

Visual Reasoning over Time Series via Multi-Agent Systems

Weilin Ruan¹ Yuxuan Liang^{*,1}

¹The Hong Kong University of Science and Technology (Guangzhou), China

rwlinno@gmail.com, yuxliang@outlook.com

Abstract

Time series analysis underpins many real-world applications, yet existing time-series-specific methods and pretrained large-model-based approaches remain limited in integrating intuitive visual reasoning and generalizing across tasks with adaptive tool usage. To address these limitations, we propose **MAS4TS**, a tool-driven multi-agent system for general time series tasks, built upon an **Analyzer-Reasoner-Executor** paradigm that integrates agent communication, visual reasoning, and latent reconstruction within a unified framework. MAS4TS first performs visual reasoning over time series plots with structured priors using a Vision-Language Model to extract temporal structures, and subsequently reconstructs predictive trajectories in latent space. Three specialized agents coordinate via shared memory and gated communication, while a router selects task-specific tool chains for execution. Extensive experiments on multiple benchmarks demonstrate that MAS4TS achieves state-of-the-art performance across a wide range of time series tasks, while exhibiting strong generalization and efficient inference.

1 Introduction

Time series analysis is fundamental to a wide range of real-world applications, including energy systems (Deb et al., 2017), transportation (Zheng and Huang, 2020), and other critical domains (Idrees et al., 2019; Karevan and Suykens, 2020; Schneider and Dickinson, 1974). Core tasks such as forecasting, classification, imputation, and anomaly detection (Lim and Zohren, 2021) play a crucial role in supporting data-driven decision-making in these scenarios.

Recent advances in pretrained Large Language Models (LLMs) have demonstrated strong capabilities in pattern understanding and general reasoning (Brown et al., 2020;

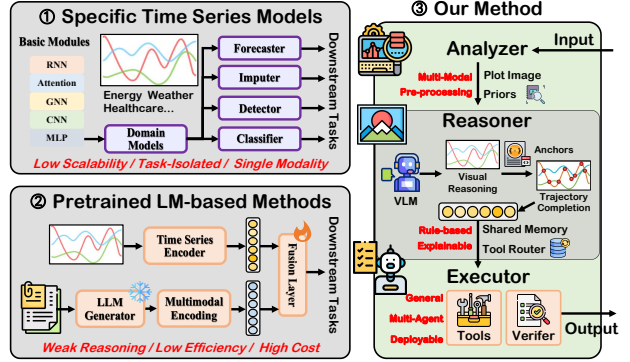


Figure 1. The motivation of MAS4TS. Existing approaches (left) struggle to support intuitive reasoning and task-adaptive execution, while MAS4TS (right) organizes time series analysis into a unified paradigm that collaboratively analyzes, reasons, and executes.

Touvron et al., 2023). Motivated by this progress, a growing body of work explores leveraging pretrained models for time series analysis (Jin et al., 2024a; Liu et al., 2024a; Zhou et al., 2023), aiming to improve generalization and contextual modeling beyond traditional neural architectures (Wu et al., 2023b; Zhou et al., 2021; Nie et al., 2023a). In parallel, agentic methods and multi-agent systems (MAS) (Ravuru et al., 2024; Gu et al., 2025; Chang et al., 2025a; Wu et al., 2024; Hong et al., 2023; Qian et al., 2023; Huang et al., 2025) are emerging as a practical paradigm for tool orchestration, task decomposition, and coordinated execution.

Despite these advances, the application of multi-agent paradigms to general time series analysis remains largely unexplored. Most existing approaches primarily rely on implicit historical pattern matching in numerical or embedding spaces, lacking explicit and interpretable reasoning mechanisms that capture latent temporal structures. Even recent vision-enhanced time series methods (Chen et al., 2024; Zhong et al., 2025; Ruan et al., 2025) often treat visual representations as auxiliary features, rather than leveraging the ability of modern Vision-Language Models (VLMs) to perform *direct reasoning on time series plots* that naturally encode trends, regime shifts, and structural transitions. Moreover, current methods are typically designed for isolated tasks under fixed inference pipelines, with limited support for adaptive tool orchestration and coordinated exe-

cution across heterogeneous time series workloads. These limitations can be summarized by two conceptual gaps:

- **From Representation Learning to Structural Reasoning.** Existing methods operate in numerical or latent embedding spaces where morphological patterns remain implicit and entangled. Recent vision-enhanced methods (e.g., TimeVLM, VisionTS) convert time series into vision primarily for *representation learning*, but fail to exploit visual reasoning capabilities to extract explicit structural knowledge from plots like human analysts naturally use for temporal understanding. *Can we shift from learning representations to performing intuitive structural reasoning over time series visualizations?*
- **From Task-Specific Models to Tool-Driven Execution.** Current paradigms treat forecasting, classification, imputation, and anomaly detection as separate problems, typically requiring task-specific heads and individual procedures even when sharing a common backbone. Unlike multi-agent systems that leverage tool orchestration for adaptive problem-solving, time series methods operate under the *fixed model then single-task finetuning* paradigms. It prevents pretrained models from functioning as general-purpose, task-adaptive time series solvers and necessitates multiplicative scaling costs. *Can we unify diverse time series tasks through shared structural reasoning while enabling adaptive tool selection for task-specific execution?*

General time series tasks fundamentally require an understanding of morphological semantics that are explicit in visual plots yet implicit in numerical embeddings. This reveals that the bottleneck of time series intelligence lies in structural reasoning rather than representation learning, suggesting a paradigm shift: unify tasks through visual reasoning while enabling adaptive tool-driven execution.

Based on this insight, we propose **MAS4TS**, the first tool-driven Multi-Agent System for general Time Series analysis, built upon an **Analyzer–Reasoner–Executor** paradigm. MAS4TS decomposes time series analysis into complementary roles: the *Analyzer* performs statistical analysis and data preprocessing; the *Reasoner* integrates visual and numerical reasoning to extract high-level temporal priors and derive latent trajectory constraints; and the *Executor* invokes appropriate tools and models to generate task-specific outputs. Specifically, MAS4TS performs visual reasoning directly on time series plots using VLMs to obtain structured anchors, which subsequently guide trajectory reconstruction in latent space. Through shared memory and inter-agent communication, specialized agents collaborate to solve diverse time series tasks within a unified framework. Our key contributions are summarized as follows:

- We propose **MAS4TS**, the first tool-driven multi-agent system for general time series analysis, establishing an

Analyzer–Reasoner–Executor paradigm for task-adaptive time series modeling.

- We introduce an intuitive reasoning approach that utilizes VLMs to extract visual anchors from plots and priors then transforms them into numerical constraints for latent trajectory reconstruction.
- We conduct extensive experiments and demonstrate the superior performance of MAS4TS among forecasting, classification, imputation, and anomaly detection tasks.

2 Related Work

2.1 Time Series Analysis

Classical time series models (Box et al., 2015; Cleveland et al., 1990) remain effective for simple or linear patterns. Deep learning has become the dominant paradigm in modern time series analysis, enabled by increasingly specialized architectures. Recurrent models (Medsker et al., 2001; Hochreiter and Schmidhuber, 1997b; Cho et al., 2014) capture short-term temporal dependencies but struggle with gradient decay over long sequences, while convolutional models (Bai et al., 2018; Van Den Oord et al., 2016; Wu et al., 2023b) efficiently extract local and multi-scale patterns, offering strong computational advantages. Transformer-based models (Zhou et al., 2021; Wu et al., 2021; Liu et al., 2022b; Zhou et al., 2022; Nie et al., 2023a) have become the mainstream due to their flexible receptive fields and modeling power, with innovations such as auto-correlation attention, frequency-domain decomposition, and patch-level representations. Linear architectures (Zeng et al., 2023; Xu et al., 2023; Zhou et al., 2025) have shown surprising competitiveness on large forecasting benchmarks, renewing interest in inductive bias and architectural simplicity.

2.2 Pretrained Large Models for Time Series

Inspired by the success of LLMs, recent work explores foundation models for time series (Liang et al., 2024; Jin et al., 2024b). **LLM-based approaches** (Jin et al., 2024a; Liu et al., 2024a; Zhou et al., 2023; Chang et al., 2025b) align time series with text embeddings, leveraging pretrained language models for cross-domain generalization. However, these methods face inherent modality gaps and computational overhead. **Vision-enhanced approaches** convert time series to images and use pretrained vision models (Chen et al., 2024; Zhong et al., 2025; Shen et al., 2025; Ruan et al., 2025; Lyu et al., 2025), but primarily treat visual representations as auxiliary features for representation learning. *Our work differs by leveraging VLMs to perform direct visual reasoning on time series plots, extracting explicit structural patterns that guide downstream task execution.*

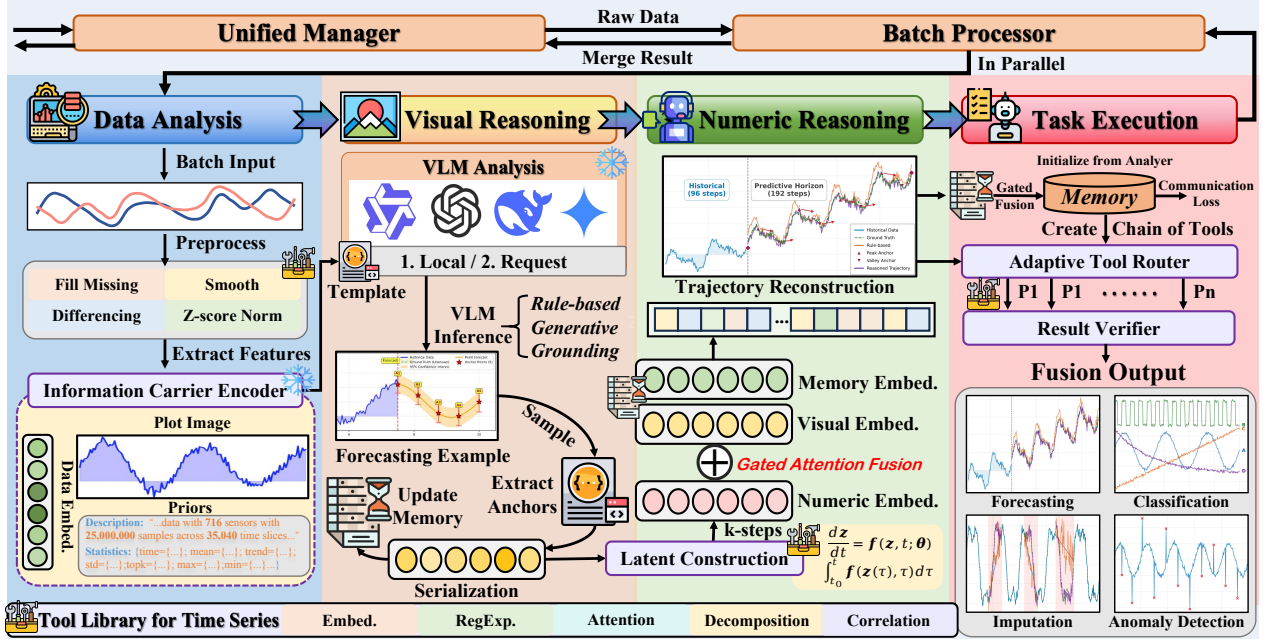


Figure 2. Overview of the MAS4TS framework following an Analyzer-Reasoner-Executor paradigm. The Reasoner agent performs two sequential stages: Visual Reasoning (VLM-based anchor extraction) and Numeric Reasoning (latent trajectory reconstruction). A unified tool library supports the entire workflow, with primary integration at the Executor.

