

MOMEMTO: PATCH-BASED MEMORY GATE MODEL IN TIME SERIES FOUNDATION MODEL

Anonymous authors

Paper under double-blind review

ABSTRACT

Recently reconstruction-based deep models have been widely used for time series anomaly detection, but as their capacity and representation capability increase, these models tend to over-generalize, often reconstructing unseen anomalies accurately. Prior works have attempted to mitigate this by incorporating a memory architecture that stores prototypes of normal patterns. Nevertheless, these approaches suffer from high training costs and have yet to be effectively integrated with time series foundation models (TFMs). To address these challenges, we propose **MOMEMTO**, an improved variant of TFM for anomaly detection, enhanced with a patch-based memory module to mitigate over-generalization. The memory module is designed to capture representative normal patterns from multiple domains and enables a single model to be jointly fine-tuned across multiple datasets through a multi-domain training strategy. **MOMEMTO** initializes memory items with latent representations from a pre-trained encoder, organizes them into patch-level units, and updates them via an attention mechanism. We evaluate our method using 23 univariate benchmark datasets. Experimental results demonstrate that **MOMEMTO**, as a single model, achieves higher scores on AUC and VUS metrics compared to baseline methods, and further enhances the performance of its backbone TFM, particularly in few-shot learning scenarios.

1 INTRODUCTION

Complex cyber-physical systems operate continuously, accumulating vast volumes of sequential time series data from multiple sensors. Monitoring these systems and detecting anomalies at an early stage is highly beneficial in terms of operational cost. In real-world anomaly detection scenarios, it is common to simultaneously handle multi-domain datasets that originate from different application areas (e.g., distinct sensor types or process configurations). Such scenarios frequently involve data imbalance and limited labeled data. We formulate this problem as an unsupervised learning task to tackle these challenges.

Reconstruction-based deep models are widely used for unsupervised time series anomaly detection. In this approach, an encoder-decoder network is trained with a self-supervised objective to reconstruct input sequences. These models aim to reproduce normal patterns precisely, while generating higher reconstruction errors for anomalies. Various models have been proposed within this paradigm, including OmniAnomaly (Su et al., 2019), Anomaly Transformer (Xu et al., 2022), TranAD (Tuli et al., 2022), and TimesNet (Wu et al., 2023). Recent models have achieved strong performance due to the representation learning capability of neural networks. However, existing reconstruction-based approaches often suffer from over-generalization, where anomalous inputs are reconstructed with high accuracy (Gong et al., 2019; Park et al., 2020). This issue arises when the encoder captures distinctive features of anomalies or when the decoder has sufficient capacity to reconstruct abnormal representations.

Recent research has increasingly focused on TFMs, which are pre-trained on large-scale time series data from a variety of sources. TFMs can be adapted to new tasks with minimal additional training through approaches such as multi-task learning, few-shot learning, and prompting. However, applying TFMs to anomaly detection presents two main challenges. First, most TFMs, such as MOIRAI (Woo et al., 2024), TimesFM (Das et al., 2024), Chronos (Ansari et al., 2024), Time-MoE (Shi et al., 2024), and Sundial (Liu et al., 2025) are primarily designed for time series forecasting. They often

054 employ decoder-only architectures that limit their suitability for detection tasks. Second, TFMs used
055 for anomaly detection typically adopt reconstruction-based approaches, leaving them vulnerable to
056 over-generalization.

057 To address these challenges, we adapt the Gate memory module of MEMTO (Song et al., 2023),
058 a reconstruction-based deep model for time series anomaly detection. This module stores memory
059 items that represent the prototypical features of normal patterns. When anomalies are reconstructed
060 using these normal memory items, the resulting outputs resemble normal samples, thus alleviat-
061 ing over-generalization. However, MEMTO is highly sensitive to memory initialization, requiring
062 encoder pre-training to obtain informative and stable representations. It also follows a one-model-
063 per-dataset framework, which leads to high training costs when handling multiple datasets. In ad-
064 dition, MEMTO uses a point-level memory rather than a patch-based structure, which reduces its
065 effectiveness in detecting interval or periodic anomalies (Shen, 2025).

066 Furthermore, we adopt the pre-trained encoder from MOMENT (Goswami et al., 2024), a TFM
067 trained on large-scale time series datasets. MOMENT employs a patch-level masked representa-
068 tion learning strategy, where input sequences are divided into fixed-length patches and trained to
069 reconstruct the masked patches. Its encoder, built on the T5 architecture (Raffel et al., 2020), can be
070 coupled with task-specific decoders, providing highly transferable time series representations and
071 flexibility across diverse tasks. However, like other TFMs, MOMENT adopts multi-domain pre-
072 training, but achieving competitive performance typically requires fine-tuning a separate model for
073 each target dataset. This limitation prevents it from fully exploiting the benefits of multi-domain
074 representations. Moreover, as a reconstruction-based model for anomaly detection, MOMENT still
075 remains vulnerable to over-generalization.

076 Building upon previous methods, we introduce MOMEMTO, an improved variant of MOMENT
077 adapted for reconstruction-based time series anomaly detection. The model integrates the pre-trained
078 encoder from MOMENT with the patch-based memory module. Unlike MEMTO, our model ben-
079 efits from MOMENT’s encoder, which already possesses sufficient representational capacity. This
080 provides well-initialized features for the memory module and alleviates MEMTO’s sensitivity to
081 memory initialization. The patch-based memory module stores encoder outputs that are organized
082 at the patch level. Whereas the original memory items stored prototypical features of normal patterns
083 restricted to a single dataset, the patch-based memory items are designed to capture representative
084 normal patterns from multiple domains. This design alleviates over-generalization while enabling
085 the model to learn effectively in diverse domains. Specifically, multi-domain training refers to jointly
086 fine-tuning a single model across datasets from diverse application areas, rather than training a sep-
087 arate model for each dataset. This strategy facilitates knowledge sharing across domains and allows
088 more effective use of limited data, while also reducing computational cost in both training time
089 and memory usage. Under the multi-domain training strategy, a single MOMEMTO model achieves
090 strong results on 23 univariate benchmark datasets compared to baseline methods. In few-shot learn-
091 ing scenarios, it demonstrates significant improvements over its backbone TFM. Our contributions
092 can be summarized as follows:

- 092 • We introduce a patch-based memory architecture to support TFMs, which stores features
093 from multi-domain datasets and is broadly applicable to various patch-based methods.
- 094 • MOMEMTO extends MOMENT into a TFM specialized for time series anomaly detec-
095 tion, supporting multi-domain training and leveraging the patch-based memory module to
096 mitigate over-generalization.
- 097 • MOMEMTO achieves strong results on 23 univariate benchmark datasets and demonstrates
098 superior few-shot anomaly detection performance.

101 2 RELATED WORK

102 **Transformers and patching for time series** Transformers (Vaswani et al., 2017) have demon-
103 strated strong performance on sequential data processing. For time series analysis, Transformers
104 benefit from the self-attention mechanism and are used to capture reliable long-range temporal de-
105 pendencies (Kitaev et al., 2020; Li et al., 2021; Zhou et al., 2021; Liu et al., 2024). Among them,
106 PatchTST (Nie et al., 2023) is a Transformer-based model specifically designed for time series anal-
107

108 ysis. Its first strategy is patching, in which input sequences are segmented into subsequence-level
 109 patches, preserving local semantic information while reducing computational complexity. Its second
 110 strategy is a channel-independent processing scheme, where each univariate series is processed in-
 111 dependently with shared Transformer weights, minimizing inter-channel interference and effectively
 112 capturing distinct temporal dynamics. Together, these designs enhance both the efficiency and the
 113 representation capability of Transformers for time series tasks.

114
 115 **Memory-guided deep models** Recently, memory architectures have been introduced to enhance
 116 neural models by enabling external storage and retrieval of long-term information. They have been
 117 applied in diverse domains such as natural language processing, including large-scale language mod-
 118 eling (Lample et al., 2019), retrieval-augmented generation (Lewis et al., 2020), and long-context un-
 119 derstanding (Wang et al., 2023), and in computer vision tasks like video captioning (Lei et al., 2020),
 120 video object segmentation (Oh et al., 2019) as well as in reinforcement learning for episodic memory
 121 and sample-efficient decision making (Le et al., 2021). For anomaly detection, several approaches
 122 in computer vision employ memory modules to store features of normal patterns (Zhou et al.,
 123 2023). MemAE (Gong et al., 2019) integrates a memory module into an autoencoder. MNAD (Park
 124 et al., 2020) introduces a weighted memory-update strategy to capture temporal anomalies in videos.
 125 Based on these concepts, MEMTO is the first model to introduce a memory architecture into time
 126 series anomaly detection, employing a data-driven update mechanism. More recently, H-PAD (Shen,
 127 2025) extends this idea by learning hybrid patch- and period-level prototypes to capture both interval
 and periodic anomalies.

128 MOMEMTO extends MEMTO by refining its memory mechanism to explicitly capture patch-level
 129 patterns. Our patch-based memory module is designed based on patching and channel-independent
 130 strategies, allowing selective updates of memory items rather than updating all items as in MEMTO.
 131 This design enables the memory module to learn more informative features of normal patterns while
 132 also improving the computational efficiency of memory updates compared to the original memory
 133 module (Nie et al., 2023; Song et al., 2023).

