

STAN: Spectral-Temporal Attention Networks for Multivariate Time-Series Anomaly Detection

Anonymous Author(s)
Affiliation
anonymous@email.com

Abstract

We introduce STAN (Spectral-Temporal Attention Networks), a dual-branch architecture that detects anomalies in multivariate time series by jointly modeling frequency-domain and temporal-domain representations. Existing methods process time series in a single domain, missing anomalies that manifest differently in spectral versus temporal views—periodic pattern disruptions are salient in the frequency domain while point anomalies are easier to capture temporally. STAN addresses this via (i) a spectral branch that applies attention over the top- k FFT components, (ii) a causal temporal transformer branch, and (iii) a cross-attention gate that learns an adaptive, per-timestep fusion of both representations, trained with a joint reconstruction and spectral contrastive objective. An Extreme Value Theory calibrator provides statistically grounded per-channel thresholds. STAN achieves new state-of-the-art results on all five standard benchmarks (SMD, MSL, SMAP, SWaT, PSM), with an average point-adjust F1 of 92.54%—a +2.71 point improvement over the best prior method on each dataset—while using 32% fewer parameters and 43% less inference time than the closest transformer baseline.

1 Introduction

Multivariate time-series anomaly detection is critical for monitoring server fleets [9], spacecraft telemetry [4], and industrial control systems [6]. The task is challenging because anomalies are rare, diverse in morphology, and often involve subtle inter-channel dependencies.

Recent transformer-based approaches [12–14] have advanced state of the art by modeling long-range temporal dependencies. However, they operate primarily in the time domain. This is a fundamental limitation: many real-world anomalies—such as disrupted periodicity in server CPU cycles or shifted resonance frequencies in mechanical systems—are most visible in the frequency domain, while contextual point anomalies are best captured temporally. No single-domain representation suffices.

FEDformer [15] demonstrated the value of frequency-enhanced features for time-series forecasting, but its frequency components serve as auxiliary features rather than a first-class processing path. We argue that anomaly detection requires *equal-footing* spectral and temporal processing with a learned fusion mechanism that can adapt to each timestep.

37 We propose STAN, a dual-branch architecture with three key design choices:

- 38 • **Spectral branch.** An FFT-based attention module that selects the top- k frequency
39 components by magnitude and applies multi-head attention over them, capturing
40 periodic structure and spectral deviations.
- 41 • **Temporal branch.** A causal transformer encoder that models sequential dependencies
42 with autoregressive masking, ensuring anomaly scores reflect only past and present
43 context.
- 44 • **Cross-attention gate.** A bidirectional cross-attention mechanism with a learned
45 sigmoid gate that produces a per-timestep, per-dimension adaptive fusion of spectral
46 and temporal representations.

47 Training combines a reconstruction objective with a spectral contrastive loss (NT-Xent [3])
48 that encourages discriminative frequency representations. At inference, an Extreme Value
49 Theory (EVT) calibrator [8] fits a Generalized Extreme Value distribution to tail scores,
50 providing statistically principled per-channel thresholds.

51 STAN achieves new state of the art on all five standard multivariate anomaly detection
52 benchmarks, with an average F1 of 92.54% under the point-adjust protocol [13], while being
53 the most parameter-efficient transformer-based method at 2.8M parameters.

54 2 Related Work

55 **Deep learning for time-series anomaly detection.** Early deep approaches used recurrent
56 architectures: OmniAnomaly [9] combines stochastic recurrent networks with normalizing
57 flows, and LSTM-VAE [7] uses variational autoencoders with LSTM backbones. USAD [2]
58 introduced adversarial training for unsupervised anomaly detection on multivariate series.
59 More recently, transformer-based methods have dominated: Anomaly Transformer [13]
60 introduces association discrepancy to distinguish anomalies from normal patterns, DCdeter-
61 tor [14] uses dual attention with contrastive learning, and TimesNet [12] models temporal
62 2D variations. TranAD [10] applies deep transformers with adversarial training. All these
63 methods operate primarily in the time domain.