2.3 Multi-Agent Systems

Multi-agent systems (MAS) have shown strong capabilities in complex reasoning, planning, and tool orchestration (Wu et al., 2024; Hong et al., 2023; Qian et al., 2023). By assigning agents specialized roles and enabling structured communication, MAS can decompose and tackle tasks beyond the capacity of a single model. Tool-augmented approaches (Schick et al., 2023; Qin et al., 2023) further extend agent capabilities through external function invocation. Recent work explores agentic approaches for specific time series tasks such as forecasting (Huang et al., 2025) and anomaly detection (Gu et al., 2025; Ravuru et al., 2024), but these methods address isolated problems rather than providing a unified analytical framework. *Our work fills this gap by introducing the tool-driven multi-agent framework that integrates visual reasoning for general time series analysis.*

3 Methodology

As shown in Figure 2, MAS4TS adopts a plan-and-execute multi-agent architecture (Huang et al., 2024), where a manager dynamically decomposes time series tasks into specialized stages executed by core agents. A processor then performs task-aware adaptive planning and enables parallel execution among agents to improve inference efficiency. For each parallel batch, the workflow consists of three agents: ① *Analyzer*, which extracts statistical features \mathbf{F}_{stat} and a task-aware plot image representation; ② *Reasoner*, which

leverages a VLM to generate anchors \mathcal{A} (visual reasoning stage) and reconstruct the sequence in a latent space (numeric reasoning stage); and ③ *Executor*, which produces and verifies final predictions $\hat{\mathbf{Y}}$ through adaptive tool usage.

3.1 Problem Formulation

Let $\mathbf{X} = \{x_1, \dots, x_L\} \in \mathbb{R}^{L \times D}$ denote an input multi-variate time series, where L is the length of the observed sequence and D is the number of variables (features). We consider four general tasks: ① *Forecasting* to predict future values $\mathbf{Y} = \{x_{L+1}, \dots, x_{L+H}\} \in \mathbb{R}^{H \times D}$ over a forecasting horizon of length H ; ② *Classification* to assign a categorical label $y \in \{1, \dots, C\}$ to the entire sequence, where C is the number of classes; ③ *Imputation* to estimate missing entries $\{\hat{x}_{i,j} : (i, j) \in \Omega\}$, with Ω denoting the index set of missing time-variable pairs; and ④ *Anomaly Detection* to identify anomalous timestamps $s \subset \{1, \dots, L\}$ within the observed window. MAS4TS provides a unified agent-driven framework that shares common perception, reasoning, and memory components across these tasks.

3.2 Multi-Agent Collaboration Framework

Time series analysis requires heterogeneous reasoning that is difficult to unify in a single model, motivating our design of specialized agents that collaborate through shared memory to enable selective knowledge transfer while preserving task-specific representations.

Shared Memory Mechanism. We define a set of agents $\{\text{Agent}^{(i)}\}_{i=1}^N$ for batch data processing, where in our instantiation $N=3$ corresponds to *Analyzer*, *Reasoner*, and *Executor*. Each agent has (i) an embedding function $\text{Embed}^{(i)}(\cdot)$, (ii) a local memory view $\mathcal{M}^{(i)}$, and (iii) a communication module $\text{Com}(\cdot)$. Cross-agent coordination is achieved through a shared memory matrix $\mathbf{M} \in \mathbb{R}^{L \times d_m}$, which integrates knowledge across agents via gated fusion. Given an input time series \mathbf{X} and current shared memory \mathbf{M} , each agent computes:

$$\mathbf{h}^{(i)}, \alpha_i = \text{Agent}^{(i)}(\text{Embed}^{(i)}(\mathbf{X}), \mathcal{M}^{(i)}), \quad (1)$$

$$\text{Com}(\mathbf{h}^{(i)} | \mathcal{M}^{(i)}) = \mathbf{h}^{(i)} + \text{LN}(\text{Enc}_\ell(\mathbf{h}^{(i)}, \mathcal{M}^{(i)})), \quad (2)$$

$$\mathbf{M} \leftarrow \mathbf{M} \oplus \alpha_i \cdot \text{Com}(\mathbf{h}^{(i)} | \mathcal{M}^{(i)}), \quad (3)$$

where $\mathbf{h}^{(i)}$ denotes the agent-specific hidden state, $\alpha_i \in [0, 1]$ is the confidence score controlling the update magnitude, and $\text{LN}(\cdot)$ denotes the layer normalization for $\ell = 1, \dots, L_{\text{enc}}$ encoder layers conditioned on the agent view. The gating mechanism \oplus enables iterative memory refinement while reducing the impact of noisy agent outputs.

Information Carrier Encoder. Agents communicate through structured information carriers $\{\mathcal{I}, \mathcal{E}, \mathcal{P}\}$ that are first generated by the **Analyzer**, detailed as follows:

- \mathcal{I} : Plot images rendered directly from the raw time series via a rule-based operator $\mathcal{I} = \mathcal{R}(\{(t, x_{t,d})\}_{t=1}^L)$;
- \mathcal{E} : Normalized embeddings from tool-integrated time series preprocessing and encoding;
- \mathcal{P} : Structured priors including statistical features, semantic descriptions, and shared memory from communication.

These carriers form a compact yet expressive interface for inter-agent communication and downstream reasoning.

3.3 Latent Visual Reasoning

The **Reasoner** aims to bridge intuitive visual pattern recognition and precise numerical inference for time series analysis. Rather than directly fitting historical trajectories, we leverage visual reasoning to identify sparse but informative structural cues and subsequently enforce them through continuous-time latent dynamics.

VLM-based Anchoring. Many salient temporal structures in time series, including trends, regime shifts, and boundary behaviors, are often easier to identify from visual representations than from raw 1-D sequences. Vision-Language Models (VLMs) perform high-level pattern recognition directly on plot images. Formally, given plot images

\mathcal{I} and structured priors \mathcal{P} , we query the VLM as:

$$\mathcal{A} = \text{VLM}(\text{Enc}_{\text{vis}}(\mathcal{I}), \text{Enc}_{\text{text}}(\text{Prompt}(\mathcal{P}))) \quad (4)$$

$$= \{(t_k, v_k, \tau_k)\}_{k=1}^K, \quad (5)$$

where the prompt (see Appendix D) specifies the required output schema. Each tuple consists of a critical time index $t_k \in \{1, \dots, L\}$, an estimated value $v_k \in \mathbb{R}$, and an indicator τ_k (default: $\{-1, 0, 1\}$) denoting local temporal behavior (e.g., falling/stable/rising). These sparse anchors encode high-level structural knowledge extracted through visual reasoning, serving as semantic constraints that guide the latent reconstruction.

Latent Construction. While visual anchors are sparse and discrete, downstream time series tasks require dense and temporally coherent representations. To bridge this gap, we introduce a *latent construction* module that reconstructs a continuous latent trajectory by treating anchors as conditioning signals that guide latent-space evolution. Formally, we define a continuous-time latent state $\mathbf{z}(t)$ and parameterize its dynamics by:

$$\frac{d\mathbf{z}(t)}{dt} = f_\phi(\mathbf{z}(t), t | \mathcal{A}, \mathcal{P}), \quad \mathbf{z}(t_0) = \mathbf{z}_0, \quad (6)$$

where f_ϕ is a learnable dynamics function conditioned on the anchor set \mathcal{A} and priors \mathcal{P} , and the initial state \mathbf{z}_0 is obtained from the embedding \mathcal{E} . We employ a finite-step latent reconstruction operator $\mathcal{S}(\cdot)$ to obtain latent states $\{\mathbf{z}_t\}_{t=1}^{L+H}$, covering the observed window and the horizon:

$$\mathbf{u}^{(k)} = \mathcal{S}^{(k)}(\dots \mathcal{S}^{(1)}(\mathbf{u}^{(0)}, \mathcal{A}, \mathcal{P})), \quad \mathbf{z}_t = [\mathbf{u}^{(k)}]_t, \quad (7)$$

where $\mathbf{u}^{(i)} \in \mathbb{R}^{(L+H) \times D}$ denotes the *sequence-level* latent representation after the i -th solver step, with $\mathbf{u}^{(0)} = \text{Proj}(\mathcal{E})$ initialized from the data embedding and fused by communication, and \mathbf{z}_t being the t -th row of $\mathbf{u}^{(k)}$.

In practice, we instantiate $\mathcal{S}(\cdot)$ using a finite-step ($k \leq 20$) explicit latent ODE solver, as anchor conditioning substantially restricts the feasible solution space and reduces the effective curvature of the latent trajectory. Intuitively, and consistent with prior work on continuous latent modeling (Dormand and Prince, 1980; Rubanova et al., 2019), anchor-conditioned latent dynamics act as a form of *soft boundary regularization*, biasing the reconstructed trajectory toward semantically meaningful temporal structures while preserving flexibility in local variations.

Latent Fusion with Data Embedding. The reconstructed latent trajectory provides high-level structural guidance but does not replace the original data representation. Instead, it is fused with the data embedding to form a unified latent representation using attention fusion:

$$\mathbf{Z} = \text{Attn}(\mathcal{E} \oplus \text{Proj}(\mathbf{u}^{(k)})), \quad (8)$$

where \oplus denotes a gated fusion operator. The resulting latent representation \mathbf{Z} is dense, differentiable, and semantically grounded, serving as a bridge between training-free visual reasoning and data-driven numerical inference within a unified latent framework.

3.4 Tool-Integrated Execution and Optimization

The **Executor** converts latent representations into task-specific outputs via tool-integrated execution and result verification, while optimizing only coordination-critical components and keeping the reasoning modules training-free. Instead of relying on a single monolithic predictor, MAS4TS dynamically composes a chain-of-tools that applies specialized time-series operators to refine intermediate results.

Tool Library for Time Series. At execution time, an adaptive router selects a tool chain from a predefined library $\mathcal{T} = \{\mathcal{T}_1, \dots, \mathcal{T}_M\}$, where each tool implements a specialized operation (e.g., trend extrapolation, decomposition, smoothing, statistical prediction, constraint projection, and residual correction). Tool routing adapts to the task type κ , structured priors \mathcal{P} , and the shared memory \mathbf{M} maintained in earlier stages. Formally, the executor produces a routing policy π_Θ over tools and constructs a tool sequence $\mathbf{m} = (m_1, \dots, m_J)$ accordingly. The final output is obtained by sequential composition:

$$\pi_\Theta(m \mid \kappa, \mathcal{P}, \mathbf{M}) = \text{Softmax}(g_\psi(\kappa, \mathcal{P}, \mathbf{M})), \quad (9)$$

$$\hat{\mathbf{Y}} = \mathcal{T}_{m_J} \circ \dots \circ \mathcal{T}_{m_1}(\mathbf{Z}), \quad (10)$$

where g_ψ is a gated controller network, and the tool sequence with length J is dynamically determined by policy π_Θ and connected by \circ operation. This design enhances raw predictions through plug-and-play expert operators without retraining a monolithic backbone for a specific task.

