..

74 F Broader Impacts

75 The detection of anomalies in time series data can minimize downtime and avert financial losses.
 76 Utilizing real-time monitoring of system conditions, anomaly detection techniques for time series data
 77 can automatically detect deviations from the expected system behavior, thereby avoiding potential
 78 risks and financial harm. This has the potential to reduce the need for manual monitoring of faults
 79 and to expedite decision-making processes. Additionally, it can promote the sustainability of AI by
 80 preventing energy wastage and system malfunction.

Table 3: Data sources for algorithms and datasets. **Reproduced by using the squared value of channel 1 as the anomaly score.

Dataset	SWaT		WADI		PSM		MSL		SMAP		SMD	
	F1 _{PA} *	F1*										
Simple Heuristic	[5] - 1	[5] - 1	[5] - 1	[5] - 1	[6] - 1	**	[6] - 1	[5] - 1	[6] - 1	[5] - 1	[6] - 1	[5] - 1
DAGMM	[5] - 2	[18] - 3	[5] - 2	[5] - 2	[18] - 3	[18] - 3	[5] - 2	[5] - 2	[5] - 2	[5] - 2	[5] - 2	[5] - 2
LSTM-VAE	[13] - 2	[13] - 2	[13] - 2	[13] - 2	[17] - 2	[19] - 3	[13] - 2	[5] - 2	[13] - 2	[5] - 2	[13] - 2	[5] - 2
MSCRED	[18] - 2	[18] - 3	[18] - 2	[18] - 3	[18] - 3	[18] - 3	[18] - 2	[18] - 3	[18] - 2	[18] - 3	[18] - 2	[18] - 3
OmniAnomaly	[5] - 2	[5] - 2	[5] - 2	[5] - 2	[18] - 3	[18] - 3	[10] - 1	[5] - 2	[10] - 1	[5] - 2	[10] - 1	[5] - 2
MAD-GAN	[18] - 3	[11] - 1	[18] - 3	[11] - 1	[18] - 3	[18] - 3	[18] - 2	[18] - 3	[18] - 2	[18] - 3	[18] - 2	[18] - 3
MTAD-GAT	[20] - 2	[12] - 3	[20] - 2	[12] - 3	[12] - 3	[12] - 3	[12] - 1	[12] - 3	[12] - 1	[12] - 3	[20] - 2	[12] - 3
USAD	[13] - 1	[13] - 1	[13] - 1	[13] - 1	[18] - 3	[18] - 3	[13] - 1	[5] - 2	[13] - 1	[5] - 2	[13] - 1	[5] - 2
THOC	[14] - 1	[5] - 2	[5] - 2	[5] - 2	[17] - 2	-	[14] - 1	[5] - 2	[14] - 1	[5] - 2	[5] - 2	[5] - 2
UAE	[4] - 1	[4] - 1	[4] - 1	[4] - 1	[4] - 3	[4] - 3	[4] - 1	[4] - 1	[4] - 1	[4] - 1	[4] - 1	[4] - 1
GDN	[5] - 2	[15] - 1	[5] - 2	[15] - 1	[15] - 3	[15] - 3	[5] - 2	[5] - 2	[5] - 2	[5] - 2	[5] - 2	[5] - 2
GTA	[16] - 1	[16] - 3	[16] - 1	[16] - 3	[16] - 3	[16] - 3	[16] - 1	[16] - 3	[16] - 1	[16] - 3	[16] - 3	[16] - 3
AnomalyTransformer	[17] - 1	[17] - 3	[17] - 3	[17] - 3	[17] - 1	[17] - 3	[17] - 1	[17] - 3	[17] - 1	[17] - 3	[17] - 1	[17] - 3
TranAD	[18] - 1	[18] - 3	[18] - 1	[18] - 3	[18] - 3	[18] - 3	[18] - 1	[18] - 3	[18] - 1	[18] - 3	[18] - 1	[18] - 3

81 G Limitations

82 Despite exhibiting competitive performance against other models, the proposed NPSR algorithm has
 83 a few limitations. The point-based model used in training does not incorporate temporal informa-
 84 tion, which makes it challenging to effectively reconstruct low-dimensional datasets. This issue is
 85 particularly challenging for univariate time series since raw inputs would not work for point-based
 86 models. One possible solution to this problem is to increase dimensionality by aggregating multiple
 87 time points. However, the effectiveness of this approach is yet to be confirmed.

88 Another limitation of NPSR is the absence of an automatic threshold-finding method, which makes it
 89 difficult to determine a suitable threshold when deploying the model. To address this issue, one can
 90 define a target false positive rate and estimate the threshold that achieves this target rate using the
 91 validation set since only normal data is needed.