135 3 METHOD

136
 137 We consider multiple univariate time series datasets, which are categorized into distinct domains.
 138 Each domain exhibits unique temporal characteristics such as length, sampling frequency, trend,
 139 seasonality, and noise. A single univariate time series from a dataset X is denoted as x . When
 140 divided into non-overlapping patches of fixed length L , the sequence is represented as $x \in \mathbb{R}^{P \times L}$,
 141 where P is the number of observed patches. If the number of observed patches P is smaller than
 142 the maximum number of patches N that the model can process, we apply zero-padding so that the
 143 sequence matches the shape $\mathbb{R}^{N \times L}$.

145 3.1 MOMEMTO

146
 147 Figure 1 illustrates the overall architecture of MOMEMTO, which mainly consists of an en-
 148 coder–decoder with the patch-based memory module. The input time series x is first processed
 149 by the encoder. The encoder output representation $q \in \mathbb{R}^{P \times d_{model}}$ is used as queries to interact
 150 with the memory module, where d_{model} denotes the embedding dimension. The query vectors
 151 are compared with the stored memory items to retrieve the most relevant ones, and the selected
 152 memory items are subsequently updated conditioned on the queries. The queries are then refined
 153 via memory-based attention. The updated queries $\tilde{q} \in \mathbb{R}^{P \times d_{model}}$ are combined with the original
 154 queries q (i.e., $[q; \tilde{q}] \in \mathbb{R}^{P \times 2d_{model}}$) and passed to the decoder. The decoder maps the combined
 155 query representation to the input space to reconstruct the time series $\hat{x} \in \mathbb{R}^{P \times L}$.

156 3.1.1 PRE-TRAINED ENCODER AND DECODER

157
 158 We adopt the pre-trained encoder from MOMENT-large as the backbone for our model. The encoder
 159 receives a univariate time series $x \in \mathbb{R}^{N \times L}$ together with a mask vector of length N , where each
 160 entry is 0 or 1 to indicate unobserved and observed patches, respectively. Each input patch is treated
 161 as a token and embedded into a vector of length d_{model} . Patch-level self-attention is then applied to
 capture temporal dependencies and generate contextualized representations. For reconstruction, we

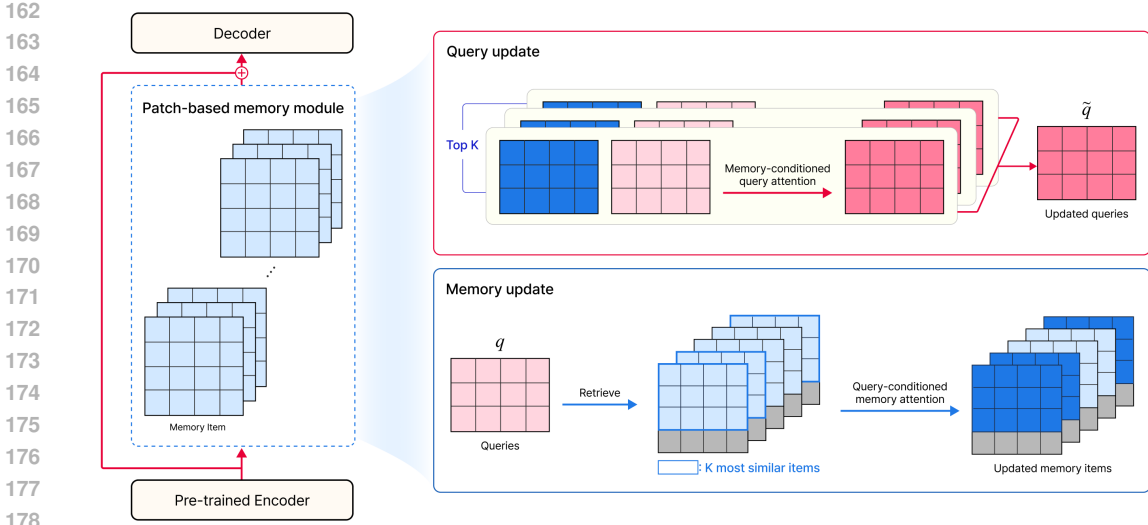


Figure 1: Architecture of MOMEMTO.

use a shared lightweight decoder implemented as two fully-connected layers. This design mitigates the risk of over-generalization and promotes reliance on the encoder’s representational capacity, while enabling the model to generalize across datasets with varying sequence lengths.

3.1.2 PATCH-BASED MEMORY MODULE

To alleviate over-generalization in reconstruction-based models, we introduce the patch-based memory module as a redesigned variant of MEMTO’s original memory module. It naturally aligns with the encoder’s patch representations. The module consists of M memory items $\{m_i\}_{i=1}^M$, where each memory item $m_i \in \mathbb{R}^{N \times d_{model}}$ stores a prototypical normal pattern at the patch level. The module operates in a data-driven manner through five stages: memory initialization, memory alignment, memory update, query update, and query alignment. Through these stages, the memory items and queries are progressively refined and aligned with representative normal patterns.

Memory initialization Each memory item m_i is initialized with a specific domain, which may encompass multiple datasets. It is stored as the mean of encoder representations from S_i instances randomly sampled across the datasets belonging to that domain. These encoder representations are denoted by $\{q_s\}_{s=1}^{S_i}$. The initial value of m_i is then computed as

$$m_i = \frac{1}{S_i} \sum_{s=1}^{S_i} q_s, \quad i = 1, \dots, M.$$

To ensure stable updates, memory items are consistently L2-normalized at the patch level.

Memory alignment After initializing the memory items, we first identify the observed patches in the encoder output q using the input mask and retrieve the corresponding subset of representations. To ensure consistent interaction, each memory item m_i is aligned with the same subset of observed patches, and the queries are likewise L2-normalized, denoted as \hat{q} .

$$\begin{array}{ccc} q \in \mathbb{R}^{N \times d_{model}} & \xrightarrow{\text{retrieve}} & q \in \mathbb{R}^{P \times d_{model}} \\ m_i \in \mathbb{R}^{N \times d_{model}} & & m_i \in \mathbb{R}^{P \times d_{model}} \end{array}$$

This alignment process allows the model to handle time series of varying lengths flexibly and selectively update memory items, enhancing adaptability across datasets.

Memory update We then compute the similarity between the normalized queries \hat{q} and the corresponding memory slices using dot products. Instead of updating every memory item, we select the top- K memory items with the highest similarity scores. These selected memory items $m_{1st}, m_{2nd}, \dots, m_{K-th}$ are updated through query-based patch-level attention, where m_{i-th} denotes the memory item with the i -th highest similarity to \hat{q} . The resulting updates are proportionally applied to the selected memory items in a data-driven way, with the remaining update procedure following MEMTO’s standard strategy (details in Appendix B.4).

Query update In the query update stage, we refine \hat{q} into an updated representation \tilde{q} using memory-based patch-level attention. This process is performed only on the top- K memory items selected in the previous stage. For each selected memory item m_{i-th} , we compute an intermediate representation \tilde{q}_{i-th} . The final refined representation \tilde{q} is obtained by taking the weighted sum of the K intermediate representations:

$$\tilde{q} = w_{1st} \cdot \tilde{q}_{1st} + w_{2nd} \cdot \tilde{q}_{2nd} + \dots + w_{K-th} \cdot \tilde{q}_{K-th}$$

where w_{i-th} denotes the similarity score between \hat{q} and m_{i-th} , computed in the previous stage.

Query alignment The resulting \tilde{q} is then combined with the original q in a patch-aligned manner, where each patch-level memory item contributes only to the corresponding patch position. This combined representation is then used as the input to the shared decoder. This ordered memory structure preserves the positional structure of the sequence and captures the underlying patterns in each patch.

3.2 MULTI-DOMAIN TRAINING

Unlike the conventional one-model-per-dataset paradigm, we train a single unified model that is jointly optimized on the entire dataset collection, enabling it to detect anomalies in time series with heterogeneous lengths and characteristics. During multi-domain training, the patch-based memory module is updated in a data-driven manner, allowing memory items to adapt to representative normal patterns from different domains. At initialization, the number of memory items is set equal to the number of user-defined domains, which provides a balanced starting point but does not enforce a strict one-to-one correspondence. Through data-driven updates, a single memory item may accumulate information from multiple domains and evolve into a domain-general feature. This training strategy facilitates knowledge sharing while maintaining a single model architecture in which the memory captures diverse normal patterns.

4 EXPERIMENTS

Datasets We evaluate our model on the TSB-AD-U benchmark (Liu & Paparrizos, 2024), consisting of 870 univariate time series across 23 datasets. The benchmark is constructed through an organized curation process: both univariate and multivariate datasets are collected, with multivariate data converted into univariate channels. Channels unrelated to anomaly labels, such as categorical, binary, or noisy channels, are discarded. Mislabelled or low-quality series are removed through algorithmic tests and human review, and balanced sampling is applied to ensure fairness. The collection also includes sequences from public benchmarks such as MSL (Mars Science Laboratory rover), SMD (Server Machine Dataset), and the UCR Anomaly Archive (Wu & Keogh, 2022). Further dataset details are provided in Appendix A.