Result Verifier. To ensure numerical validity and task consistency, MAS4TS employs a lightweight verifier that performs: (i) format and shape correction, (ii) constraint validation, and (iii) task-specific normalization. If the candidate output fails verification, the executor either recomputes or falls back to a conservative basic tool for a valid prediction. Given a candidate output $\hat{\mathbf{Y}}$, the verifier produces the final prediction:

$$\mathbf{Y}^* = \text{Verifier}(\hat{\mathbf{Y}} \mid \mathcal{A}, \mathcal{P}), \quad (11)$$

where visual anchors \mathcal{A} and structured priors \mathcal{P} serve as external consistency checks. For continuous outputs, violations of anchor-implied bounds are softly corrected via constraint-aware projection; for discrete or categorical tasks, the verifier enforces schema consistency and removes logically invalid predictions. This decouples high-level reasoning from strict constraint satisfaction, improving robustness without sacrificing flexibility.

Optimization Strategy. MAS4TS supports end-to-end differentiable execution while selectively training only coordination-critical components. Static preprocessing, VLM prompting/inference, and selected tools in \mathcal{T} are kept *training-free*, serving as fixed reasoning and execution primitives. Optimization focuses on learnable components that govern information flow and coordination, including: (i) embedding/communication functions, (ii) shared memory representations, and (iii) gating and routing parameters Θ . We optimize a task-dependent objective:

$$\mathcal{L} = \mathcal{L}_{\text{task}}(\mathbf{Y}^*, \mathbf{Y}) + \lambda \cdot \mathcal{L}_{\text{com}}(\mathbf{M}), \quad (12)$$

where $\mathcal{L}_{\text{task}}$ corresponds to the downstream objective (e.g., \mathcal{L}_{MSE} for forecasting), and \mathcal{L}_{com} regularizes memory updates to encourage stable inter-agent communication.

MAS4TS integrates multi-agent collaboration, intuitive visual reasoning, and tool-augmented execution into a unified framework. By explicitly separating structural reasoning, continuous latent reconstruction, and tool-integrated execution, MAS4TS achieves strong adaptability across diverse time series tasks while preserving numerical consistency and training efficiency.

4 Experiments

We conduct extensive experiments to evaluate MAS4TS on four fundamental time series tasks: forecasting, classification, imputation, and anomaly detection. We benchmark on 21 widely-used datasets and compare against more than 20 competitive baselines, including task-specific time series models and recent LM-based approaches. For fair comparison, we follow standard protocols for each task (data splits, input/output lengths, and evaluation metrics) and use the same splits and evaluation settings across methods. Experimental details are provided in Appendix A.

4.1 Forecasting

Setups. For long-term forecasting, we evaluate on 7 widely-adopted benchmarks: ETT (4 subsets: ETTh1, ETTh2, ETTm1, ETTm2), Weather, Solar-Energy, and Exchange-Rate. Following established protocols (Wu et al., 2023b; Liu et al., 2024b), we use a fixed look-back window of 96 and evaluate prediction horizons $H \in \{96, 192, 336, 720\}$.

Results. As shown in Table 1, MAS4TS achieves competitive performance across all benchmarks. On the Weather dataset, MAS4TS obtains an average MSE of 0.232, outperforming MAFS (0.252) and TimeVLM (0.262) by 7.9% and 11.5%, respectively. Temporal dependencies exhibit non-stationary patterns with regime shifts on the Exchange dataset, while MAS4TS achieves 0.343 average MSE, demonstrating that visual anchoring effectively captures abrupt trend changes that challenge purely numerical ap-

Table 1. Summarized results for the long-term forecasting task. We report the averaged metrics (Avg) across four prediction lengths {96, 192, 336, 720} for each dataset. Full results are available in Table 9. **Red bold**: the best results; **blue**: the second best.

Models	MAS4TS (Ours)		MAFS 2025		TimeVLM 2025		UniTime 2024a		iTransformer 2024b		PatchTST 2023b		Crossformer 2023		TiDE 2023		TimesNet 2023a		DLinear 2023		SCINet 2022a		FEDformer 2022		Stationary 2022c		Autoformer 2021	
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	0.232	0.270	0.252	0.275	0.262	0.281	0.253	0.276	0.258	0.278	0.265	0.285	0.264	0.320	0.271	0.320	0.259	0.287	0.265	0.315	0.292	0.363	0.309	0.360	0.288	0.314	0.338	0.382
Solar-Energy	0.212	0.276	0.234	0.265	0.259	0.293	0.270	0.315	0.233	0.262	0.287	0.333	0.406	0.442	0.347	0.417	0.403	0.374	0.330	0.401	0.282	0.375	0.328	0.383	0.350	0.390	0.586	0.557
Exchange	0.343	0.397	0.396	0.417	0.355	0.399	0.465	0.462	0.360	0.403	0.367	0.404	0.940	0.707	0.370	0.413	0.416	0.443	0.354	0.414	0.750	0.626	0.519	0.429	0.461	0.454	0.613	0.539
ETTh1	0.440	0.438	0.450	0.440	0.457	0.443	0.442	0.448	0.454	0.447	0.516	0.484	0.529	0.522	0.541	0.507	0.458	0.450	0.461	0.457	0.747	0.647	0.498	0.484	0.570	0.537	0.496	0.487
ETTh2	0.379	0.402	0.380	0.403	0.384	0.405	0.379	0.403	0.383	0.407	0.391	0.411	0.942	0.684	0.611	0.550	0.414	0.427	0.563	0.519	0.954	0.723	0.437	0.449	0.526	0.516	0.450	0.459
ETTm1	0.367	0.381	0.397	0.402	0.388	0.398	0.385	0.399	0.407	0.410	0.406	0.407	0.513	0.495	0.419	0.419	0.400	0.406	0.404	0.408	0.485	0.481	0.448	0.452	0.481	0.456	0.588	0.517
ETTm2	0.278	0.326	0.288	0.330	0.279	0.329	0.293	0.334	0.288	0.332	0.290	0.334	0.757	0.610	0.358	0.404	0.291	0.333	0.354	0.402	0.954	0.723	0.305	0.349	0.306	0.347	0.327	0.371

Table 2. Average classification accuracy (%) on 8 UEA multivariate time series datasets. **Red bold**: the best result; **blue**: the second best. The full results among all datasets are provided in Table 10.

Models	DTW (1994)	XGBoost (2016)	LSTM (1997a)	LSTNet (2018a)	LSSL (2022)	TCN (2019)	Trans. (2017)	Re. (2020)	In. (2021)	Pyra. (2022b)	Auto. (2021)	FED. (2022)	ETS. (2022)	iTrans. (2024b)	DLinear (2023)	LightTS (2023)	TiDE (2023)	MAS4TS (Ours)
Avg. Acc.	62.78	61.48	51.72	66.40	65.30	67.26	67.31	66.93	59.93	65.58	66.00	65.32	65.83	65.16	64.85	64.47	64.16	68.25

proaches. The ETT benchmarks further validate the robustness of our framework: MAS4TS consistently ranks among the top performers across all four subsets, with particularly strong results on ETTh2 (0.379 MSE), where complex multi-scale patterns require both fine-grained numerical reasoning and high-level structural understanding.

4.2 Classification

Setups. We evaluate on 8 multivariate datasets from the UEA archive (Bagnall et al., 2018): EthanolConcentration (chemical), FaceDetection (video), Handwriting (motion capture), Heartbeat (medical), JapaneseVowels (audio), SelfRegulationSCP1/SCP2 (EEG), and UWaveGestureLibrary (gesture recognition). These datasets span diverse domains with varying sequence lengths (from 29 to 1,751), dimensionality (from 3 to 144), and class counts.

Results. Table 2 presents the classification results. MAS4TS achieves the highest average accuracy of 68.25%, outperforming iTransformer (65.16%) and DLinear (64.85%) by 4.7% and 5.2%, respectively. On FaceDetection, which requires capturing subtle temporal dynamics in video sequences, MAS4TS achieves 69.38% accuracy compared to 65.3% for TiDE, demonstrating the discriminative power of VLM-extracted visual features. MAS4TS maintains robust performance across datasets with heterogeneous characteristics, whereas specialized methods like Informer show higher variance.

4.3 Imputation

Setups. We evaluate imputation on the ETT family, Electricity, Weather, and Exchanges datasets under random missing scenarios with mask ratios $r \in \{12.5\%, 25\%, 37.5\%, 50\%\}$.

Missing values are randomly sampled across all time steps and variables, following the protocol in (Wu et al., 2023b).

Results. As shown in Table 3, MAS4TS achieves state-of-the-art imputation performance across all mask ratios. At 50% missing ratio on ETTm2, MAS4TS obtains an MSE of 0.030, reducing the error by 34.8% compared to the second LightTS (0.046). The improvement margin increases with higher mask ratios, demonstrating strong robustness under severe data corruption. On the Exchange dataset, MAS4TS achieves an average of 0.005 MSE, outperforming all the baselines. This validates that our anchor-constrained completion is particularly effective when missing segments are extensive: visual anchors provide reliable structural waypoints, while the latent construction treats these anchors as boundary constraints to produce smooth and temporally coherent trajectories. Although gains on low missing ratios are more modest, the multi-agent coordination mechanism ensures consistent improvements, and the foundation model’s generalization capability offers potential for further enhancement on domain-specific missing patterns.

4.4 Anomaly Detection

Setups. We evaluate on five server and infrastructure monitoring benchmarks: SMD (Su et al., 2019) (server machine dataset with 38 dimensions), MSL and SMAP (Hundman et al., 2018) (spacecraft telemetry from NASA), SWaT (Mathur and Tippenhauer, 2016) (water treatment plant with 51 sensors), and PSM (Abdulaal et al., 2021) (server metrics from eBay). These datasets contain both point anomalies and contextual anomalies with varying contamination ratios.

Results. Table 4 presents the anomaly detection results.

Visual Reasoning over Time Series via Multi-Agent Systems

Table 3. Results for the imputation task. We randomly mask 12.5%, 25%, 37.5% and 50% time points and report the average performance. Full results are available in Table 11. The best and second-best results are **red** and **blue**, respectively.