64 **Frequency-domain methods.** Autoformer [11] introduced decomposition with auto-
65 correlation for time-series forecasting, and FEDformer [15] extended this with frequency-
66 enhanced decomposition, demonstrating that spectral features improve temporal modeling.
67 However, in both cases frequency components are auxiliary to a primarily temporal archi-
68 tecture. In contrast, STAN treats spectral processing as a first-class branch with dedicated
69 attention and learned fusion.

70 **Contrastive learning for anomaly detection.** DCdetector [14] applies contrastive learning
71 to temporal representations. SimCLR [3] established the NT-Xent framework for visual
72 contrastive learning. STAN adapts contrastive learning specifically to spectral embeddings,
73 using frequency-domain augmentations (frequency masking) alongside temporal jitter to
74 create complementary views.

75 **Threshold calibration.** Most anomaly detectors use fixed percentile thresholds, which
 76 are sensitive to score distribution shape. EVT-based approaches [8] provide statistically
 77 grounded thresholds by modeling the tail of the score distribution. STAN uses a per-channel
 78 GEV fit that adapts to heterogeneous score distributions across channels.

79 3 Method

80 Given a multivariate time series $\mathbf{X} \in \mathbb{R}^{T \times C}$ with T timesteps and C channels, STAN
 81 processes sliding windows $\mathbf{x} \in \mathbb{R}^{W \times C}$ (window size W) through dual branches, fuses them
 82 via cross-attention gating, and produces per-timestep anomaly scores $\mathbf{s} \in \mathbb{R}^W$.

83 3.1 Spectral Branch

84 The spectral branch converts the input to the frequency domain via the real FFT, selects
 85 informative frequency components, and applies attention over them.

86 **Frequency selection.** Given a projected input $\mathbf{H} = \mathbf{x}\mathbf{W}_{\text{in}} \in \mathbb{R}^{W \times D}$ where $\mathbf{W}_{\text{in}} \in \mathbb{R}^{C \times D}$,
 87 we compute the real FFT along the time axis:

$$\hat{\mathbf{H}} = \text{FFT}(\mathbf{H}) \in \mathbb{C}^{F \times D}, \quad F = \lfloor W/2 \rfloor + 1 \quad (1)$$

88 We retain the top- k frequency components by mean magnitude across dimensions:

$$\mathcal{I}_k = \text{top-}k \left(\frac{1}{D} \sum_{d=1}^D |\hat{\mathbf{H}}_{:,d}| \right), \quad \hat{\mathbf{H}}_k = \hat{\mathbf{H}}[\mathcal{I}_k] \in \mathbb{R}^{k \times D} \quad (2)$$

89 where we use magnitudes $|\hat{\mathbf{H}}_k|$ as the real-valued representation.

90 **Spectral attention.** Multi-head self-attention is applied over the k selected frequency
 91 components:

$$\mathbf{S} = \text{LayerNorm} \left(\text{MHA}(\hat{\mathbf{H}}_k, \hat{\mathbf{H}}_k, \hat{\mathbf{H}}_k) + \hat{\mathbf{H}}_k \right) \in \mathbb{R}^{k \times D} \quad (3)$$

92 The spectral embedding is projected back to the temporal dimension via linear interpolation:
 93 $\mathbf{E}_s = \text{Interp}(\mathbf{S}, W) \in \mathbb{R}^{W \times D}$.

94 3.2 Temporal Branch

95 The temporal branch is a standard transformer encoder with causal masking. The input is
 96 projected and augmented with learnable positional encodings:

$$\mathbf{H}_t = \mathbf{x}\mathbf{W}_{\text{in}}^{(t)} + \mathbf{P}_{:,W}, \quad \mathbf{P} \in \mathbb{R}^{W_{\text{max}} \times D} \quad (4)$$

97 A causal mask $\mathbf{M} \in \{0, -\infty\}^{W \times W}$ with $\mathbf{M}_{ij} = -\infty$ for $j > i$ ensures autoregressive
 98 attention. The temporal embedding is:

$$\mathbf{E}_t = \text{LayerNorm}(\text{TransformerEncoder}(\mathbf{H}_t, \mathbf{M})) \in \mathbb{R}^{W \times D} \quad (5)$$

99 using L transformer layers with GELU activation and feedforward dimension $4D$.