Implementation details We generate subsequences by applying a non-overlapping window with a length of 512 to obtain fixed-length inputs for each time series. We evaluate anomaly detection performance using threshold-independent metrics, such as the Area Under Curve (AUC) and Volume Under the Surface (VUS) (Paparrizos et al., 2022). More detailed information on hyperparameter settings can be found in Appendix B.

4.1 MAIN RESULTS

In this experiment, we evaluate the performance of MOMEMTO on time series anomaly detection by comparing it with 13 baseline methods. The baselines consist of three categories: classic approaches (OCSVM (Schölkopf et al., 1999), LOF (Breunig et al., 2000), Isolation Forest (Liu et al., 2008)), Deep learning models (LSTM-AD (Malhotra et al., 2015), Donut (Xu et al., 2018), OmniAnomaly (Su et al., 2019), TranAD (Tuli et al., 2022), Anomaly Transformer (Xu et al., 2022), MEMTO, TimesNet (Wu et al., 2023), KAN-AD (Zhou et al., 2025), DADA (Shentu et al., 2025)), and pre-trained TFMs (TimesFM, Chronos, MOMENT).

Table 1: Model performance comparison on the benchmark. Metrics are reported as percentages (%). Subscripts indicate training settings: *md* = multi-domain joint fine-tuning of a single model across all datasets; models without a subscript denote the one-model-per-dataset setting.

Category	Model	AUC-PR	AUC-ROC	VUS-PR	VUS-ROC
Classical	LOF	14.24	58.17	16.87	67.74
	OCSVM	15.85	65.13	22.51	73.49
	Isolation Forest	29.04	<u>71.18</u>	29.80	77.51
Deep learning	MEMTO	6.90	52.05	11.30	56.30
	Donut	15.28	55.97	20.78	68.56
	TimesNet	18.10	61.19	25.91	72.03
	Anomaly Transformer	19.02	62.53	23.61	73.05
	TranAD	24.67	63.66	29.95	73.87
	KAN-AD	29.20	69.54	30.25	76.30
	OmniAnomaly	29.87	69.45	32.15	76.49
	LSTM-AD	31.77	67.70	32.66	75.74
	DADA	<u>33.59</u>	70.42	<u>34.06</u>	<u>77.62</u>
Foundation	Chronos	26.43	65.64	27.18	72.91
	TimesFM	28.43	67.29	29.88	74.20
	MOMENT	29.45	69.97	29.73	76.84
	MOMEMTO	32.83	70.00	33.23	77.49
	MOMEMTO_{md}	36.35	74.83	37.62	80.95

Table 1 shows the evaluation results on the benchmark. Higher values in these metrics indicate more accurate anomaly detection. Among all compared models, **MOMEMTO_{md}**, as a single model, achieves the best performance across all evaluation criteria. Moreover, reconstruction-based TFMs such as MOMENT and MOMEMTO yield stronger results than other TFMs that are primarily designed for forecasting tasks. A direct comparison with its backbone TFM, MOMENT, shows that MOMEMTO consistently surpasses its backbone across metrics and settings. This demonstrates that the patch-based memory module effectively mitigates over-generalization and improves detection capability, even when built upon the same pre-trained encoder, highlighting the importance of architectural design tailored to anomaly detection.

Few-shot learning To evaluate the robustness of our approach under limited data scenarios, we conduct a few-shot learning experiment. In this setting, the ratio of training data is gradually varied from 10% to 90%, and the performance of MOMEMTO is compared with MOMENT under the condition that only a small fraction of samples is available for fine-tuning.

Figure 2 shows the effect of training data ratio on AUC-PR. Across all ratios, MOMEMTO consistently outperforms MOMENT, highlighting the advantage of the patch-based memory module. The gap is substantial in the few-shot regime (10%–30%), where both the one-model-per-dataset and multi-domain variants achieve noticeably higher scores than the corresponding MOMENT variants. As the proportion of training data increases, MOMEMTO continues to improve in both settings, while MOMENT shows marginal or no further improvement. Moreover, MOMEMTO_{md} achieves consistently higher performance than MOMEMTO without multi-domain training, indicating that multi-domain training further strengthens the memory module by improving data efficiency. Overall, these results demonstrate that MOMEMTO is robust under limited data and benefits even more from additional data when trained in the multi-domain setting.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

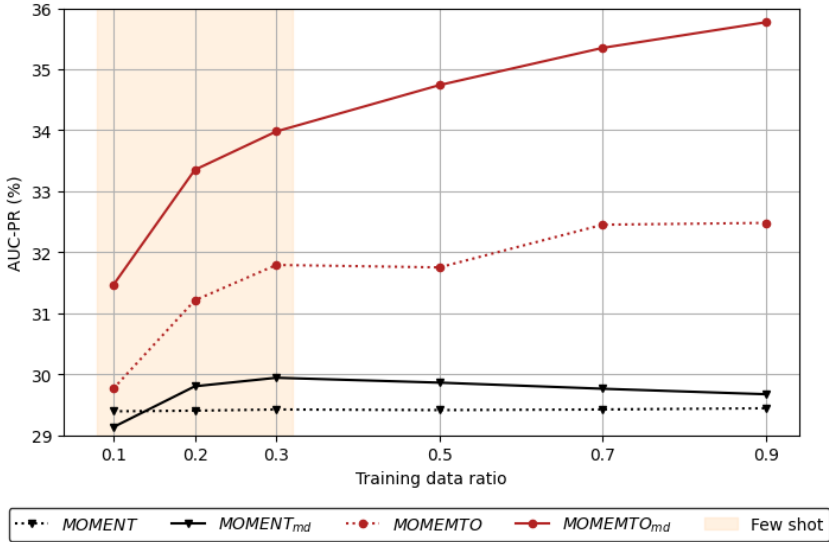


Figure 2: Effect of the training data ratio on AUC-PR.

4.2 MODEL ANALYSIS

Ablation study To investigate the contribution of each component in MOMEMTO, we conduct an ablation study by varying the encoder initialization (scratch vs. pre-trained) and the use of the patch-based memory module.

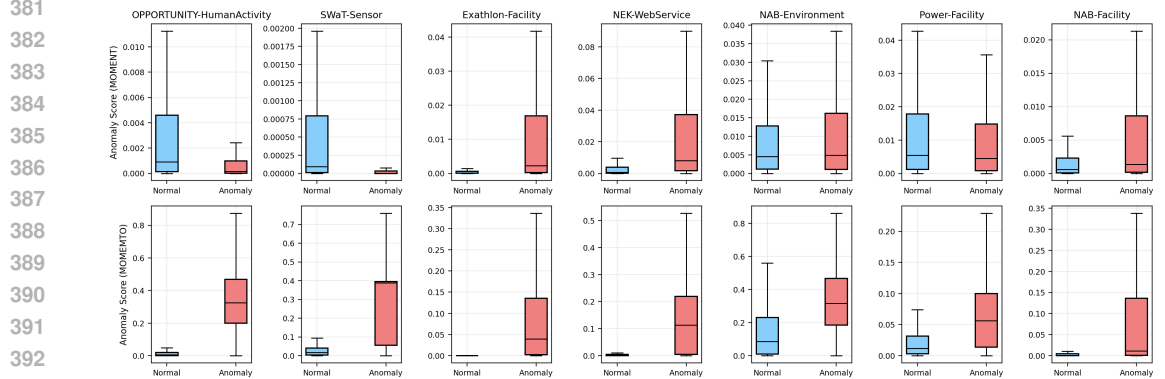
As shown in Table 2, both components contribute significantly to performance improvement. Replacing a scratch encoder with a pre-trained one improves AUC-PR (22.44 → 29.70), while adding the memory module to a scratch encoder yields an even larger gain (22.44 → 33.63). The combination of both components achieves the best performance, reaching an AUC-PR of 36.35. Similar trends are consistently observed across other metrics. These results indicate that the two components provide complementary benefits: the patch-based memory module mitigates over-generalization, and the pre-trained encoder provides a strong representational foundation that amplifies the memory module’s effectiveness.

Table 2: Ablation results of pre-trained encoder and patch-based memory module. ‘PMM’ denotes the patch-based memory module.