Models	MAS4TS (Ours)		ETSformer (2022)		LightTS (2022)		DLinear (2023)		FEDformer (2022)		Pyraformer (2021a)		Informer (2021)		Reformer (2020)		LSTM (1997)		TCN (2019)		LSSL (2022)	
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	0.054	0.147	0.120	0.253	0.104	0.218	0.093	0.206	0.062	0.177	0.717	0.570	0.071	0.188	0.050	0.154	0.989	0.786	0.516	0.497	0.113	0.254
ETTh2	0.030	0.103	0.208	0.327	0.046	0.151	0.096	0.208	0.101	0.215	0.465	0.508	0.156	0.292	0.157	0.280	1.027	0.800	0.266	0.407	0.175	0.324
ETTh1	0.129	0.235	0.202	0.329	0.284	0.373	0.201	0.306	0.117	0.246	0.842	0.682	0.161	0.279	0.122	0.245	1.225	0.873	0.621	0.571	0.424	0.481
ETTh2	0.067	0.166	0.367	0.436	0.119	0.250	0.142	0.259	0.163	0.279	1.079	0.792	0.337	0.452	0.234	0.352	2.039	1.114	0.431	0.503	0.495	0.475
Electricity	0.086	0.202	0.214	0.339	0.131	0.262	0.132	0.260	0.130	0.259	0.297	0.382	0.222	0.328	0.200	0.313	0.277	0.365	0.582	0.597	0.222	0.293
Weather	0.036	0.070	0.076	0.171	0.055	0.117	0.052	0.110	0.099	0.203	0.152	0.235	0.045	0.104	0.038	0.087	0.365	0.434	0.183	0.291	0.045	0.108
Exchange	0.005	0.041	0.128	0.262	0.042	0.142	0.046	0.146	0.241	0.324	0.398	0.532	0.439	0.554	0.302	0.439	0.032	0.123	0.043	0.146	0.040	0.139
1st Count	5	6	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0

Table 4. Anomaly detection results measured by F1-score (%). A higher F1 indicates better performance. Full results are available in Table 12. **Red bold**: the best results; **blue**: the second best.

Datasets	Transformer 2017	Reformer 2020	Informer 2021	Autoformer 2021	Pyraformer 2022b	Anomaly T. 2021	Stationary 2022c	ETSformer 2022	DLinear 2023	TiDe 2023	iTransformer 2024b	PatchTST 2023b	MAS4TS (Ours)
SMD	79.56	75.32	81.65	85.11	83.04	85.49	84.62	83.13	77.10	68.91	71.15	84.62	85.10
MSL	78.68	84.40	84.06	79.05	84.86	83.31	77.50	85.03	84.88	70.18	72.54	78.70	82.06
SMAP	69.70	70.40	69.92	71.12	71.09	71.18	71.09	69.50	69.26	64.00	66.87	68.82	68.58
SWaT	80.37	82.80	81.43	92.74	91.78	83.10	79.88	84.91	87.52	76.73	79.18	85.72	93.14
PSM	76.07	73.61	77.10	93.29	82.08	79.40	97.29	91.76	93.55	92.50	95.17	96.08	95.04
Avg F1	76.88	77.31	78.83	84.26	82.57	80.50	82.08	82.87	82.46	74.46	76.98	82.79	84.78

MAS4TS achieves the highest average F1-score of 84.78%, improving over ETSformer (82.87%) and iTransformer (76.98%) by 1.5% and 9.3%, respectively. On SWaT, MAS4TS achieves 93.14% F1-score with balanced precision (91.61%) and recall (94.74%), indicating effective discrimination between normal fluctuations and genuine anomalies. The **Reasoner** with visual reasoning proves particularly valuable for anomaly detection: by directly “seeing” anomalous patterns from time series visualizations, it can identify salient deviations before the model fully converges, pre-locating suspicious regions that guide subsequent numerical analysis. Visual anchors encode expected temporal structures learned from normal patterns, enabling sensitive detection when observations deviate from these structural priors. In addition, the **Executor**’s verification step filters false positives via anchor-consistency checks, improving precision without sacrificing recall.

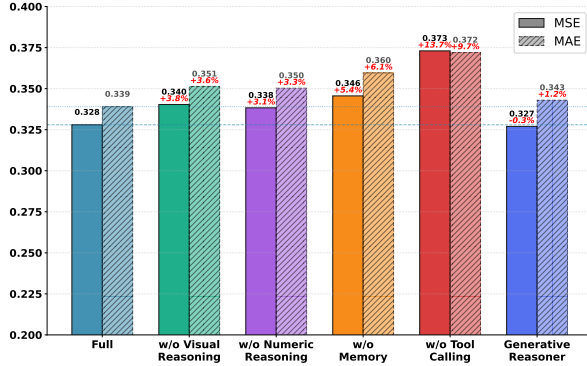


Figure 3. The ablation results on the Weather dataset with prediction horizon $H = 720$. Each component contributes to the performance, while generative reasoner does not improve further.

4.5 Model Analysis

Ablation Studies. We conduct systematic ablation studies to quantify the contribution of key components in MAS4TS. All variants are evaluated on the Weather dataset with prediction horizon $H = 720$. Figure 3 reports results for the full MAS4TS and the following variants: (a) without Visual Reasoning; (b) without Numeric Reasoning; (c) without Shared Memory; (d) without Tools and use a simple predictor in the **Executor**; and (e) *Generative Reasoner*, which replaces the default rule-based rendering pipeline with a generative model (*Nano-Banana*) to generate the visualization, followed by the same anchor extraction procedure.

The ablation results reveal several insights. Removing Tool Calling causes the largest performance degradation, increasing MSE by 13.7% (0.328 \rightarrow 0.373) and MAE by 9.7% (0.339 \rightarrow 0.372), highlighting that adaptive tool invocation and integration are crucial for handling forecasting scenarios. Disabling Visual Reasoning, Numeric Reasoning, or Shared Memory leads to a certain degree of MSE increases (3.1%–5.4%), indicating that these components provide complementary signals and collectively improve robustness. The *Generative Reasoner* variant does not yield a clear improvement over the default instruct-style anchor extraction; nevertheless, this result does not preclude the potential of directly reasoning with generated visualizations under more suitable settings. Overall, these findings suggest that MAS4TS benefits from the coordinated use of multiple reasoning modalities, where tool adaptation plays a central role in effectively leveraging the available components.

Hyperparameter Studies. We analyze the impact of key hyperparameters on ETTh1 with prediction horizon $H =$

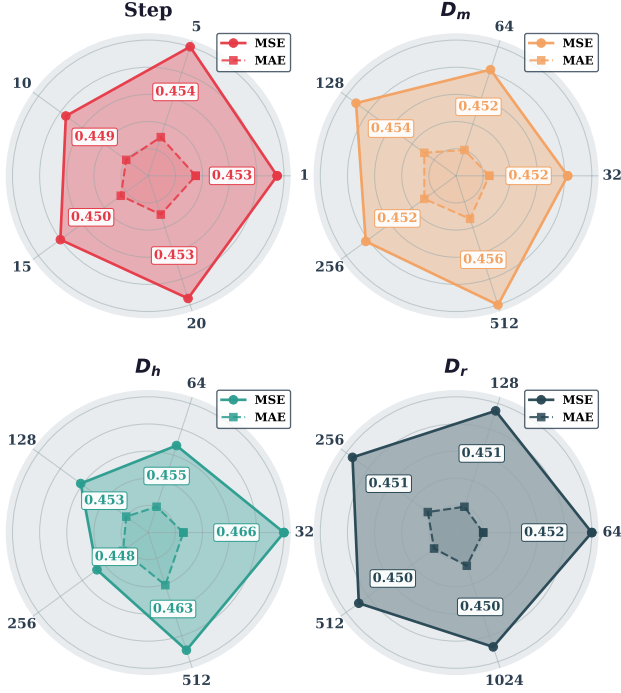


Figure 4. Hyperparameter sensitivity analysis on ETTh1 ($H = 720$). Each radar chart shows MSE (solid) and MAE (dashed) for: a) reasoning steps; b) memory dimension D_m ; c) hidden dimension D_h ; and d) feedforward dimension D_r .

720, as shown in Figure 4. The radar charts visualize MSE and MAE variations across different parameter settings.

① Reasoning Steps: Performance improves as steps increase from 1 to 20, achieving an optimal MSE of 0.449 at 10 steps, then slightly degrades with further increases. This suggests that 10 integration steps provide sufficient numerical precision for the latent ODE-solver while avoiding computational overhead. **② Dimension of Communication Memory D_m :** The model exhibits remarkable robustness to memory dimension variations, with MSE ranging narrowly from 0.452 to 0.456 across all settings (32-512). We adopt $D_m = 64$ as the default for computational efficiency. **③ Dimension of Hidden State D_h :** Performance consistently improves as D_h increases from 32 (MSE=0.466) to 256 (MSE=0.448), but degrades at 512 due to overfitting. The optimal $D_h = 256$ balances model capacity and generalization. **④ Feedforward Dimension D_r :** Similar to D_m , the model shows low sensitivity to D_r , with MSE stabilizing around 0.450 for $D_r \geq 128$. We set $D_r = 512$ as the default, which achieves marginally better performance without high computational cost.

4.6 Efficiency Study

To evaluate the computational efficiency of MAS4TS, we compare it with representative pretrained LM-based time

Table 5. Efficiency comparison between MAS4TS and pretrained LM-based time-series baselines on ETTh1 forecasting ($L=96$, $H=96$). We report trainable parameters, peak GPU memory, training time per epoch, and inference latency per batch. MAS4TS variants use *Qwen3-VL-235B-Thinking* (Thinking) and *Nano-Banana* (Generative), and the baselines use GPT-2 as the text encoding.

Models	Params (M)	Mem (GB)	Train (s)	Infer (ms)
MAS4TS	5.38	0.05	24.55	90.49
- Thinking	5.38	0.05	35.98	133.75
- Generative	5.38	0.05	915.92	2767.01
Time-VLM	151.58	0.72	319.67	1067.14
UniTime [†]	107.32	5.18	53.49	165.77
TimeLLM [†]	134.90	15.96	853.75	2296.28

* Experiments conducted on NVIDIA 4090 with batch size 32.

* The [†] represents the deployment and inference of local LMs, while we reason through the EAS API.

series methods under a unified experimental setting. All measurements are conducted on the forecasting task of the ETTh1 dataset with an input length of 96 and a prediction length of 96. We report the number of trainable parameters, peak GPU memory consumption, end-to-end training time per epoch, and inference latency per batch. MAS4TS uses an instruction-following VLM interface (*Qwen3-VL-235B-Instruct* by default) by default for anchor extraction to balance output quality and efficiency.

Table 5 summarizes the efficiency comparison results. Although MAS4TS integrates a multi-agent workflow with foundation models, it maintains a compact backbone (5.38M parameters) and extremely low peak GPU memory footprint (0.05 GB), which are orders of magnitude smaller than pretrained LM-based methods such as Time-VLM, UniTime, and TimeLLM. Concretely, these baselines involve over 100M parameters and consume several to tens of gigabytes of GPU memory, whereas MAS4TS keeps the learnable core lightweight and invokes the VLM only when needed, in a mini-batched and parallel manner across samples. Default *Instruct-style* anchoring provides the best efficiency-quality trade-off and extracts anchors using a structured, schema-constrained interface that avoids long, free-form decoding, while both thinking and generative variants introduce heavier reasoning overhead. The results indicate that the practical efficiency of MAS4TS comes primarily from lightweight coordination plus structured tool usage, rather than scaling a monolithic architecture.

5 Conclusion

We presented MAS4TS, a tool-driven multi-agent framework for general time-series analysis built on an Analyzer-Reasoner-Executor paradigm. MAS4TS combines VLM-based visual anchoring with anchor-guided latent trajectory reconstruction. It further enables adaptive routing and tool-chain execution to solve forecasting, classification, imputation, and anomaly detection within a unified protocol.