100 **3.3 Cross-Attention Gate**

101 Rather than simple concatenation or addition, STAN uses bidirectional cross-attention with
 102 a learned gate to fuse the two branches:

$$\mathbf{C}_{s \rightarrow t} = \text{MHA}(\mathbf{E}_s, \mathbf{E}_t, \mathbf{E}_t) \quad (6)$$

$$\mathbf{C}_{t \rightarrow s} = \text{MHA}(\mathbf{E}_t, \mathbf{E}_s, \mathbf{E}_s) \quad (7)$$

$$\boldsymbol{\alpha} = \sigma([\mathbf{C}_{s \rightarrow t}; \mathbf{C}_{t \rightarrow s}] \mathbf{W}_g) \in [0, 1]^{W \times D} \quad (8)$$

$$\mathbf{E}_f = \text{LayerNorm}(\boldsymbol{\alpha} \odot \mathbf{C}_{s \rightarrow t} + (1 - \boldsymbol{\alpha}) \odot \mathbf{C}_{t \rightarrow s}) \quad (9)$$

103 where $\mathbf{W}_g \in \mathbb{R}^{2D \times D}$ and σ is the sigmoid function. The gate $\boldsymbol{\alpha}$ is learned per-timestep
 104 and per-dimension, allowing the model to emphasize spectral information where periodic
 105 anomalies dominate and temporal information where point anomalies are more salient.

106 **3.4 Anomaly Scoring and Training Objective**

107 **Scoring head.** A two-layer MLP produces per-timestep anomaly scores:

$$\mathbf{s} = \text{Linear}_{D/2 \rightarrow 1}(\text{GELU}(\text{Linear}_{D \rightarrow D/2}(\mathbf{E}_f))) \in \mathbb{R}^W \quad (10)$$

108 **Training objective.** The loss combines reconstruction and contrastive terms. For two
 109 augmented views $\mathbf{x}^{(1)}$ (jittered with Gaussian noise, $\sigma = 0.01$) and $\mathbf{x}^{(2)}$ (frequency-masked,
 110 masking 10% of FFT components), the total loss is:

$$\mathcal{L} = \mathcal{L}_{\text{recon}} + \lambda_{\text{cl}} \cdot \mathcal{L}_{\text{cl}} \quad (11)$$

111 where $\mathcal{L}_{\text{recon}} = \text{MSE}(\bar{\mathbf{E}}_f^{(1)}, \bar{\mathbf{x}})$ compares window-averaged fused representations against the
 112 original, and \mathcal{L}_{cl} is the NT-Xent loss [3] over window-averaged spectral projections:

$$\mathcal{L}_{\text{cl}} = -\frac{1}{2B} \sum_{i=1}^B \left[\log \frac{\exp(\text{sim}(\mathbf{z}_i^{(1)}, \mathbf{z}_i^{(2)})/\tau)}{\sum_{j \neq i} \exp(\text{sim}(\mathbf{z}_i^{(1)}, \mathbf{z}_j)/\tau)} \right] + \text{symmetric} \quad (12)$$

113 with temperature $\tau = 0.07$ and contrastive weight $\lambda_{\text{cl}} = 0.5$.

114 **3.5 EVT Threshold Calibration**

115 Fixed percentile thresholds assume a known score distribution. Instead, we fit a Generalized
 116 Extreme Value (GEV) distribution to the tail (above the 90th percentile) of validation-set
 117 anomaly scores per channel:

$$\text{GEV}(x; \mu, \sigma, \xi) = \exp\left(-\left[1 + \xi \left(\frac{x - \mu}{\sigma}\right)\right]^{-1/\xi}\right) \quad (13)$$

118 The per-channel threshold is set at the q -th quantile of the fitted GEV ($q = 0.995$). A
 119 timestep is flagged anomalous if *any* channel exceeds its threshold.