Encoder	PMM	AUC-PR	AUC-ROC	VUS-PR	VUS-ROC
scratch	×	22.44	67.35	18.90	66.50
scratch	○	33.63	72.17	34.65	78.78
pre-trained	×	29.70	69.60	30.69	76.58
pre-trained	○	36.35	74.83	37.62	80.95

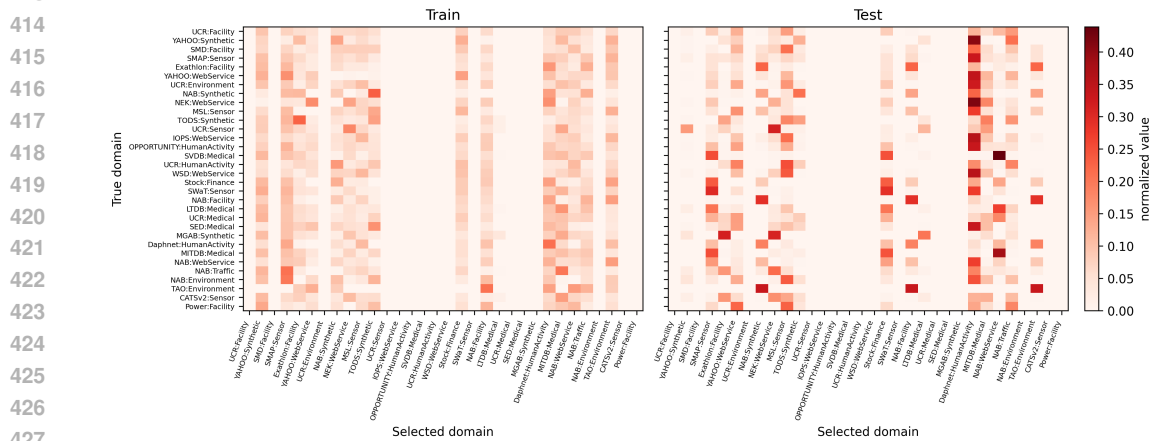
Comparison of anomaly score distributions Figure 3 presents the distributions of anomaly scores for normal and anomalous samples across a subset of domains. Each boxplot shows the distributions of normal and anomalous samples within a domain, enabling direct comparison of the models’ discriminative ability. Additional results for the remaining domains are provided in the Appendix C.5. Overall, MOMENT exhibits limited separation between normal and anomalous samples, with substantial overlap in their score distributions. In contrast, MOMEMTO often yields a clearer margin between the two groups: anomaly scores for anomalous samples tend to shift toward higher values, whereas those of normal samples remain clustered around low values. This tendency

378 suggests that MOMEMTO improves the distinction between normal and anomalous patterns in the
 379 majority of domains.
 380



381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394 Figure 3: Comparison of anomaly score distributions between normal (blue) and anomalous (red)
 395 samples across a subset of domains. The upper row shows the results obtained with the backbone
 396 model MOMENT, while the lower row corresponds to our proposed model.

399 **Data-driven memory utilization** We visualize in Figure 4 how memory items are accessed for
 400 subsequences from different domains. Each heatmap displays the accumulated similarity between
 401 the true domain of the input and the selected memory domains, normalized such that each row sums
 402 to one. The left illustrates updates during the entire training, and the right shows references by new
 403 inputs at test time. With the data-driven strategy, no domain labels are provided during training
 404 or testing. Updates tend to concentrate on a subset of memory items, rather than being evenly
 405 distributed across domains, and these items are consistently referenced at test time. To validate
 406 the effectiveness of our strategy, we compare it with two alternative strategies: (i) freezing the
 407 memory items without updates and (ii) updating only the memory item initially assigned to the
 408 input’s domain during training and testing. Table 3 summarizes the performance comparison. Even
 409 though the initial allocation of memory items to domains does not remain aligned during training,
 410 our data-driven strategy still achieves higher performance while utilizing fewer memory items than
 411 alternative strategies. This indicates that the memory module has learned a stable set of domain-
 412 general features, which are effectively reused across unseen inputs.



413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428 Figure 4: Visualization of row-normalized accumulated domain-to-memory similarity over the en-
 429 tire training process (left) and at test time (right). Each heatmap shows how subsequences from a
 430 true domain reference memory items across domains.
 431

Table 3: Performance comparison of different memory update strategies. *No update*: memory frozen after initialization; *Own-domain update*: update only the item initially assigned to the input’s domain; *Data-driven update (ours)*: update the similarity-retrieved item regardless of initial domain. At test time, both baselines use the memory item corresponding to the input’s domain.

Strategy	AUC-PR	AUC-ROC	VUS-PR	VUS-ROC
No update	34.97	73.77	36.91	80.19
Own-domain update	35.16	73.64	36.60	80.06
Data-driven update (ours)	36.35	74.83	37.62	80.95

Computational efficiency We measure the training time, inference time, and model size under different configurations of the patch-based memory module (PMM) and the multi-domain training paradigm, as summarized in Table 4. Here, *context switching* refers to the additional overhead incurred when multiple separately trained models must be loaded and executed sequentially, such as repeated model loading, GPU memory allocation, and parameter initialization. While the inference time remains stable at around 54 seconds across all configurations, the total training time and model size differ drastically. Training separate models introduces excessive context switching overhead (around 1,300 seconds) and substantially enlarges the total model size beyond 450 GB, which is computationally expensive. In contrast, our unified model with PMM achieves efficient total training (87.39 seconds) and maintains a compact size (1.3 GB).

Table 4: Comparison of training, context switching, inference, and model size under different configurations of the patch-based memory module (PMM) and multi-domain training. The configuration without PMM is equivalent to MOMENT.

	Time (sec)				Model size (GB)
	Training	Context switching	Total training	Inference	
MOMEMTO	82.22	5.17	87.39	54.61	1.30
w/o PMM	75.83	4.59	80.42	54.23	1.29
w/o multi-domain training	135.77	1311.51	1447.28	54.54	455.74
w/o PMM & multi-domain training	108.96	1294.26	1403.22	54.29	451.64

5 CONCLUSION AND FUTURE WORK

We introduce MOMEMTO, a time series foundation model for anomaly detection enhanced with a patch-based memory module. By leveraging a pre-trained encoder and the patch-based memory module, our approach mitigates over-generalization and achieves superior performance across 23 benchmark datasets, with notable gains in few-shot and multi-domain settings. While MOMEMTO shows promising results, several important aspects remain unexplored. The experiments focus primarily on a univariate benchmark, and the study does not include a theoretical analysis of MOMEMTO. Future work will investigate these aspects and broaden the applicability of our approach.

REFERENCES

Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Syndar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, Jasper Zschiegner, Danielle C. Maddix, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. Chronos: Learning the language of time series. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=gerNCVqqqTR>.

- 486 Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, and Jörg Sander. Lof: identifying density-
487 based local outliers. In *Proceedings of the 2000 ACM SIGMOD International Conference on*
488 *Management of Data*, SIGMOD '00, pp. 93–104, New York, NY, USA, 2000. Association for
489 Computing Machinery. ISBN 1581132174. doi: 10.1145/342009.335388. URL [https://](https://doi.org/10.1145/342009.335388)
490 doi.org/10.1145/342009.335388.
- 491 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model
492 for time-series forecasting. In *Proceedings of the 41st International Conference on Machine*
493 *Learning*, ICML'24. JMLR.org, 2024.
- 494 Dong Gong, Lingqiao Liu, Lê Vương, Budhaditya Saha, Moussa Mansour, Svetha Venkatesh, and
495 Anton Hengel. Memorizing normality to detect anomaly: Memory-augmented deep autoencoder
496 for unsupervised anomaly detection. pp. 1705–1714, 10 2019. doi: 10.1109/ICCV.2019.00179.
- 497
498 Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.
499 Moment: a family of open time-series foundation models. In *Proceedings of the 41st International*
500 *Conference on Machine Learning*, ICML'24. JMLR.org, 2024.
- 501
502 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In
503 *International Conference on Learning Representations*, 2020. URL [https://openreview.](https://openreview.net/forum?id=rkgNKkHtvB)
504 [net/forum?id=rkgNKkHtvB](https://openreview.net/forum?id=rkgNKkHtvB).
- 505 Guillaume Lample, Alexandre Sablayrolles, Marc'Aurelio Ranzato, Ludovic Denoyer, and Herve
506 Jegou. *Large memory layers with product keys*. Curran Associates Inc., Red Hook, NY, USA,
507 2019.
- 508
509 Hung Le, Thommen Karimpanal George, Majid Abdolshah, Truyen Tran, and Svetha Venkatesh.
510 Model-based episodic memory induces dynamic hybrid controls. In *Proceedings of the 35th*
511 *International Conference on Neural Information Processing Systems*, NIPS '21, Red Hook, NY,
512 USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- 513
514 Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara Berg, and Mohit Bansal. MART: Memory-
515 augmented recurrent transformer for coherent video paragraph captioning. In Dan Jurafsky,
516 Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meet-*
517 *ing of the Association for Computational Linguistics*, pp. 2603–2614, Online, July 2020. As-
518 sociation for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.233. URL [https:](https://aclanthology.org/2020.acl-main.233/)
[//aclanthology.org/2020.acl-main.233/](https://aclanthology.org/2020.acl-main.233/).
- 519
520 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
521 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
522 Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the*
523 *34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook,
524 NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 524
525 Zhihan Li, Youjian Zhao, Jiaqi Han, Ya Su, Rui Jiao, Xidao Wen, and Dan Pei. Multivariate time
526 series anomaly detection and interpretation using hierarchical inter-metric and temporal embed-
527 ding. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data*
Mining, pp. 3220–3230, 2021.
- 528
529 Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 Eighth IEEE International*
530 *Conference on Data Mining*, pp. 413–422, 2008. doi: 10.1109/ICDM.2008.17.
- 531
532 Qinghua Liu and John Paparrizos. The elephant in the room: Towards a reliable time-
533 series anomaly detection benchmark. In A. Globerson, L. Mackey, D. Belgrave, A. Fan,
534 U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Process-*
535 *ing Systems*, volume 37, pp. 108231–108261. Curran Associates, Inc., 2024. URL
536 [https://proceedings.neurips.cc/paper_files/paper/2024/file/](https://proceedings.neurips.cc/paper_files/paper/2024/file/c3f3c690b7a99fba16d0efd35cb83b2c-Paper-Datasets_and_Benchmarks_Track.pdf)
537 [c3f3c690b7a99fba16d0efd35cb83b2c-Paper-Datasets_and_Benchmarks_](https://proceedings.neurips.cc/paper_files/paper/2024/file/c3f3c690b7a99fba16d0efd35cb83b2c-Paper-Datasets_and_Benchmarks_Track.pdf)
[Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/c3f3c690b7a99fba16d0efd35cb83b2c-Paper-Datasets_and_Benchmarks_Track.pdf).
- 538
539 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.
itransformer: Inverted transformers are effective for time series forecasting, 2024. URL [https:](https://arxiv.org/abs/2310.06625)
[//arxiv.org/abs/2310.06625](https://arxiv.org/abs/2310.06625).