References

- A. Abdulaal, Z. Liu, and T. Lancewicki. Practical approach to asynchronous multivariate time series anomaly detection and localization. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pages 2485–2494, 2021.
- A. Bagnall, H. A. Dau, J. Lines, M. Flynn, J. Large, A. Bostrom, P. Southam, and E. Keogh. The uea multivariate time series classification archive, 2018. *arXiv preprint arXiv:1811.00075*, 2018.
- S. Bai, J. Z. Kolter, and V. Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- D. J. Berndt and J. Clifford. Using dynamic time warping to find patterns in time series. In *KDD Workshop*, 1994.
- G. E. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung. *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. In *Advances in neural information processing systems*, pages 1877–1901, 2020.
- D. Campos, M. Zhang, B. Yang, T. Kieu, C. Guo, and C. S. Jensen. Lightts: Lightweight time series classification with adaptive ensemble distillation. *Proceedings of the ACM on Management of Data*, 1(2):1–27, 2023.
- C. Chang, Y. Shi, D. Cao, W. Yang, J. Hwang, H. Wang, J. Pang, W. Wang, Y. Liu, W.-C. Peng, et al. A survey of reasoning and agentic systems in time series with large language models. *arXiv preprint arXiv:2509.11575*, 2025a.
- C. Chang, W.-Y. Wang, W.-C. Peng, and T.-F. Chen. Llm4ts: Aligning pre-trained llms as data-efficient time-series forecasters. *ACM Transactions on Intelligent Systems and Technology*, 16(3):1–20, 2025b.
- M. Chen, L. Shen, Z. Li, X. J. Wang, J. Sun, and C. Liu. Visions: Visual masked autoencoders are free-lunch zero-shot time series forecasters, 2024. URL <https://arxiv.org/abs/2408.17253>.
- T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. *KDD*, 2016.
- K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- R. B. Cleveland, W. S. Cleveland, J. E. McRae, I. Terpenning, et al. Stl: A seasonal-trend decomposition. *J. off. Stat*, 6(1):3–73, 1990.
- A. Das, W. Kong, A. Leach, S. Mathur, R. Sen, and R. Yu. Long-term forecasting with tide: Time-series dense encoder. *arXiv preprint arXiv:2304.08424*, 2023.
- C. Deb, F. Zhang, J. Yang, S. E. Lee, and K. W. Shah. A review on time series forecasting techniques for building energy consumption. *Renewable and Sustainable Energy Reviews*, 74:902–924, 2017.
- J. R. Dormand and P. J. Prince. A family of embedded runge-kutta formulae. *Journal of computational and applied mathematics*, 6(1):19–26, 1980.
- J.-Y. Franceschi, A. Dieuleveut, and M. Jaggi. Unsupervised scalable representation learning for multivariate time series. In *NeurIPS*, 2019.
- A. Gu, K. Goel, and C. Ré. Efficiently modeling long sequences with structured state spaces. In *ICLR*, 2022.
- Y. Gu, Y. Xiong, J. Mace, Y. Jiang, Y. Hu, B. Kasikci, and P. Cheng. Argos: Agentic time-series anomaly detection with autonomous rule generation via large language models. *arXiv preprint arXiv:2501.14170*, 2025.
- Y. He and J. Zhao. Temporal convolutional networks for anomaly detection in time series. In *Journal of Physics: Conference Series*, volume 1213, page 042050. IOP Publishing, 2019.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 1997a.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997b.
- S. Hong, M. Zhuge, J. Chen, X. Zheng, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2023.
- Q. Huang, Z. Zhou, Y. Li, K. Yang, B. Wang, and Y. Wang. Many minds, one goal: Time series forecasting via sub-task specialization and inter-agent cooperation. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- X. Huang, W. Liu, X. Chen, X. Wang, H. Wang, D. Lian, Y. Wang, R. Tang, and E. Chen. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.
- K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Soderstrom. Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 387–395, 2018.
- S. M. Idrees, M. A. Alam, and P. Agarwal. A prediction approach for stock market volatility based on time series data. *IEEE Access*, 7:17287–17298, 2019.
- M. Jin, S. Wang, L. Ma, Z. Chu, J. Y. Zhang, X. Shi, P.-Y. Chen, Y. Liang, Y.-F. Li, S. Pan, and Q. Wen. Time-LLM: Time series forecasting by reprogramming large language models. In *International Conference on Learning Representations (ICLR)*, 2024a.
- M. Jin, Y. Zhang, W. Chen, K. Zhang, Y. Liang, B. Yang, J. Wang, S. Pan, and Q. Wen. Position paper: What can large language models tell us about time series analysis. *arXiv preprint arXiv:2402.02713*, 2024b.
- Z. Karevan and J. A. Suykens. Transductive lstm for time-series prediction: An application to weather forecasting. *Neural Networks*, 125:1–9, 2020.
- N. Kitaev, L. Kaiser, and A. Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.

- G. Lai, W.-C. Chang, Y. Yang, and H. Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *SIGIR*, 2018a.
- G. Lai, W.-C. Chang, Y. Yang, and H. Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 95–104, 2018b.
- S. Li, X. Jin, Y. Xuan, X. Zhou, W. Chen, Y.-X. Wang, and X. Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances in Neural Information Processing Systems*, 32, 2019.
- Y. Liang, H. Wen, Y. Nie, Y. Jiang, M. Jin, D. Song, S. Pan, and Q. Wen. Foundation models for time series analysis: A tutorial and survey. In *Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining*, pages 6555–6565, 2024.
- B. Lim and S. Zohren. Time-series forecasting with deep learning: a survey. *Philosophical transactions of the royal society a: mathematical, physical and engineering sciences*, 379(2194), 2021.
- M. Liu, A. Zeng, M. Chen, Z. Xu, Q. Lai, L. Ma, and Q. Xu. Scinet: Time series modeling and forecasting with sample convolution and interaction. *Advances in Neural Information Processing Systems*, 35:5816–5828, 2022a.
- S. Liu, H. Yu, C. Liao, J. Li, W. Lin, A. X. Liu, and S. Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International Conference on Learning Representations*, 2022b.
- X. Liu, J. Hu, Y. Li, S. Diao, Y. Liang, B. Hooi, and R. Zimmermann. Unitime: A language-empowered unified model for cross-domain time series forecasting. In *Proceedings of the ACM on Web Conference 2024*, pages 4095–4106, 2024a.
- Y. Liu, H. Wu, J. Wang, and M. Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. *Advances in Neural Information Processing Systems*, 35:9881–9893, 2022c.
- Y. Liu, T. Hu, H. Zhang, H. Wu, S. Wang, L. Ma, and M. Long. itransformer: Inverted transformers are effective for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024b.
- S. Lyu, S. Zhong, W. Ruan, Q. Liu, Q. Wen, H. Xiong, and Y. Liang. Occamvts: Distilling vision models to 1 URL <https://arxiv.org/abs/2508.01727>.
- A. P. Mathur and N. O. Tippenhauer. Swat: A water treatment testbed for research and training on ics security. In *2016 international workshop on cyber-physical systems for smart water networks (CySWater)*, pages 31–36. IEEE, 2016.
- L. R. Medsker, L. Jain, et al. Recurrent neural networks. *Design and Applications*, 5(64-67):2, 2001.
- Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations*, 2023a.
- Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations*, 2023b.
- C. Qian, W. Liu, H. Liu, N. Chen, Y. Dang, J. Li, C. Yang, W. Chen, Y. Su, X. Cong, et al. Chatdev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.
- Y. Qin, S. Liang, Y. Ye, K. Zhu, L. Yan, Y. Lu, Y. Lin, X. Cong, X. Tang, B. Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- C. Ravuru, S. S. Sakhinana, and V. Runkana. Agentic retrieval-augmented generation for time series analysis. *arXiv preprint arXiv:2408.14484*, 2024.
- W. Ruan, S. Zhong, H. Wen, and Y. Liang. Vision-enhanced time series forecasting via latent diffusion models. *arXiv preprint arXiv:2502.14887*, 2025.
- Y. Rubanova, R. T. Chen, and D. K. Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. *Advances in neural information processing systems*, 32, 2019.
- T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer, N. Cancedda, and T. Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- S. H. Schneider and R. E. Dickinson. Climate modeling. *Reviews of Geophysics*, 12(3):447–493, 1974.
- C. Shen, W. Yu, Z. Zhao, D. Song, W. Cheng, H. Chen, and J. Ni. Multi-modal view enhanced large vision models for long-term time series forecasting. *arXiv preprint arXiv:2505.24003*, 2025.
- Y. Su, Y. Zhao, C. Niu, R. Liu, W. Sun, and D. Pei. Robust anomaly detection for multivariate time series through stochastic recurrent neural network. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2828–2837, 2019.
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- A. Van Den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu, et al. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 12:1, 2016.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- G. Woo, C. Liu, D. Sahoo, A. Kumar, and S. Hoi. Etsformer: Exponential smoothing transformers for time-series forecasting. In *arXiv preprint arXiv:2202.01381*, 2022.
- H. Wu, J. Xu, J. Wang, and M. Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34: 22419–22430, 2021.

- H. Wu, T. Hu, Y. Liu, H. Zhou, J. Wang, and M. Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *International Conference on Learning Representations*, 2023a.
- H. Wu, T. Hu, Y. Liu, H. Zhou, J. Wang, and M. Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *International Conference on Learning Representations*, 2023b.
- Q. Wu, G. Bansal, J. Zhang, Y. Wu, B. Li, E. Zhu, L. Jiang, X. Zhang, S. Zhang, J. Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- J. Xu, H. Wu, J. Wang, and M. Long. Anomaly transformer: Time series anomaly detection with association discrepancy. *arXiv preprint arXiv:2110.02642*, 2021.
- Z. Xu, A. Zeng, and Q. Xu. Fits: Modeling time series with 10k parameters. *arXiv preprint arXiv:2307.03756*, 2023.
- A. Zeng, M. Chen, L. Zhang, and Q. Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, pages 11121–11128, 2023.
- Y. Zhang and J. Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. *ICLR*, 2023.
- J. Zheng and M. Huang. Traffic flow forecast through time series analysis based on deep learning. *IEEE Access*, 8:82562–82570, 2020.
- S. Zhong, W. Ruan, M. Jin, H. Li, Q. Wen, and Y. Liang. Time-vlm: Exploring multimodal vision-language models for augmented time series forecasting. *arXiv preprint arXiv:2502.04395*, 2025.
- H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on Artificial Intelligence*, pages 11106–11115, 2021.
- T. Zhou, Z. Ma, Q. Wen, X. Wang, L. Sun, and R. Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning*, pages 27268–27286, 2022.
- T. Zhou, P. Niu, X. Wang, L. Sun, and R. Jin. One fits all: Power general time series analysis by pretrained lm. In *Advances in Neural Information Processing Systems*, 2023.
- Z. Zhou, G. Lyu, Y. Huang, Z. Wang, Z. Jia, and Z. Yang. Sd-former: transformer with spectral filter and dynamic attention for multivariate time series long-term forecasting.
- Z. Zhou, J. Hu, Q. Wen, J. T. Kwok, and Y. Liang. Multi-order wavelet derivative transform for deep time series forecasting. *arXiv preprint arXiv:2505.11781*, 2025.

A Experiment Details

A.1 Dataset Details

Table 6. Dataset descriptions. The dataset size is organized in (Train, Validation, Test).