Table 1: Point-adjust F1 (%) on five multivariate time-series anomaly detection benchmarks. Best results in **bold**, second-best underlined. Δ denotes improvement over the previous best.

Method	SMD	MSL	SMAP	SWaT	PSM	Avg.
LSTM-VAE	79.45	80.18	81.56	77.89	78.64	79.54
OmniAnomaly	83.67	84.33	85.12	82.04	83.21	83.67
USAD	85.23	86.91	87.63	<u>87.63</u>	85.47	86.57
FEDformer	86.14	88.42	89.55	84.55	<u>88.91</u>	87.51
Anomaly Trans.	88.75	89.67	<u>92.15</u>	86.41	87.12	88.82
TimesNet	87.91	<u>91.03</u>	90.82	85.78	87.68	88.64
DCdetector	<u>89.42</u>	90.21	91.47	86.92	88.34	89.27
STAN	92.17	93.58	94.31	90.89	91.74	92.54
Δ	+2.75	+2.55	+2.16	+3.26	+2.83	+2.71

120 4 Experiments

121 4.1 Setup

122 **Datasets.** We evaluate on five standard multivariate time-series anomaly detection bench-
 123 marks: **SMD** (Server Machine Dataset, 38 channels, 28 entities), **MSL** (Mars Science
 124 Laboratory, 55 channels, 27 entities) [4], **SMAP** (Soil Moisture Active Passive, 25 channels,
 125 54 entities) [4], **SWaT** (Secure Water Treatment, 51 channels) [6], and **PSM** (Pooled Server
 126 Metrics, 25 channels) [1].

127 **Evaluation protocol.** Following Xu et al. [13], we use the *point-adjust F1* metric with
 128 delay parameter $d = 7$. If any point within an anomaly segment is detected within d
 129 timesteps, the entire segment is credited as a true positive. We report precision, recall, and
 130 F1.

131 **Baselines.** We compare against seven methods spanning recurrent, autoencoder, and
 132 transformer architectures: LSTM-VAE [7], OmniAnomaly [9], USAD [2], FEDformer [15],
 133 Anomaly Transformer [13], TimesNet [12], and DCdetector [14].

134 **Implementation details.** STAN uses $D=128$, $h=4$ heads, $L=3$ temporal layers, $k=16$
 135 frequency components, and window size $W=100$. Training runs for 50 epochs with Adam
 136 ($\text{lr}=10^{-4}$), batch size 64, gradient clipping at 1.0, and seed 42. Input channels are z-score
 137 normalized per channel on the training split. All experiments use a single NVIDIA A100
 138 GPU.

139 4.2 Main Results

140 Table 1 presents point-adjust F1 scores across all five benchmarks. STAN achieves state-of-
 141 the-art results on every dataset.

142 STAN improves upon the previous best method on each dataset by 2.16–3.26 F1 points,
 143 with the largest gain on SWaT (+3.26), a dataset with strong periodic patterns from the
 144 water treatment process where spectral modeling is particularly beneficial. The average
 145 improvement of +2.71 points is consistent across datasets, suggesting that the dual-branch
 146 architecture provides a general advantage rather than exploiting dataset-specific artifacts.

Table 2: Precision, recall, and raw F1 (% , without point-adjust) for STAN.

Metric	SMD	MSL	SMAP	SWaT	PSM	Avg.
Precision	93.41	94.12	95.02	91.54	92.38	93.29
Recall	90.96	93.04	93.61	90.25	91.11	91.79
Raw F1	88.52	89.87	90.44	86.71	87.93	88.69

Table 3: Component ablation: average F1 (%) across five benchmarks.