- 540 Yong Liu, Guo Qin, Zhiyuan Shi, Zhi Chen, Caiyin Yang, Xiangdong Huang, Jianmin Wang, and
541 Mingsheng Long. Sundial: A family of highly capable time series foundation models. *arXiv*
542 *preprint arXiv:2502.00816*, 2025.
- 543 Pankaj Malhotra, Lovekesh Vig, Gautam Shroff, and Puneet Agarwal. Long short term memory
544 networks for anomaly detection in time series. 04 2015.
- 546 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth
547 64 words: Long-term forecasting with transformers. In *International Conference on Learning*
548 *Representations*, 2023.
- 549 Seoung Wug Oh, Joon Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using
550 space-time memory networks. In *Proceedings - 2019 International Conference on Computer*
551 *Vision, ICCV 2019*, Proceedings of the IEEE International Conference on Computer Vision, pp.
552 9225–9234, United States, October 2019. Institute of Electrical and Electronics Engineers Inc.
553 doi: 10.1109/ICCV.2019.00932. Publisher Copyright: © 2019 IEEE.; 17th IEEE/CVF Inter-
554 national Conference on Computer Vision, ICCV 2019 ; Conference date: 27-10-2019 Through
555 02-11-2019.
- 556 John Paparrizos, Paul Boniol, Themis Palpanas, Ruey S. Tsay, Aaron Elmore, and Michael J.
557 Franklin. Volume under the surface: a new accuracy evaluation measure for time-series
558 anomaly detection. *Proc. VLDB Endow.*, 15(11):2774–2787, July 2022. ISSN 2150-8097. doi:
559 10.14778/3551793.3551830. URL <https://doi.org/10.14778/3551793.3551830>.
- 561 Hyunjong Park, Jongyoun Noh, and Bumsub Ham. Learning memory-guided normality for anomaly
562 detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-*
563 *nition*, pp. 14372–14381, 2020.
- 564 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
565 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text
566 transformer. *J. Mach. Learn. Res.*, 21(1), January 2020. ISSN 1532-4435.
- 567 Bernhard Schölkopf, Robert Williamson, Alex Smola, John Shawe-Taylor, and John Platt. Support
568 vector method for novelty detection. In *Proceedings of the 13th International Conference on*
569 *Neural Information Processing Systems, NIPS’99*, pp. 582–588, Cambridge, MA, USA, 1999.
570 MIT Press.
- 571 Ke-Yuan Shen. Learn hybrid prototypes for multivariate time series anomaly detection. In
572 *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=8TBGdH3t6a>.
- 573 Qichao Shentu, Beibu Li, Kai Zhao, Yang Shu, Zhongwen Rao, Lujia Pan, Bin Yang, and Chenjuan
574 Guo. Towards a general time series anomaly detector with adaptive bottlenecks and dual adver-
575 sarial decoders. In *The Thirteenth International Conference on Learning Representations*, 2025.
576 URL <https://openreview.net/forum?id=aKcd7ImG5e>.
- 577 Xiaoming Shi, Shiyu Wang, Yuqi Nie, Dianqi Li, Zhou Ye, Qingsong Wen, and Ming Jin. Time-
578 mae: Billion-scale time series foundation models with mixture of experts, 2024. URL <https://arxiv.org/abs/2409.16040>.
- 579 Junho Song, Keonwoo Kim, Jeonglyul Oh, and Sungzoon Cho. Memto: Memory-
580 guided transformer for multivariate time series anomaly detection. In A. Oh, T. Nau-
581 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural*
582 *Information Processing Systems*, volume 36, pp. 57947–57963. Curran Associates, Inc.,
583 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/b4c898eb1fb556b8d871f9e9ead92256-Paper-Conference.pdf)
584 [file/b4c898eb1fb556b8d871f9e9ead92256-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/b4c898eb1fb556b8d871f9e9ead92256-Paper-Conference.pdf).
- 585 Ya Su, Youjian Zhao, Chenhao Niu, Rong Liu, Wei Sun, and Dan Pei. Robust anomaly detec-
586 tion for multivariate time series through stochastic recurrent neural network. In *Proceedings*
587 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining,*
588 *KDD ’19*, pp. 2828–2837, New York, NY, USA, 2019. Association for Computing Machin-
589 ery. ISBN 9781450362016. doi: 10.1145/3292500.3330672. URL [https://doi.org/10.](https://doi.org/10.1145/3292500.3330672)
590 [1145/3292500.3330672](https://doi.org/10.1145/3292500.3330672).

- 594 Shreshth Tuli, Giuliano Casale, and Nicholas R. Jennings. Tranad: deep transformer networks
595 for anomaly detection in multivariate time series data. *Proc. VLDB Endow.*, 15(6):1201–1214,
596 February 2022. ISSN 2150-8097. doi: 10.14778/3514061.3514067. URL <https://doi.org/10.14778/3514061.3514067>.
- 598 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
599 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of the 31st Inter-*
600 *national Conference on Neural Information Processing Systems*, NIPS’17, pp. 6000–6010, Red
601 Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- 603 Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Aug-
604 menting language models with long-term memory. In *Proceedings of the 37th International Con-*
605 *ference on Neural Information Processing Systems*, NIPS ’23, Red Hook, NY, USA, 2023. Curran
606 Associates Inc.
- 607 Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo.
608 Unified training of universal time series forecasting transformers. In *Forty-first International*
609 *Conference on Machine Learning*, 2024.
- 610 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:
611 Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International*
612 *Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=ju_Uqw384Oq)
613 [id=ju_Uqw384Oq](https://openreview.net/forum?id=ju_Uqw384Oq).
- 615 Renjie Wu and Eamonn J. Keogh. Current time series anomaly detection benchmarks are flawed
616 and are creating the illusion of progress (extended abstract). In *2022 IEEE 38th International*
617 *Conference on Data Engineering (ICDE)*, pp. 1479–1480, 2022. doi: 10.1109/ICDE53745.2022.
618 00116.
- 619 Haowen Xu, Wenxiao Chen, Nengwen Zhao, Zeyan Li, Jiahao Bu, Zhihan Li, Ying Liu, Youjian
620 Zhao, Dan Pei, Yang Feng, et al. Unsupervised anomaly detection via variational auto-encoder
621 for seasonal kpis in web applications. In *Proceedings of the 2018 World Wide Web Conference on*
622 *World Wide Web*, pp. 187–196. International World Wide Web Conferences Steering Committee,
623 2018.
- 624 Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer: Time series
625 anomaly detection with association discrepancy. In *International Conference on Learning Repre-*
626 *sentations*, 2022. URL https://openreview.net/forum?id=LzQQ89U1qm_.
- 628 Hang Zhou, Junqing Yu, and Wei Yang. Dual memory units with uncertainty regulation for
629 weakly supervised video anomaly detection. In *Proceedings of the Thirty-Seventh AAAI Con-*
630 *ference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of*
631 *Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intel-*
632 *ligence*, AAAI’23/IAAI’23/EAAI’23. AAAI Press, 2023. ISBN 978-1-57735-880-0. doi:
633 10.1609/aaai.v37i3.25489. URL <https://doi.org/10.1609/aaai.v37i3.25489>.
- 634 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
635 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *The Thirty-*
636 *Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Virtual Conference*, volume 35, pp.
637 11106–11115. AAAI Press, 2021.
- 638 Quan Zhou, Changhua Pei, Fei Sun, HanJing, Zhengwei Gao, Haiming Zhang, Gaogang Xie, Dan
639 Pei, and Jianhui li. KAN-AD: Time series anomaly detection with kolmogorov–arnold net-
640 works. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=LWQ4zu9SdQ>.
- 641
642
643
644
645
646
647

A TSB-AD-U BENCHMARK

In this section, we provide details of the TSB-AD-U benchmark (Liu & Paparrizos, 2024) used in our experiments. The benchmark comprises 23 univariate time series datasets with a total of 870 time series, collected from diverse domains. Table 5 summarizes the dataset composition, including the number of series, average sequence length, anomaly counts, and anomaly ratios. The benchmark is divided into two partitions: the Eval set, which contains 350 time series and is designated for evaluation, and the Tuning set, which contains 48 time series and is used for hyperparameter optimization. Table 6 summarizes the statistics of the Eval and Tuning sets.