Tasks	Dataset	Dim	Series Length	Dataset Size	Information (Frequency)
Long-term Forecasting	ETTM1, ETTm2	7	{96, 192, 336, 720}	(34465, 11521, 11521)	Electricity (15 mins)
	ETTh1, ETTh2	7	{96, 192, 336, 720}	(8545, 2881, 2881)	Electricity (Hourly)
	Weather	21	{96, 192, 336, 720}	(36792, 5271, 10540)	Weather (10 mins)
	Solar-Energy	137	{96, 192, 336, 720}	(36345, 5065, 10321)	Solar (10 mins)
	Exchange	8	{96, 192, 336, 720}	(5120, 665, 1422)	Exchange rate (Daily)
Imputation	ETTM1, ETTm2	7	96	(34465, 11521, 11521)	Electricity (15 mins)
	ETTh1, ETTh2	7	96	(8545, 2881, 2881)	Electricity (Hourly)
	Electricity	321	96	(18317, 2633, 5261)	Electricity (Hourly)
	Weather	21	96	(36792, 5271, 10540)	Weather (10 mins)
	Exchange	8	96	(5120, 665, 1422)	Exchange rate (Daily)
Classification (UEA)	EthanolConcentration	3	1751	(261, 0, 263)	Alcohol Industry
	FaceDetection	144	62	(5890, 0, 3524)	Face (250Hz)
	Handwriting	3	152	(150, 0, 850)	Handwriting
	Heartbeat	61	405	(204, 0, 205)	Heart Beat
	JapaneseVowels	12	29	(270, 0, 370)	Voice
	SelfRegulationSCP1	6	896	(268, 0, 293)	Health (256Hz)
	SelfRegulationSCP2	7	1152	(200, 0, 180)	Health (256Hz)
	UWaveGestureLibrary	3	315	(120, 0, 320)	Gesture
Anomaly Detection	SMD	38	100	(566724, 141681, 708420)	Server Machine
	MSL	55	100	(44653, 11664, 73729)	Spacecraft
	SMAP	25	100	(108146, 27037, 427617)	Spacecraft
	SWaT	51	100	(396000, 99000, 449919)	Infrastructure
	PSM	25	100	(105984, 26497, 87841)	Server Machine

As shown in Table 6, the benchmark datasets span diverse domains and temporal granularities:

- **Long-term Forecasting:** 7 datasets covering electricity transformer temperature (ETT series with 4 subsets), electricity consumption (Electricity), transportation flow (Traffic), meteorological indicators (Weather), solar power generation (Solar-Energy), and currency exchange rates (Exchange).
- **Imputation:** The same forecasting datasets with artificially induced random missing patterns at 12.5%, 25%, 37.5%, and 50% mask ratios, following the protocol established in (Wu et al., 2023b).
- **Classification:** 8 representative datasets from the UEA multivariate time series archive (Bagnall et al., 2018), covering chemical sensing (EthanolConcentration), video analysis (FaceDetection), motion capture (Handwriting, UWaveGestureLibrary), medical diagnostics (Heartbeat), audio recognition (JapaneseVowels), and brain-computer interface (SelfRegulationSCP1, SelfRegulationSCP2).
- **Anomaly Detection:** 5 real-world datasets from server machines (SMD, PSM), spacecraft telemetry systems (MSL, SMAP), and water treatment infrastructure (SWaT).

A.2 Baseline Methods

We reproduce all baselines using the official codebases and/or settings reported in the original papers. To ensure a fair comparison, we keep the *input embedding* and *final projection* interfaces consistent across neural models when applicable, and focus on comparing their core sequence modeling or task-specific mechanisms. Note that some baselines are only evaluated on specific tasks.

(1) Time-series specialized models.

• Transformer-based time-series models:

- *iTransformer* (Liu et al., 2024b): performs self-attention across variates (inverted tokenization) to better model multivariate dependencies.
- *PatchTST* (Nie et al., 2023b): uses patching and channel-independence to reduce sequence length and stabilize long-horizon learning.
- *Crossformer* (Zhang and Yan, 2023): models cross-dimension dependencies via a two-stage attention design (intra-/inter-dimension).
- *Non-stationary Transformer* (Liu et al., 2022c) (reported as *Stationary*): mitigates distribution shift with series stationarization and de-stationary attention.
- *Autoformer* (Wu et al., 2021): replaces point-wise attention with auto-correlation and progressive series decomposition for long-term forecasting.
- *FEDformer* (Zhou et al., 2022): conducts frequency-enhanced decomposition and attention in the spectral domain for efficient long-sequence modeling.
- *ETSformer* (Woo et al., 2022): injects exponential smoothing (ETS-style) components to explicitly capture level/trend/seasonality.
- *Informer* (Zhou et al., 2021): adopts ProbSparse attention to reduce attention cost for long sequence forecasting.
- *Reformer* (Kitaev et al., 2020): uses LSH attention to approximate softmax attention with lower memory/computation.
- *Pyraformer* (Liu et al., 2022b): builds pyramidal attention across multiple temporal resolutions to model long-range structure.

• CNN/MLP-style time-series models:

- *TimesNet* (Wu et al., 2023b): reshapes 1D series into 2D representations guided by multi-periodicity and applies 2D convolution blocks.
- *DLinear* (Zeng et al., 2023): decomposes series into trend and residual parts and forecasts with lightweight linear layers.
- *LightTS* (Campos et al., 2023): leverages efficient MLP mixing with interval/continuous sampling to handle long contexts.
- *TiDE* (Das et al., 2023): encodes the look-back with covariates via dense residual blocks and decodes future values with a per-time-step temporal decoder.
- *SCINet* (Liu et al., 2022a): recursively splits and interacts subsequences to capture hierarchical temporal dependencies.
- *TCN* (He and Zhao, 2019): uses dilated causal convolutions to model long receptive fields without recurrence.

• RNN/State-space style models:

- *LSTM* (Hochreiter and Schmidhuber, 1997b): recurrent gating mechanism for learning long/short-term temporal dependencies.
- *LSTNet* (Lai et al., 2018b): combines CNN feature extraction with RNN temporal modeling and skip connections for periodic patterns.
- *LSSL* (Gu et al., 2022): reparameterizes HiPPO-initialized state space models as diagonal-plus-low-rank, enabling near-linear-time long-sequence modeling via efficient Cauchy-kernel convolution..

(2) Pretrained LM-based baselines.

- *MAFS* (Huang et al., 2025): employs a multi-agent framework where each agent specializes in a forecasting sub-task, with inter-agent communication and a voting aggregator for final prediction. (Note: we include it here as it shares the pre-train then fine-tune paradigm, though it pre-trains *iTransformer* on sub-tasks rather than leveraging large language models.)
- *Time-VLM* (Zhong et al., 2025): integrates temporal, visual, and textual modalities through pretrained Vision-Language Models (e.g., ViLT, CLIP). Time series are transformed into images via frequency/periodicity encoding and multi-scale convolution, while contextual text prompts are generated from statistical features. A frozen VLM encodes both modalities, and cross-modal attention fuses them with temporal memory features for forecasting.
- *UniTime* (Liu et al., 2024a): employs pretrained language models with domain instructions. It adopts channel-independence and patch-based tokenization, uses natural language instructions to provide domain identification and alleviate domain confusion, and applies masking to balance convergence speeds across domains. The Language-TS Transformer aligns time series from various input spaces to the common latent space of language models for cross-domain generalization.

(3) Classical baselines and others.

- *Transformer* (Vaswani et al., 2017): the vanilla encoder-decoder attention architecture as a general-purpose sequence baseline.
- *DTW* (Berndt and Clifford, 1994): uses Dynamic Time Warping distance (typically with 1-NN) as a strong classical baseline for time-series classification.

- *XGBoost* (Chen and Guestrin, 2016): gradient-boosted trees trained on engineered temporal features/statistics.

For anomaly detection, we additionally compare against *Anomaly Transformer* (Xu et al., 2021), a method specifically designed for this task that identifies anomalies via association discrepancy between series-wise and prior associations. Beyond task-specific models, we also consider classical methods that remain highly competitive on experimental results. For instance, in classification, *DTW*, *XGBoost*, and the vanilla *Transformer* have all demonstrated strong performance and are thus included as additional baselines.

A.3 Evaluation Metrics

For the metrics, we adopt different evaluation criteria tailored to each task:

Long-term Forecasting and Imputation. We use mean square error (MSE) and mean absolute error (MAE):

$$\text{MSE} = \frac{1}{H \cdot C} \sum_{i=1}^H \sum_{j=1}^C (\mathbf{X}_{i,j} - \hat{\mathbf{X}}_{i,j})^2, \quad \text{MAE} = \frac{1}{H \cdot C} \sum_{i=1}^H \sum_{j=1}^C |\mathbf{X}_{i,j} - \hat{\mathbf{X}}_{i,j}|,$$

where $\mathbf{X}, \hat{\mathbf{X}} \in \mathbb{R}^{H \times C}$ denote the ground truth and predicted values with H time steps and C feature dimensions. Lower values indicate better performance.

Classification. We adopt accuracy (ACC) as the primary evaluation metric:

$$\text{ACC} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Samples}} \times 100\%.$$

Higher accuracy indicates better classification performance.

Anomaly Detection. We adopt precision (P), recall (R), and F1-score (F1) following the point-adjust protocol (Xu et al., 2021):

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{F1} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}},$$

where TP, FP, and FN denote true positives, false positives, and false negatives, respectively. The F1-score provides a balanced measure between precision and recall, with higher values indicating better detection performance.

A.4 Model and Agent Configuration

This section provides a comprehensive overview of the model hyperparameters, agent-specific configurations, and shared memory mechanism in MAS4TS. The framework adopts a modular design where each agent maintains specialized parameter sets while coordinating through a unified memory interface.

Training Configuration. Table 7 summarizes the task-specific training configurations. All experiments employ the AdamW optimizer with default momentum parameters $(\beta_1, \beta_2) = (0.9, 0.999)$ and weight decay of 10^{-2} . We apply early stopping with a patience of 3 epochs to prevent overfitting and use cosine annealing for learning rate scheduling. Mixed-precision training (FP16) is enabled by default to accelerate computation while maintaining numerical stability.

Table 7. Task-specific training configurations for MAS4TS. The embedding dimension d_{model} and network depth are kept consistent across tasks to ensure fair comparison, while learning rates and batch sizes are tuned per task for optimal convergence.

Task	Architecture				Training			
	d_{model}	Layers	Patch	Stride	LR	Loss	Batch	Epochs
Long-term Forecast	64	2	16	8	10^{-4}	MSE	32	10
Imputation	64	2	16	8	10^{-3}	MSE	16	10
Classification	64	2	16	8	10^{-3}	CE	16	30
Anomaly Detection	64	3	16	8	10^{-4}	MSE	128	10

LR: initial learning rate with cosine annealing decay. CE: Cross-Entropy loss.