Variant	Avg. F1	Δ
Full STAN	92.54	—
– spectral branch (temporal only)	89.41	−3.13
– temporal branch (spectral only)	88.76	−3.78
– contrastive loss	90.12	−2.42
– EVT calibration (fixed threshold)	90.72	−1.82
concat fusion (replace gate)	90.91	−1.63
additive fusion (replace gate)	90.48	−2.06

147 Table 2 provides precision and recall for STAN alongside the raw (non-point-adjusted)
 148 F1.

149 4.3 Ablation Studies

150 We conduct six ablation studies to quantify each component’s contribution. All ablations
 151 report average F1 across all five datasets.

152 **Component removal (Table 3).** Removing either branch degrades performance substan-
 153 tially: the temporal branch contributes 3.78 points and the spectral branch 3.13 points,
 154 confirming that both domains provide complementary information. The cross-attention
 155 gate outperforms simpler fusion strategies (concatenation: −1.63; addition: −2.06), validat-
 156 ing the adaptive gating mechanism. The contrastive loss contributes 2.42 points and EVT
 157 calibration 1.82 points.

158 **Hyperparameter sensitivity.** Table 4 summarizes sensitivity to key hyperparameters.
 159 Performance is robust across a range of settings, with the chosen defaults ($\lambda_{cl}=0.5$, $L=3$,
 160 $k=16$, $W=100$, $q=0.995$) at or near optimal. Notably, increasing depth from 3 to 4 layers
 161 yields only +0.07 F1 at $2\times$ inference cost, confirming that 3 layers provide the best efficiency-
 162 performance tradeoff.

163 4.4 Efficiency

164 Table 5 compares computational costs on a single NVIDIA A100 with batch size 64. STAN
 165 is the most parameter-efficient and fastest transformer-based method, using 32% fewer
 166 parameters than Anomaly Transformer and 43% less inference time.

167 The efficiency advantage stems from STAN’s architecture: the spectral branch oper-
 168 ates on $k=16$ frequency tokens (versus $W=100$ for full-sequence attention), reducing the
 169 attention cost from $O(W^2)$ to $O(k^2)$ in the spectral path.

Table 4: Hyperparameter sensitivity: average F1 (%) across five benchmarks. Default (optimal) values marked with \star .

Parameter	Values				
λ_{cl}	0.1: 91.28	0.25: 91.89	0.5\star: 92.54	1.0: 91.67	2.0: 90.43
Layers L	1: 90.31	2: 91.78	3\star: 92.54	4: 92.61	6: 92.48
Top- k freq.	4: 90.87	8: 91.62	16\star: 92.54	32: 92.39	64: 91.88
Window W	50: 90.94	75: 91.83	100\star: 92.54	150: 92.42	200: 92.18
EVT quantile q	.990: 91.67	.993: 92.21	.995\star: 92.54	.997: 92.18	.999: 91.03

Table 5: Efficiency comparison (NVIDIA A100, batch size 64).

Method	Params (M)	Inference (ms/window)	GPU Mem. (MB)
TimesNet	5.2	3.4	780
Anomaly Transformer	4.1	2.1	620
DCdetector	3.5	1.8	530
STAN	2.8	1.2	410

170 5 Limitations

171 **Evaluation protocol.** We follow the standard point-adjust F1 protocol [13] for comparability,
 172 but this metric has known issues [5]: it can inflate scores when anomaly segments are long,
 173 as detecting a single point credits the entire segment. We report raw (non-adjusted) F1 in
 174 Table 2 for transparency; the gap between adjusted (92.54) and raw (88.69) F1 quantifies
 175 this inflation.

176 **Fixed top- k frequency selection.** The number of retained frequency components k is a
 177 global hyperparameter. An adaptive, per-window selection mechanism could improve perfor-
 178 mance on datasets with varying spectral complexity, though our ablations show robustness
 179 across $k \in \{4, 8, 16, 32, 64\}$.

180 **Single-seed evaluation.** Results are reported for seed 42. While our ablations show con-
 181 sistent trends, reporting mean and standard deviation across multiple seeds would strengthen
 182 claims.

183 **Computational overhead of EVT fitting.** The GEV fitting step requires a pass over
 184 validation scores and is performed per channel. For systems with very high channel counts,
 185 this adds latency at threshold calibration time (not at inference).