Table 5: Summary characteristics of 23 datasets included in TSB-AD-U. '-' in the 2nd column indicates this dataset is transformed from the multivariate dataset. The 'Category' column indicates whether the datasets feature point anomalies (P) or sequence anomalies (Seq).

Name	# TS Collected	# TS Curated	Avg Dim	Avg TS Len	Avg # Anomaly	Avg Anomaly Len	Anomaly Ratio	Category
UCR	250	228	1	67818.7	1	198.9	0.6%	P&Seq
NAB	58	28	1	5099.7	1.6	370.1	10.6%	Seq
YAHOO	367	259	1	1560.2	5.5	2.5	0.6%	P&Seq
IOPS	58	17	1	72792.3	25.6	48.7	1.3%	Seq
MGAB	10	9	1	97777.8	9.7	20.0	0.2%	Seq
WSD	210	111	1	17444.5	5.1	25.4	0.6%	Seq
SED	6	3	1	23332.3	14.7	64.0	4.1%	Seq
TODS	15	15	1	5000.0	97.3	18.7	6.3%	P&Seq
NEK	48	9	1	1073.0	2.9	51.1	8.0%	P&Seq
Stock	90	20	1	15000.0	1246.9	1.1	9.4%	P&Seq
Power	1	1	1	35040.0	4	750	8.5%	Seq
Daphnet (U)	-	1	1	38774.0	6	384.3	5.9%	Seq
CATSV2 (U)	-	1	1	300000.0	19.0	778.9	4.9%	Seq
SWaT (U)	-	1	1	419919.0	27.0	1876.0	12.1%	Seq
LTDB (U)	-	9	1	99700.0	127.5	144.5	18.6%	Seq
TAO (U)	-	3	1	10000.0	838.7	1.1	9.4%	P&Seq
Exathlon (U)	-	32	1	44075.8	3.1	1577.3	11.0%	Seq
MITDB (U)	-	8	1	631250.0	68.7	451.9	4.2%	Seq
MSL (U)	-	9	1	3492.0	1.3	130.0	5.8%	Seq
SMAP (U)	-	19	1	7700.2	1.2	210.1	2.8%	Seq
SMD (U)	-	38	1	24207.7	2.4	173.7	2.0%	Seq
SVDB (U)	-	20	1	171380.0	36.4	292.5	3.6%	Seq
OPP (U)	-	29	1	16544.8	1.4	653.4	6.4%	Seq

Table 6: Statistics of the TSB-AD-U benchmark splits.

	Split	# TS	Avg Length	Avg Anomaly Length	Avg # Anomalies	Anomaly Ratio
	All	870	38814.1	179.5	39.7	2.4%
TSB-AD-U	Eval	350	51886.7	321.3	46.6	4.5%
	Tuning	48	47143.3	185.9	82.6	3.5%

B IMPLEMENTATION DETAILS

All experiments in this paper are conducted on the Eval set of the TSB-AD-U benchmark, which contains 350 univariate time series. We describe below how the domain is defined, the training settings we consider, and the training configurations we adopt in our experiments.

B.1 DOMAIN PARTITION

We assign a domain label to each time series in the Eval set. A domain label is defined as a tuple (Dataset, Sub-domain), where the Dataset indicates the source collection and the Sub-domain specifies the actual domain of the time series. Time series characteristics may vary across domains within the same dataset, and they may also differ across datasets even when they share the same domain. To capture this heterogeneity, domain labels serve as a heuristic for balanced memory initialization.

Each file in the benchmark follows a structured naming convention.

For example:

```
001_NAB_id_1_Facility_tr_1007_1st_2014.csv
```

In this case:

- **Dataset:=** NAB
- **Sub-domain:=** Facility
- **Domain label:=** (NAB, Facility)

In total, the 350 univariate time series in the Eval set are grouped into 32 domains according to this labeling. These domain labels are used only during the memory initialization stage.

B.2 TRAINING SETTINGS

We distinguish between two training settings used in our experiments:

- **One-model-per-dataset:** In this conventional scheme, a separate model is trained for each dataset in the Eval set. Since the Eval set contains 350 time series, this setting requires 350 independently trained models.
- **Multi-domain training:** In contrast, our proposed approach jointly optimizes a single model across all datasets in the Eval set.

Unless otherwise specified, all baseline models follow the one-model-per-dataset setting, while MOMEMTO adopts the multi-domain training setting.

B.3 BASELINES AND HYPERPARAMETER SETTINGS

All baseline models are implemented by training on the train set and predicting the test set as defined in the benchmark. Their hyperparameters are primarily tuned on the tuning set following the benchmark protocol. Most models segment each time series into fixed-length windows using a sliding window with stride 1. A subset of models, such as MOMENT, Anomaly Transformer, and MEMTO, use non-overlapping windows to respect the original implementations in their respective papers.

Table 7: Hyperparameter settings for baseline models.

Model	Hyperparameter
LOF	n_neighbors = 50, metric = minkowski
OCSVM	periodicity = 2, kernel = rbf
Isolation Forest	n_estimators = 200
MEMTO	win_size = 100, lr = 0.00005, λ = 0.01, n_memory = 10
Donut	win_size = 60, lr = 0.0001
TimesNet	win_size = 32, lr = 0.0001
AnomalyTransformer	win_size = 100, lr = 0.001
TranAD	win_size = 10, lr = 0.0001
OmniAnomaly	win_size = 5, lr = 0.002
LSTM-AD	win_size = 100, lr = 0.0008
Chronos	win_size = 100
TimesFM	win_size = 96
MOMENT	win_size = 512, lr = 0.0001

For MOMEMTO, we use MOMENT-large (Goswami et al., 2024) as the backbone model, and most of the hyperparameters remain consistent with the original MOMENT configuration. Only the parameters related to the patch-based memory module are newly introduced. Table 8 summarizes the hyperparameter settings and environment used in our experiments.

Table 8: Hyperparameter settings for MOMEMTO and experimental environment.

Category	Setting
Window size	512
Patch length (L)	8
Patch stride length	8
Number of patches (N)	64
d_{model}	1024
Number of referenced items (K)	3
Temperature parameter (τ)	0.3
Learning rate	0.0001
Epochs	2
Optimizer	Adam
Loss function	Mean Squared Error (MSE)
GPU	NVIDIA GeForce RTX 4090 (24GB)
Framework	PyTorch 2.7.0

B.4 ALGORITHMS

Algorithms 1 and 2 illustrate the overall mechanism of our model. They present the matrix operation version of the forward process when a single input subsequence is given.

Algorithm 1 Proposed Method MOMEMTO

Input $x \in \mathbb{R}^{N \times L}$: input time series, $[\text{mask}] = \{0, 1\}^N$, $\mathcal{M} \in \mathbb{R}^{M \times N \times d_{model}}$: domain-specific initialized memory items

Training params f_e : encoder, f_d : decoder, f_m : patch-based memory module

- 1: $q \leftarrow f_e(x, [\text{mask}])$ $\triangleright q \in \mathbb{R}^{P \times d_{model}}$
 - 2: $\hat{q} \leftarrow \text{L2Norm}(q, \text{dim} = 1)$
 - 3: $\tilde{q} \leftarrow f_m(\hat{q}, \mathcal{M}, [\text{mask}])$
 - 4: $q \leftarrow \text{concat}([q, \tilde{q}], \text{dim} = 1)$ $\triangleright q \in \mathbb{R}^{P \times (2 \cdot d_{model})}$
 - 5: $\hat{x} \leftarrow f_d(q)$
 - 6: **return** \hat{x} \triangleright Reconstructed time series
-

Algorithm 2 Proposed Patch-based memory module

Input $\hat{q} \in \mathbb{R}^{P \times d_{model}}$: normalized queries, $[\text{mask}] = \{0, 1\}^N$, $\mathcal{M} \in \mathbb{R}^{M \times N \times d_{model}}$: memory items

Training params $U_\psi, W_\psi \in \mathbb{R}^{d_{model} \times d_{model}}$: linear projection matrices