Ablation Flags. MAS4TS supports fine-grained ablation studies through configurable flags that enable/disable individual components:

- `enable_visual_reasoner`: Toggle VLM-based anchor extraction (default: True)
- `enable_numeric_reasoner`: Toggle Neural ODE sequence reconstruction (default: True)

Visual Reasoning over Time Series via Multi-Agent Systems

- `enable_shared_memory`: Toggle cross-agent memory sharing (default: True)
- `enable_gated_attention`: Toggle gated vs. simple MLP communication (default: True)
- `enable_tools`: Toggle tool-augmented execution vs. simple MLP heads (default: True)
- `completion_strategy`: ODE solver alternative: {'ode', 'linear', 'quadratic', 'repeat'}

Table 8. MAS4TS Model Architecture Parameters

Parameter	Value	Description
<i>Global Configuration</i>		
<code>d_model</code>	64	Embedding dimension for all agents
<code>d_memory</code>	64	Shared memory dimension
<code>patch_len</code>	16	Patch length for tokenization
<code>stride</code>	8	Stride between patches
<i>Data Analyzer Agent</i>		
<code>top_k_features</code>	10	Top-k features for covariance selection
<code>feature_method</code>	covariance	Feature selection criterion
<i>Visual Reasoner Agent</i>		
<code>vlm_model</code>	Qwen3-VL-235B-Instruct	Vision-Language Model backbone
<code>max_anchors</code>	20	Maximum anchor points per sample
<code>confidence_threshold</code>	0.7	Minimum confidence for anchor acceptance
<i>Numeric Reasoner Agent</i>		
<code>hidden_dim</code>	128	Hidden dimension for ODE dynamics
<code>ode_solver</code>	RK4	ODE integration method
<code>ode_step</code>	0.05	Integration step size
<i>Task Executor Agent</i>		
<code>e_layers</code>	2	Encoder layers in task tools
<code>n_heads</code>	4	Attention heads
<code>dropout</code>	0.1	Dropout rate

A.5 Implementation Details

We provide the dataset descriptions and experiment configurations in Table 6 and 7. All experiments are conducted on a single NVIDIA GeForce RTX 4090 GPU (48GB version) with a fixed random seed for reproducibility. The MAS4TS framework is implemented in PyTorch 2.9.1+cu128 with CUDA 12.8. For the Vision-Language Model component, we utilize the *Qwen2-VL-235B-Instruct* model accessed through EAS API services of the Alibaba Cloud Platform.

B Complete Results

This appendix provides comprehensive experimental results for all four evaluation tasks: long-term forecasting, time series classification, imputation, and anomaly detection.

B.1 Long-term Forecasting

Table 9 presents complete forecasting results across all datasets and prediction horizons $H \in \{96, 192, 336, 720\}$. The input sequence length is fixed at $L = 96$ for all methods following prior work (Wu et al., 2023a). MAS4TS consistently achieves competitive or superior performance compared to both task-specific models (e.g., PatchTST, iTransformer) and pretrained LM-based methods (e.g., UniTime, TimeVLM).

B.2 Classification

Table 10 shows classification accuracy on 8 representative datasets from the UEA multivariate time series archive (Bagnall et al., 2018). These datasets span diverse application domains including chemical sensing, video analysis, motion capture, medical diagnostics, and audio recognition.

Visual Reasoning over Time Series via Multi-Agent Systems

Table 9. Full results for the long-term forecasting task. We compare extensive competitive models under different prediction lengths. The input sequence length is set to 96 for all datasets. Avg is averaged from all four prediction lengths, that {96, 192, 336, 720}. **Red bold**: the best results; underlined in blue: the second best.

Models	MAS4TS (Ours)	MAFS 2025	TimeVLM 2025	UniTime 2024a	iTransformer 2024b	PatchTST 2023b	Crossformer 2023	TiDE 2023	TimesNet 2023a	DLinear 2023	SCINet 2022a	FEDformer 2022	Stationary 2022c	Autoformer 2021
Metric	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE
Weather	96	0.147 0.197	<u>0.166 0.208</u>	0.181 0.220	0.171 0.214	0.174 0.214	0.186 0.227	0.195 0.271	0.202 0.261	0.172 0.220	0.195 0.252	0.221 0.306	0.217 0.296	0.173 0.223
	192	0.192 0.244	<u>0.218 0.253</u>	0.227 0.260	0.217 0.254	0.221 0.254	0.234 0.265	<u>0.209 0.277</u>	0.242 0.298	0.219 0.261	0.237 0.295	0.261 0.340	0.276 0.336	0.245 0.285
	336	0.245 0.286	0.274 0.295	0.281 0.298	0.274 <u>0.293</u>	0.278 0.296	0.284 0.301	<u>0.273 0.332</u>	0.287 0.335	0.280 0.306	0.282 0.331	0.309 0.378	0.339 0.380	0.321 0.338
	720	0.328 0.339	0.351 0.345	0.358 0.348	0.351 <u>0.343</u>	0.358 0.347	0.356 0.349	0.379 0.401	0.351 0.386	0.365 0.359	<u>0.345 0.382</u>	0.377 0.427	0.403 0.428	0.414 0.410
	Avg	0.228 0.267	<u>0.252 0.275</u>	0.262 0.281	0.253 0.276	0.258 0.278	0.265 0.285	0.264 0.320	0.271 0.320	0.259 0.287	0.265 0.315	0.292 0.363	0.309 0.360	0.288 0.314
Solar-Energy	96	0.192 0.259	<u>0.198 0.237</u>	0.220 0.269	0.239 0.301	0.203 <u>0.237</u>	0.265 0.323	0.232 0.302	0.312 0.399	0.373 0.358	0.290 0.378	0.237 0.344	0.286 0.341	0.321 0.380
	192	0.207 0.271	<u>0.231 0.263</u>	0.257 0.293	0.269 0.316	0.233 0.261	0.288 0.332	0.371 0.410	0.339 0.416	0.397 0.376	0.320 0.398	0.280 0.380	0.291 0.337	0.346 0.369
	336	0.219 0.283	0.250 <u>0.278</u>	0.271 0.303	0.288 0.324	<u>0.248 0.273</u>	0.301 0.339	0.495 0.515	0.368 0.430	0.420 0.380	0.353 0.415	0.304 0.389	0.354 0.416	0.357 0.387
	720	0.231 0.293	0.255 0.281	0.287 0.307	0.285 0.320	<u>0.249 0.275</u>	0.295 0.336	0.526 0.542	0.370 0.425	0.420 0.381	0.357 0.413	0.308 0.388	0.380 0.437	0.375 0.424
	Avg	0.212 0.276	0.234 <u>0.265</u>	0.259 0.293	0.270 0.315	<u>0.233 0.262</u>	0.287 0.333	0.406 0.442	0.347 0.417	0.403 0.374	0.330 0.401	0.282 0.375	0.328 0.383	0.350 0.390
Exchange	96	0.088 0.208	<u>0.084 0.204</u>	0.082 0.198	0.125 0.245	0.086 0.206	0.088 0.205	0.256 0.367	0.094 0.218	0.107 0.234	0.088 0.218	0.267 0.396	0.148 0.278	0.111 0.237
	192	0.181 0.303	0.180 0.304	0.175 0.296	0.236 0.343	0.177 <u>0.299</u>	<u>0.176 0.299</u>	0.470 0.509	0.184 0.307	0.226 0.344	0.176 0.315	0.351 0.459	0.271 0.315	0.219 0.335
	336	0.293 0.396	0.316 0.406	0.324 0.412	0.420 0.470	0.331 <u>0.417</u>	<u>0.301 0.397</u>	1.268 0.883	0.349 0.431	0.367 0.448	0.313 0.427	1.324 0.853	0.460 0.427	0.421 0.476
	720	0.810 0.682	1.002 0.755	0.840 <u>0.689</u>	1.078 0.791	0.847 0.691	0.901 0.714	1.767 1.068	0.852 0.698	0.964 0.746	<u>0.839 0.695</u>	1.058 0.797	1.195 0.695	1.092 0.769
	Avg	0.343 0.397	0.396 0.417	0.355 <u>0.399</u>	0.465 0.462	0.360 0.403	0.367 0.404	0.940 0.707	0.370 0.413	0.416 0.443	<u>0.354 0.414</u>	0.750 0.626	0.519 0.429	0.461 0.454
ETTh1	96	<u>0.378 0.399</u>	0.389 0.403	0.376 0.395	0.397 0.418	0.386 <u>0.405</u>	0.460 0.447	0.423 0.448	0.479 0.464	0.384 0.402	0.397 0.412	0.654 0.599	0.395 0.424	0.513 0.491
	192	0.419 0.428	0.450 0.438	<u>0.426 0.424</u>	0.434 0.439	0.441 0.512	0.477 0.429	0.471 0.474	0.525 0.492	0.436 0.429	0.446 0.441	0.719 0.631	0.469 0.470	0.534 0.504
	336	0.459 0.445	0.475 <u>0.447</u>	0.469 0.449	<u>0.468 0.457</u>	0.487 0.458	0.546 0.496	0.570 0.546	0.565 0.515	0.491 0.469	0.489 0.467	0.778 0.659	0.530 0.499	0.588 0.535
	720	0.503 0.501	<u>0.487 0.472</u>	0.559 0.504	0.469 0.477	0.503 0.491	0.544 0.517	0.653 0.621	0.594 0.558	0.521 0.500	0.513 0.510	0.836 0.699	0.598 0.544	0.643 0.616
	Avg	0.440 0.438	0.450 <u>0.440</u>	0.457 0.443	0.442 0.448	0.454 0.447	0.516 0.484	0.529 0.522	0.541 0.507	0.458 0.450	0.461 0.457	0.747 0.647	0.498 0.484	0.570 0.537
ETTm2	96	<u>0.296 0.345</u>	0.298 0.346	0.290 0.339	0.296 0.345	0.297 0.349	0.308 0.355	0.745 0.584	0.400 0.440	0.340 0.374	0.340 0.394	0.707 0.621	0.358 0.397	0.476 0.458
	192	0.386 0.399	0.380 0.398	0.374 0.389	<u>0.374 0.394</u>	0.380 0.400	0.393 0.405	0.877 0.656	0.528 0.509	0.402 0.414	0.482 0.479	0.860 0.689	0.429 0.439	0.512 0.493
	336	0.407 0.422	<u>0.413 0.427</u>	0.433 0.440	0.415 0.427	0.428 0.432	0.427 0.436	1.043 0.731	0.643 0.571	0.452 0.452	0.591 0.541	1.000 0.744	0.496 0.487	0.552 0.551
	720	<u>0.427 0.442</u>	0.428 <u>0.444</u>	0.438 0.452	0.425 0.444	0.427 0.445	0.436 0.450	1.104 0.763	0.874 0.679	0.462 0.468	0.839 0.661	1.249 0.838	0.463 0.474	0.562 0.560
	Avg	0.379 0.402	0.380 <u>0.403</u>	0.384 0.405	<u>0.379 0.403</u>	0.383 0.407	0.391 0.411	0.942 0.684	0.611 0.550	0.414 0.427	0.563 0.519	0.954 0.723	0.437 0.449	0.526 0.516
ETTm1	96	0.296 0.341	0.329 <u>0.362</u>	<u>0.326 0.365</u>	0.322 0.363	0.334 0.368	0.352 0.374	0.404 0.426	0.364 0.387	0.338 0.375	0.346 0.374	0.418 0.438	0.379 0.419	0.386 0.398
	192	0.341 0.366	0.370 0.385	0.368 <u>0.384</u>	<u>0.366 0.387</u>	0.390 0.393	0.374 0.387	0.450 0.451	0.398 0.404	0.374 0.387	0.382 0.391	0.439 0.450	0.426 0.441	0.459 0.444
	336	0.385 0.390	0.409 0.410	<u>0.394 0.405</u>	0.398 0.407	0.426 0.420	0.421 0.414	0.532 0.515	0.428 0.425	0.410 0.411	0.415 0.415	0.490 0.485	0.445 0.459	0.495 0.464
	720	0.445 0.426	0.478 0.451	0.462 <u>0.440</u>	<u>0.454 0.440</u>	0.491 0.459	0.462 0.449	0.666 0.589	0.487 0.461	0.478 0.450	0.473 0.451	0.595 0.550	0.543 0.490	0.585 0.516
	Avg	0.367 0.381	0.397 0.402	0.388 <u>0.398</u>	<u>0.385 0.399</u>	0.407 0.410	0.406 0.407	0.513 0.495	0.419 0.419	0.400 0.406	0.404 0.408	0.485 0.481	0.448 0.452	0.481 0.456
ETTh2	96	0.175 0.260	0.180 0.264	<u>0.176 0.261</u>	0.183 0.266	0.180 0.264	0.183 0.270	0.287 0.366	0.207 0.305	0.187 0.267	0.193 0.293	0.286 0.377	0.203 0.287	0.192 0.274
	192	0.237 0.302	0.245 0.304	<u>0.240 0.302</u>	0.251 0.310	0.250 0.309	0.255 0.314	0.414 0.492	0.290 0.364	0.249 0.309	0.284 0.361	0.399 0.445	0.269 0.328	0.280 0.339
	336	<u>0.303 0.343</u>	0.316 0.351	0.299 0.339	0.319 0.351	0.311 0.348	0.309 0.347	0.597 0.542	0.377 0.422	0.321 0.351	0.382 0.429	0.637 0.591	0.325 0.366	0.334 0.361
	720	0.396 0.400	0.409 <u>0.403</u>	<u>0.400 0.416</u>	0.420 0.410	0.412 0.407	0.412 0.404	1.730 1.042	0.558 0.524	0.408 0.403	0.558 0.525	0.960 0.735	0.421 0.415	0.417 0.413
	Avg	0.278 0.326	0.288 0.330	<u>0.279 0.329</u>	0.293 0.334	0.288 <u>0.332</u>	0.290 0.334	0.757 0.610	0.358 0.404	0.291 0.333	0.354 0.402	0.954 0.723	0.305 0.349	0.306 0.347
1 st Count	27 22	0 2	6 7	2 0	0 4	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 0	0 0