92 References

- 93 [1] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos,
 94 Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers.
 95 *arXiv preprint arXiv:2009.14794*, 2020.
- 96 [2] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang,
 97 Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient transformers. *arXiv*
 98 *preprint arXiv:2011.04006*, 2020.
- 99 [3] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*,
 100 2016.
- 101 [4] Astha Garg, Wenyu Zhang, Jules Samaran, Ramasamy Savitha, and Chuan-Sheng Foo. An evaluation of
 102 anomaly detection and diagnosis in multivariate time series. *IEEE Transactions on Neural Networks and*
 103 *Learning Systems*, 33(6):2508–2517, jun 2022.
- 104 [5] Siwon Kim, Kukjin Choi, Hyun-Soo Choi, Byunghan Lee, and Sungroh Yoon. Towards a rigorous
 105 evaluation of time-series anomaly detection, 2022.
- 106 [6] Keval Doshi, Shatha Abudalou, and Yasin Yilmaz. Reward once, penalize once: Rectifying time series
 107 anomaly detection. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2022.
- 108 [7] Bo Zong, Qi Song, Martin Renqiang Min, Wei Cheng, Cristian Lumezanu, Daeki Cho, and Haifeng
 109 Chen. Deep autoencoding gaussian mixture model for unsupervised anomaly detection. In *International*
 110 *Conference on Learning Representations*, 2018.
- 111 [8] Daehyung Park, Yuuna Hoshi, and Charles C Kemp. A multimodal anomaly detector for robot-assisted
 112 feeding using an lstm-based variational autoencoder. *IEEE Robotics and Automation Letters*, 3(3):1544–
 113 1551, 2018.
- 114 [9] Chuxu Zhang, Dongjin Song, Yuncong Chen, Xinyang Feng, Cristian Lumezanu, Wei Cheng, Jingchao
 115 Ni, Bo Zong, Haifeng Chen, and Nitesh V Chawla. A deep neural network for unsupervised anomaly
 116 detection and diagnosis in multivariate time series data. In *Proceedings of the AAAI conference on artificial*
 117 *intelligence*, volume 33, pages 1409–1416, 2019.

- 118 [10] Ya Su, Youjian Zhao, Chenhao Niu, Rong Liu, Wei Sun, and Dan Pei. Robust anomaly detection for
119 multivariate time series through stochastic recurrent neural network. In *Proceedings of the 25th ACM*
120 *SIGKDD International Conference on Knowledge Discovery; Data Mining*, KDD '19, page 2828–2837,
121 New York, NY, USA, 2019. Association for Computing Machinery.
- 122 [11] Dan Li, Dacheng Chen, Baihong Jin, Lei Shi, Jonathan Goh, and See-Kiong Ng. Mad-gan: Multivariate
123 anomaly detection for time series data with generative adversarial networks. In *Artificial Neural Networks*
124 *and Machine Learning–ICANN 2019: Text and Time Series: 28th International Conference on Artificial*
125 *Neural Networks, Munich, Germany, September 17–19, 2019, Proceedings, Part IV*, pages 703–716.
126 Springer, 2019.
- 127 [12] Hang Zhao, Yujing Wang, Juanyong Duan, Congrui Huang, Defu Cao, Yunhai Tong, Bixiong Xu, Jing Bai,
128 Jie Tong, and Qi Zhang. Multivariate time-series anomaly detection via graph attention network. In *2020*
129 *IEEE International Conference on Data Mining (ICDM)*, pages 841–850. IEEE, 2020.
- 130 [13] Julien Audibert, Pietro Michiardi, Frédéric Guyard, Sébastien Marti, and Maria A. Zuluaga. Usad:
131 Unsupervised anomaly detection on multivariate time series. In *Proceedings of the 26th ACM SIGKDD*
132 *International Conference on Knowledge Discovery; Data Mining*, KDD '20, page 3395–3404, New York,
133 NY, USA, 2020. Association for Computing Machinery.
- 134 [14] Lifeng Shen, Zhuocong Li, and James Kwok. Timeseries anomaly detection using temporal hierarchical
135 one-class network. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances*
136 *in Neural Information Processing Systems*, volume 33, pages 13016–13026. Curran Associates, Inc., 2020.
- 137 [15] Ailin Deng and Bryan Hooi. Graph neural network-based anomaly detection in multivariate time series.
138 *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(5):4027–4035, May 2021.
- 139 [16] Zekai Chen, Dingshuo Chen, Xiao Zhang, Zixuan Yuan, and Xiuzhen Cheng. Learning graph structures
140 with transformer for multivariate time-series anomaly detection in iot. *IEEE Internet of Things Journal*,
141 9(12):9179–9189, 2021.
- 142 [17] Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer: Time series anomaly
143 detection with association discrepancy. *arXiv preprint arXiv:2110.02642*, 2021.
- 144 [18] Shreshth Tuli, Giuliano Casale, and Nicholas R Jennings. Tranad: Deep transformer networks for anomaly
145 detection in multivariate time series data. *arXiv preprint arXiv:2201.07284*, 2022.
- 146 [19] TimyadNyda. Github - timyadnyda/variational-lstm-autoencoder: Lstm variational auto-encoder for time
147 series anomaly detection and features extraction.
- 148 [20] Jun Zhan, Siqi Wang, Xiandong Ma, Chengkun Wu, Canqun Yang, Detian Zeng, and Shilin Wang. Stgat-
149 mad : Spatial-temporal graph attention network for multivariate time series anomaly detection. In *ICASSP*
150 *2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages
151 3568–3572, 2022.