- 1: $\mathcal{M} \in \mathbb{R}^{M \times P \times d_{model}} \xleftarrow{[\text{mask}]} \mathcal{M} \in \mathbb{R}^{M \times N \times d_{model}}$ ▷ Memory alignment
- 2: Reshape: $\hat{q} \in \mathbb{R}^{1 \times (P \cdot d_{model})}, \mathcal{M} \in \mathbb{R}^{M \times (P \cdot d_{model})}$
- 3: $\lambda \leftarrow \text{softmax}(\hat{q} \mathcal{M}^T)$
- 4: $(\lambda_{1st} : m_{1st}), (\lambda_{2nd} : m_{2nd}), \dots, (\lambda_{K-th} : m_{K-th}) \leftarrow \text{TopK}(\lambda)$ ▷ Find K most similar items
- 5: Reshape: $\hat{q} \in \mathbb{R}^{P \times d_{model}}, \mathcal{M} \in \mathbb{R}^{M \times P \times d_{model}}$
- 6: $\tilde{q} \leftarrow \mathbf{0}_{P \times d_{model}}$
- 7: **for** $i \leftarrow 1$ **to** K **do** ▷ Memory update
- 8: $v \leftarrow \text{softmax}(m_{i-th} q^T)$ ▷ $m_{i-th} \in \mathbb{R}^{P \times d_{model}}$: the i -th similar item
- 9: $\psi \leftarrow \text{sigmoid}(m_{i-th} U_\psi + v q W_\psi)$ ▷ $\psi \in \mathbb{R}^{P \times d_{model}}$: update rate
- 10: $\tilde{m}_{i-th} \leftarrow (1 - \psi) \odot m_{i-th} + \psi \odot v q$
- 11: $w \leftarrow \text{softmax}(q(\tilde{m}_{i-th})^T)$ ▷ Query update
- 12: $\tilde{q}_{i-th} \leftarrow w \tilde{m}_{i-th}$
- 13: $\tilde{q} \leftarrow \tilde{q} + (\lambda_{i-th} \odot \tilde{q}_{i-th})$
- 14: **end for**
- 15: **return** \tilde{q} ▷ $\tilde{q} \in \mathbb{R}^{P \times d_{model}}$

C ADDITIONAL EXPERIMENTAL RESULTS

C.1 ADDITIONAL RESULTS FOR BENCHMARK ACCURACY EVALUATION

Table 9 summarizes more detailed VUS-PR results by domain, aggregated from 350 time series across 32 domains.

Table 9: VUS-PR results by domain.

Domain	LOF	OCSVM	iForest	MEMTO	Donut	TimesNet	A.T.	TranAD	OmniAnomaly	LSTM-AD	Chronos	TimesFM	MOMENT	MOMEMTO	MOMEMTO _{mid}
(CATSV2, Sensor)	6.38	<u>26.08</u>	7.94	5.31	8.09	7.23	9.34	7.51	10.92	29.16	9.60	24.50	19.81	9.79	25.45
(Daphnet, HumanActivity)	13.06	6.09	<u>36.43</u>	6.40	9.69	40.19	12.85	21.05	27.42	13.27	31.53	35.46	34.95	20.55	11.42
(Exathlon, Facility)	19.85	29.05	66.72	13.46	51.21	53.96	64.43	78.09	<u>86.74</u>	73.28	42.24	50.19	52.79	79.88	90.48
(IOPS, WebService)	12.05	6.91	27.71	4.57	10.07	20.95	10.94	19.98	23.55	22.81	19.16	18.72	21.28	<u>25.80</u>	22.05
(LTDB, Medical)	26.23	32.71	33.79	24.72	30.85	31.69	39.20	31.43	31.62	34.66	26.11	27.60	29.37	32.14	<u>36.97</u>
(MGAB, Synthetic)	0.44	0.73	0.45	0.43	0.42	0.43	0.42	0.49	0.49	4.47	0.45	0.44	0.51	0.44	<u>1.25</u>
(MITDB, Medical)	6.12	<u>14.09</u>	9.79	6.17	9.81	7.89	8.26	9.34	9.52	12.15	5.88	6.46	7.73	9.74	17.13
(MSL, Sensor)	14.70	28.26	28.83	11.19	22.18	32.25	27.33	25.57	28.08	24.88	22.41	31.05	<u>32.19</u>	25.32	29.00
(NAB, Environment)	25.50	7.74	35.90	10.54	11.33	12.38	33.82	32.45	31.96	20.20	10.92	10.92	13.21	<u>37.32</u>	38.20
(NAB, Facility)	15.95	20.87	23.89	12.06	16.79	21.59	<u>26.51</u>	23.79	25.93	18.93	18.73	20.17	25.20	24.18	30.85
(NAB, Synthetic)	18.23	52.23	27.27	15.67	18.72	23.30	21.14	23.49	24.55	23.40	20.05	18.96	23.62	20.96	<u>31.77</u>
(NAB, Traffic)	14.06	20.66	16.66	11.46	16.67	16.04	14.40	17.82	<u>19.90</u>	11.68	15.43	15.82	16.91	16.76	17.71
(NAB, WebService)	16.34	17.25	16.53	16.03	17.76	27.59	20.13	16.05	19.06	16.15	19.79	16.28	19.21	15.98	<u>24.31</u>
(NEK, WebService)	37.70	25.66	59.30	8.98	48.00	48.48	30.15	70.90	85.24	<u>74.16</u>	33.84	34.46	42.57	73.24	66.39
(OPPORTUNITY, HumanActivity)	14.25	11.49	43.39	16.58	28.21	5.25	67.89	65.17	69.39	58.43	5.36	4.80	6.06	<u>74.48</u>	77.91
(Power, Facility)	9.05	<u>16.36</u>	8.20	9.05	8.56	7.55	9.27	14.46	14.83	6.73	8.18	7.89	8.22	10.53	23.59
(SED, Medical)	<u>10.98</u>	5.98	35.84	7.66	10.53	4.90	6.31	4.65	5.23	5.72	6.22	5.28	5.85	5.93	7.90
(SMAP, Sensor)	15.12	51.18	24.66	5.96	30.43	<u>42.65</u>	15.91	24.29	25.36	25.02	17.66	33.02	35.09	21.14	26.38
(SMD, Facility)	13.01	8.15	34.18	3.43	28.62	53.85	14.28	34.42	38.02	<u>50.31</u>	34.09	39.99	38.72	34.45	33.31
(SVDB, Medical)	5.26	20.47	9.03	4.75	7.80	8.59	11.96	7.99	8.86	13.38	5.51	6.31	7.14	8.42	<u>15.91</u>
(SWaT, Sensor)	12.12	9.35	49.69	12.71	46.75	9.60	67.15	51.45	46.08	<u>66.78</u>	20.64	17.09	7.87	46.68	47.70
(Stock, Finance)	74.99	73.15	<u>92.30</u>	76.31	78.00	78.52	74.86	82.21	92.54	85.41	98.39	98.55	99.52	98.49	88.97
(TAO, Environment)	90.94	92.49	98.52	90.95	90.68	90.49	91.26	94.19	97.79	99.87	99.85	99.35	99.99	99.92	97.15
(TODS, Synthetic)	48.55	65.07	51.81	45.38	48.97	58.01	50.07	47.78	46.41	47.23	<u>73.24</u>	75.15	67.15	49.85	61.41
(UCR, Environment)	2.26	20.86	0.86	0.87	0.97	1.64	1.31	0.79	0.92	1.02	14.27	<u>14.81</u>	3.31	1.06	4.33
(UCR, Facility)	1.01	3.33	1.06	0.81	0.70	1.14	<u>2.73</u>	1.18	1.29	1.17	0.69	0.89	1.07	1.13	2.71
(UCR, HumanActivity)	3.92	1.39	3.44	0.85	1.57	3.51	1.84	2.48	3.38	2.47	5.62	5.72	<u>5.97</u>	2.42	10.93
(UCR, Medical)	2.25	16.16	2.64	1.64	1.97	2.79	2.15	2.38	2.46	2.32	<u>7.44</u>	7.43	4.09	2.31	5.52
(UCR, Sensor)	2.29	53.98	4.14	1.45	1.58	1.91	<u>5.54</u>	2.82	2.97	5.19	3.66	3.60	5.13	2.49	4.64
(WSD, WebService)	8.61	2.81	14.41	4.29	5.19	19.16	4.78	12.03	15.47	13.38	17.67	20.62	<u>20.27</u>	15.79	17.23
(YAHOO, Synthetic)	62.91	27.06	56.39	18.36	8.47	40.14	3.13	36.24	14.70	63.67	90.57	<u>93.75</u>	97.34	63.57	82.97
(YAHOO, WebService)	17.23	19.71	33.72	12.39	7.48	13.09	11.63	24.28	31.59	30.89	<u>70.28</u>	72.16	63.98	35.62	44.04

864 C.2 ZERO-SHOT ANOMALY DETECTION
865

866 We conduct a leave-one-out evaluation across 32 domains to assess the zero-shot performance of
867 MOMEMTO. In each iteration, one domain is held out as the target domain, while the model is
868 trained on the remaining 31 domains. The trained model is then directly applied to the held-out
869 domain without any fine-tuning, under a zero-shot setting. Performance metrics are computed indi-
870 vidually for each target domain and subsequently aggregated over all 32 domains. To highlight the
871 effectiveness of our memory-augmented design, we also compare MOMEMTO against its backbone
872 MOMENT.

874 Table 10: Zero-shot anomaly detection results under leave-one-out evaluation across 32 domains.
875 MOMEMTO is compared with MOMENT.