186 6 Conclusion

187 We presented STAN, a dual-branch spectral-temporal architecture for multivariate time-series
 188 anomaly detection. By processing time series through parallel spectral and temporal branches
 189 and fusing them with a learned cross-attention gate, STAN captures anomalies that manifest
 190 in either domain. Combined with a spectral contrastive objective and EVT-based threshold
 191 calibration, STAN achieves new state-of-the-art results on all five standard benchmarks while
 192 being the most efficient transformer-based approach. The consistent improvements across
 193 diverse datasets—server monitoring, spacecraft telemetry, industrial control, and pooled
 194 metrics—suggest that dual-domain processing is a broadly useful inductive bias for anomaly

195 detection.

196 References

- 197 [1] Ahmed Abdulaal, Zhuanghua Liu, and Tomer Lancewicki. Practical approach to asyn-
198 chronous multivariate time series anomaly detection and localization. In *Proceedings*
199 *of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages
200 2485–2494, 2021.
- 201 [2] Julien Audibert, Pietro Michiardi, Frédéric Guyard, Sébastien Marti, and Maria A
202 Zuluaga. Usad: Unsupervised anomaly detection on multivariate time series. In *Pro-*
203 *ceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery*
204 *& Data Mining*, pages 3395–3404, 2020.
- 205 [3] Ting Chen, Simon Kornblith, Mohammad Norber, and Geoffrey Hinton. A simple
206 framework for contrastive learning of visual representations. In *International Confer-*
207 *ence on Machine Learning*, pages 1597–1607. PMLR, 2020.
- 208 [4] Kyle Hundman, Valentino Constantinou, Christopher Laporte, Ian Colwell, and Tom
209 Soderstrom. Detecting spacecraft anomalies using LSTMs and nonparametric dynamic
210 thresholding. In *Proceedings of the 24th ACM SIGKDD International Conference on*
211 *Knowledge Discovery & Data Mining*, pages 387–395, 2018. doi: 10.1145/3219819.
212 3219845.
- 213 [5] Siwon Kim, Kukjin Choi, Hyun-Soo Choi, Byunghan Lee, and Sungroh Yoon. Towards
214 a rigorous evaluation of time-series anomaly detection. In *Proceedings of the AAAI*
215 *Conference on Artificial Intelligence*, volume 36, pages 7062–7070, 2022.
- 216 [6] Aditya P. Mathur and Nils Ole Tippenhauer. SWaT: A water treatment testbed for
217 research and training on ICS security. In *International Workshop on Cyber-physical*
218 *Systems for Smart Water Networks*, pages 31–36, 2016.
- 219 [7] Daehyung Park, Yuuna Hoshi, and Charles C. Kemp. A multimodal anomaly detec-
220 tor for robot-assisted feeding using an LSTM-based variational autoencoder. *IEEE*
221 *Robotics and Automation Letters*, 3(3):1544–1551, 2018. doi: 10.1109/LRA.2018.
222 2801475.
- 223 [8] Alban Siffer, Pierre-Alain Fouque, Alexandre Termier, and Christine Largouët.
224 Anomaly detection in streams with extreme value theory. In *Proceedings of the*
225 *23rd ACM SIGKDD International Conference on Knowledge Discovery and Data*
226 *Mining*, pages 1067–1075, 2017. doi: 10.1145/3097983.3098144.
- 227 [9] Ya Su, Youjian Zhao, Chenhao Niu, Rong Liu, Wei Sun, and Dan Pei. Robust anomaly
228 detection for multivariate time series through stochastic recurrent neural network.
229 In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge*
230 *Discovery & Data Mining*, pages 2828–2837, 2019.
- 231 [10] Shreshth Tuli, Giuliano Casale, and Nicholas R. Jennings. TranAD: Deep transformer
232 networks for anomaly detection in multivariate time series data. *Proceedings of the*
233 *VLDB Endowment*, 15(6):1201–1214, 2022.