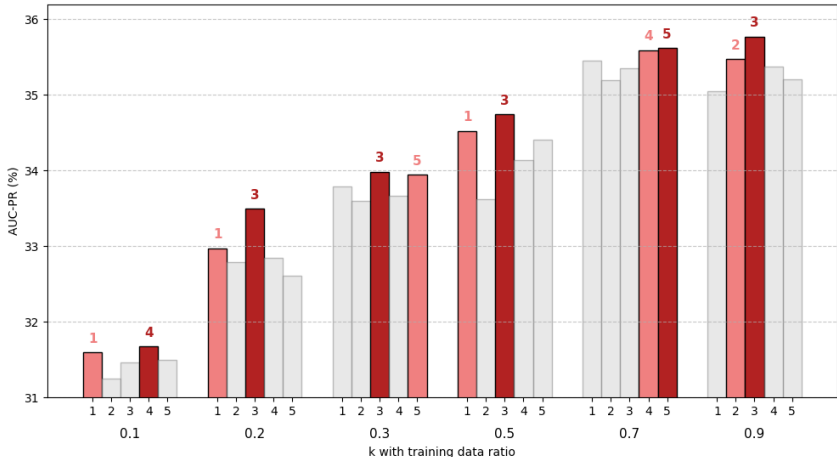
877

	AUC-PR	AUC-ROC	VUS-PR	VUS-ROC
MOMENT _{md}	29.68	69.61	30.64	76.57
MOMEMTO _{md}	34.51	72.17	35.76	78.32

881

882
883
884 C.3 NUMBER OF REFERENCED MEMORY ITEMS

885 We analyze how the number of referenced memory items (K) and the training data ratio jointly
886 affect the performance of MOMEMTO. As shown in Figure 5, when $K = 1$, the model leverages
887 only the most similar memory item, providing stable yet limited improvements. For $K \geq 2$, the
888 performance is not always superior to that with $K = 1$, but the best results are consistently achieved
889 when the model references multiple items.
890



907 Figure 5: Effect of the number of referenced memory items (K) and the training data ratio on AUC-
908 PR.

909
910
911 C.4 NUMBER OF MEMORY ITEMS

912 To evaluate how the number of memory items affects MOMEMTO, we conduct a sensitivity analysis
913 by varying the number of memory items from 64 down to 1. In this experiment, each memory item
914 is initialized using K-means centroid values.
915

916 Figure 6 visualizes the row-normalized scaled domain-to-memory similarities at the beginning and
917 end of training. Across all configurations, the model consistently relies on a small subset of memory
items, regardless of the total number of memory items.

Table 11: Sensitivity analysis of the number of memory items on the main benchmark.

	64	32	16	8	4	2	1
AUC-PR	36.11	36.05	36.26	36.07	35.97	35.69	35.81
AUC-ROC	74.59	74.22	74.18	74.19	74.30	73.25	73.50

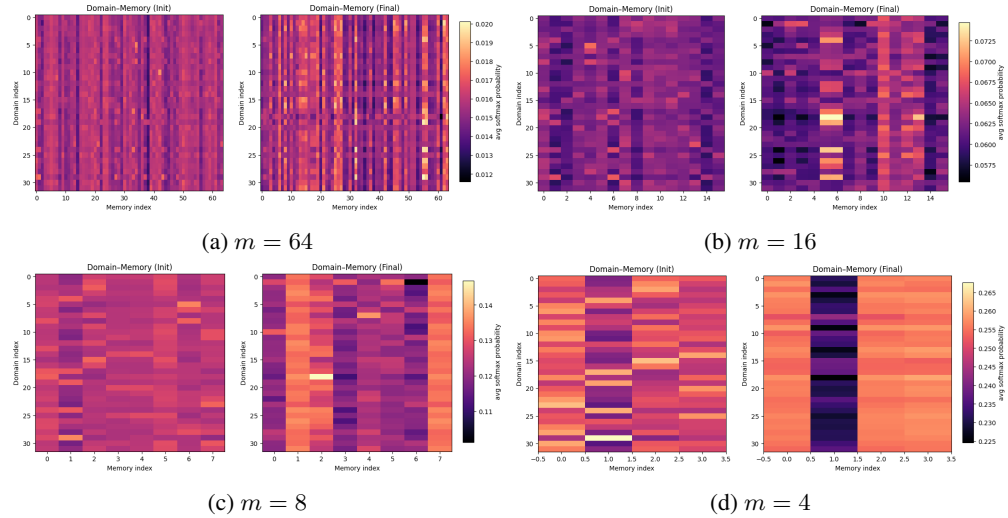


Figure 6: Scaled domain-to-memory similarity maps at the beginning (left) and end of training (right) for different numbers of memory items. Each heatmap is row-normalized.

Table 12: Few-shot performance sensitivity to the number of memory items (AUC-PR)

	0.1	0.2	0.3	0.5	0.7	0.9
32	31.02	32.84	33.40	35.04	35.23	35.63
16	30.98	32.83	33.46	35.01	35.36	35.55
8	30.95	32.75	33.48	34.81	35.21	35.16
4	30.93	32.42	33.52	34.69	35.42	35.32
2	30.85	32.80	33.52	35.08	35.67	35.65
1	30.93	32.32	33.50	34.84	35.60	35.51

These results indicate that MOMEMTO exhibits stable performance with respect to the number of memory items, consistent with the sensitivity analysis reported in MEMTO.

C.5 ANOMALY SCORE DISTRIBUTIONS

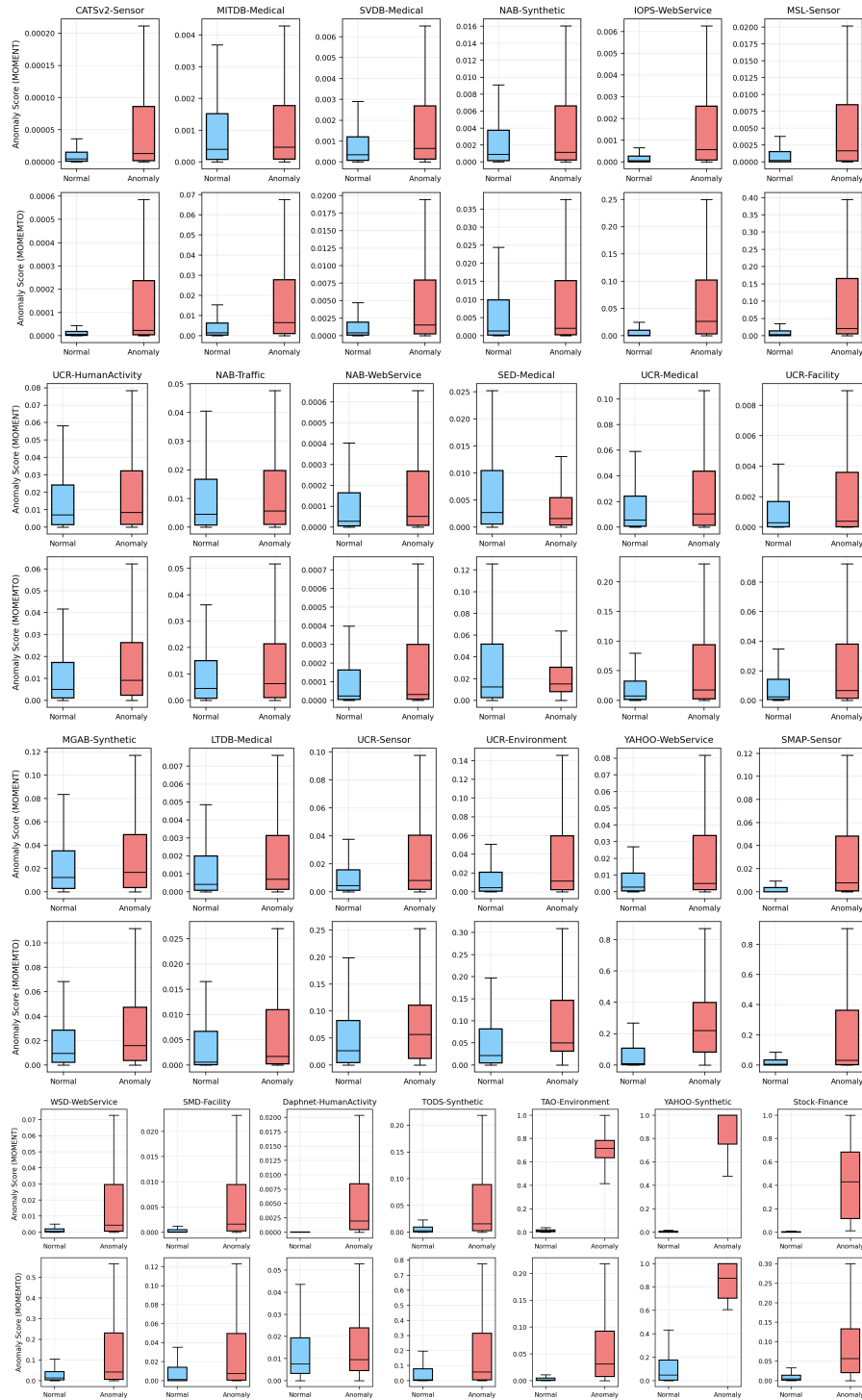


Figure 7: Comparison of anomaly score distributions between normal (blue) and anomalous (red) samples across a subset of domains. The upper row shows the results with the backbone model MOMENT, while the lower row corresponds to our proposed model.

1026 D THE USE OF LARGE LANGUAGE MODELS

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

A large language model (LLM) is used exclusively for linguistic polishing of the manuscript, including grammar correction and style refinement. The LLM does not contribute to the conception of ideas, experimental design, implementation, or analysis.