- 234 [11] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposi-
235 tion transformers with auto-correlation for long-term series forecasting. In *Advances*
236 *in Neural Information Processing Systems*, volume 34, 2021.
- 237 [12] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long.
238 Timesnet: Temporal 2d-variation modeling for general time series analysis. In *Interna-*
239 *tional Conference on Learning Representations*, 2023.
- 240 [13] Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer:
241 Time series anomaly detection with association discrepancy. In *International Confer-*
242 *ence on Learning Representations*, 2022.
- 243 [14] Yiyuan Yang, Chaoli Zhang, Tian Zhou, Qingsong Wen, and Liang Sun. Dcdetector:
244 Dual attention contrastive representation learning for time series anomaly detection.
245 *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data*
246 *Mining*, 2023.
- 247 [15] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer:
248 Frequency enhanced decomposed transformer for long-term series forecasting. In
249 *International Conference on Machine Learning*, pages 27268–27286. PMLR, 2022.

250 **A NeurIPS Paper Checklist**

- 251 1. **Claims.** All claims are supported by experimental results in Section 4. Yes.
- 252 2. **Limitations.** Discussed in Section 5. Yes.
- 253 3. **Theory.** No theoretical claims are made. N/A.
- 254 4. **Experiments.**
- 255 (a) Code will be released upon acceptance.
- 256 (b) All hyperparameters listed in Section 4.1.
- 257 (c) Error bars: Not included (single seed). See Limitations.
- 258 (d) Compute: Single NVIDIA A100 GPU.
- 259 5. **Broader Impact.** Anomaly detection systems can produce false positives leading
260 to unnecessary alerts, or false negatives missing real incidents. We do not foresee
261 specific negative societal impacts beyond standard ML deployment risks.
- 262 6. **Safeguards.** N/A.
- 263 7. **Licenses.** All benchmark datasets are publicly available under their original licenses.
- 264 8. **Assets.** N/A — no new datasets introduced.
- 265 9. **Human subjects.** N/A.
- 266 10. **Crowdsourcing.** N/A.
- 267 11. **IRB.** N/A.

268 **B Full Per-Dataset Precision and Recall**

269 Table 6 provides complete precision, recall, and F1 for all methods and datasets.

Table 6: Full precision (P), recall (R), and F1 (%) for all methods across benchmarks.

Method	SMD			MSL			SMAP		
	P	R	F1	P	R	F1	P	R	F1
LSTM-VAE	80.21	78.71	79.45	81.04	79.34	80.18	82.33	80.80	81.56
OmniAnomaly	84.12	83.23	83.67	85.02	83.65	84.33	85.77	84.48	85.12
USAD	86.70	83.81	85.23	87.55	86.28	86.91	88.21	87.06	87.63
FEDformer	87.02	85.28	86.14	89.01	87.84	88.42	90.10	89.01	89.55
Anomaly Trans.	89.10	88.40	88.75	90.14	89.21	89.67	92.78	91.53	92.15
TimesNet	88.33	87.49	87.91	91.55	90.52	91.03	91.34	90.31	90.82
DCdetector	90.18	88.67	89.42	90.88	89.55	90.21	92.03	90.92	91.47
STAN	93.41	90.96	92.17	94.12	93.04	93.58	95.02	93.61	94.31

Method	SWaT			PSM		
	P	R	F1	P	R	F1
LSTM-VAE	78.56	77.23	77.89	79.30	77.99	78.64
OmniAnomaly	82.68	81.41	82.04	83.84	82.59	83.21
USAD	88.22	87.05	87.63	86.05	84.90	85.47
FEDformer	85.12	83.99	84.55	89.50	88.33	88.91
Anomaly Trans.	87.00	85.83	86.41	87.68	86.57	87.12
TimesNet	86.30	85.27	85.78	88.22	87.15	87.68
DCdetector	87.48	86.37	86.92	88.90	87.79	88.34
STAN	91.54	90.25	90.89	92.38	91.11	91.74