Table 10. Full results for the classification task. * in the model names indicates the name of *former. We report the classification accuracy (%) as the result. The standard deviation is within 0.1%. Red Bold: the best results; underlined in blue: the second best.

Methods	DTW	XGBoost	LSTM	LSTNet	LSSL	TCN	Trans.	Re.	In.	Pyra.	Auto.	FED.	ETS.	iTrans.	DLinear	LightTS	TiDE	MAS4TS
Datasets	(1994)	(2016)	(1997a)	(2018a)	(2022)	(2019)	(2017)	(2020)	(2021)	(2022b)	(2021)	(2022)	(2022)	(2024b)	(2023)	(2023)	(2023)	(Ours)
EthanolConcentration	32.3	43.7	32.3	39.9	31.1	28.9	32.7	31.9	31.6	30.8	31.6	28.1	31.2	28.1	32.6	29.7	27.1	33.46
FaceDetection	52.9	63.3	57.7	65.7	66.7	52.8	67.3	68.6	67.0	65.7	68.4	66.0	66.3	66.3	68.0	67.5	65.3	69.38
Handwriting	28.6	15.8	15.2	25.8	24.6	53.3	32.0	27.4	32.8	29.4	36.7	28.0	32.5	24.2	27.0	26.1	23.2	33.64
Heartbeat	71.7	73.2	72.2	77.1	72.7	75.6	76.1	77.1	80.5	75.6	74.6	73.7	71.2	75.6	75.1	75.1	74.6	78.04
JapaneseVowels	94.9	86.5	79.7	98.1	98.4	98.9	98.7	97.8	98.9	98.4	96.2	98.4	95.9	96.6	96.2	96.2	95.6	97.56
SelfRegulationSCP1	77.7	84.6	68.9	84.0	90.8	84.6	92.2	90.4	90.1	88.1	84.0	88.7	89.6	90.2	87.3	89.8	89.2	91.81
SelfRegulationSCP2	53.9	48.9	46.6	52.8	52.2	55.6	53.9	56.7	53.3	53.3	50.6	54.4	55.0	54.4	50.5	51.1	53.4	55.55
UWaveGestureLibrary	90.3	75.9	41.2	87.8	85.9	88.4	85.6	85.6	85.6	83.4	85.9	85.3	85.0	85.9	82.1	80.3	84.9	86.56
Average Accuracy	62.78	61.48	51.72	66.40	65.30	67.26	67.31	66.93	59.93	65.58	66.00	65.32	65.83	65.16	64.85	64.47	64.16	68.25

Visual Reasoning over Time Series via Multi-Agent Systems

Table 11. Full results for the imputation task. We randomly mask 12.5%, 25%, 37.5% and 50% time points to compare the model performance under different missing degrees. * in the Transformers indicates the name of *former.

Models	MAS4TS (Ours)		ETS. (2022)		LightTS* (2022)		DLinear* (2023)		FED. (2022)		Pyra. (2021a)		In. (2021)		Re. (2020)		LSTM (1997)		TCN (2019)		LSSL (2022)		
Mask Ratio	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTm1	12.5%	0.046	0.137	0.067	0.188	0.075	0.180	0.058	0.162	0.035	0.135	0.670	0.541	0.047	0.155	0.032	0.126	0.974	0.780	0.510	0.493	0.101	0.231
	25%	0.049	0.143	0.096	0.229	0.093	0.206	0.080	0.193	0.052	0.166	0.689	0.553	0.063	0.180	0.042	0.146	1.032	0.807	0.518	0.500	0.106	0.235
	37.5%	0.054	0.150	0.133	0.271	0.113	0.231	0.103	0.219	0.069	0.191	0.737	0.581	0.079	0.200	0.052	0.158	0.999	0.792	0.516	0.499	0.116	0.246
	50%	0.061	0.159	0.186	0.323	0.134	0.255	0.132	0.248	0.089	0.218	0.770	0.605	0.093	0.218	0.063	0.173	0.952	0.763	0.519	0.496	0.129	0.260
	Avg	0.052	0.147	0.120	0.253	0.104	0.218	0.093	0.206	0.062	0.177	0.717	0.570	0.071	0.188	0.050	0.154	0.989	0.786	0.516	0.497	0.113	0.254
ETTm2	12.5%	0.026	0.095	0.108	0.239	0.034	0.127	0.062	0.166	0.056	0.159	0.394	0.470	0.133	0.270	0.108	0.228	1.013	0.805	0.307	0.441	0.150	0.298
	25%	0.029	0.100	0.164	0.294	0.042	0.143	0.085	0.196	0.080	0.195	0.421	0.482	0.135	0.272	0.136	0.262	1.039	0.814	0.263	0.402	0.159	0.306
	37.5%	0.031	0.105	0.237	0.356	0.051	0.159	0.106	0.222	0.110	0.231	0.478	0.521	0.155	0.293	0.175	0.300	0.917	0.744	0.250	0.396	0.180	0.321
	50%	0.034	0.111	0.323	0.421	0.059	0.174	0.131	0.247	0.156	0.276	0.568	0.560	0.200	0.333	0.211	0.329	1.140	0.835	0.246	0.389	0.210	0.353
	Avg	0.030	0.103	0.208	0.327	0.046	0.151	0.096	0.208	0.101	0.215	0.465	0.508	0.156	0.292	0.157	0.280	1.027	0.800	0.266	0.407	0.175	0.324
ETTth1	12.5%	0.104	0.211	0.126	0.263	0.240	0.345	0.151	0.267	0.070	0.190	0.857	0.609	0.114	0.234	0.074	0.194	1.265	0.896	0.599	0.554	0.422	0.461
	25%	0.119	0.227	0.169	0.304	0.265	0.364	0.180	0.292	0.106	0.236	0.829	0.672	0.140	0.262	0.102	0.227	1.262	0.883	0.610	0.567	0.412	0.456
	37.5%	0.137	0.243	0.220	0.347	0.296	0.382	0.215	0.318	0.124	0.258	0.830	0.675	0.174	0.293	0.135	0.261	1.200	0.867	0.628	0.577	0.421	0.461
	50%	0.158	0.261	0.293	0.402	0.334	0.404	0.257	0.347	0.165	0.299	0.854	0.691	0.215	0.325	0.179	0.298	1.174	0.849	0.648	0.587	0.443	0.473
	Avg	0.129	0.235	0.202	0.329	0.284	0.373	0.201	0.306	0.117	0.246	0.842	0.682	0.161	0.279	0.122	0.245	1.225	0.873	0.621	0.571	0.424	0.481
ETTth2	12.5%	0.062	0.161	0.187	0.319	0.101	0.231	0.100	0.216	0.095	0.212	0.976	0.754	0.305	0.431	0.163	0.289	2.060	1.120	0.410	0.494	0.521	0.555
	25%	0.063	0.160	0.279	0.390	0.115	0.246	0.127	0.247	0.137	0.258	1.037	0.774	0.322	0.444	0.206	0.331	2.007	1.105	0.419	0.490	0.487	0.535
	37.5%	0.068	0.168	0.400	0.465	0.126	0.257	0.158	0.276	0.187	0.304	1.107	0.800	0.353	0.462	0.252	0.370	2.033	1.111	0.429	0.498	0.487	0.529
	50%	0.074	0.176	0.602	0.572	0.136	0.268	0.183	0.299	0.232	0.341	1.193	0.838	0.369	0.472	0.316	0.419	2.054	1.119	0.467	0.529	0.484	0.523
	Avg	0.067	0.166	0.367	0.436	0.119	0.250	0.142	0.259	0.163	0.279	1.079	0.792	0.337	0.452	0.234	0.352	2.039	1.114	0.431	0.503	0.495	0.475
Electricity	12.5%	0.070	0.188	0.196	0.321	0.102	0.229	0.092	0.214	0.107	0.237	0.297	0.383	0.218	0.326	0.190	0.308	0.277	0.366	0.621	0.620	0.217	0.341
	25%	0.079	0.199	0.207	0.332	0.121	0.252	0.118	0.247	0.120	0.251	0.294	0.380	0.219	0.326	0.197	0.312	0.281	0.369	0.559	0.585	0.219	0.341
	37.5%	0.090	0.201	0.219	0.344	0.141	0.273	0.144	0.276	0.136	0.266	0.296	0.381	0.222	0.328	0.203	0.315	0.275	0.364	0.567	0.588	0.223	0.343
	50%	0.107	0.221	0.235	0.357	0.160	0.293	0.175	0.305	0.158	0.284	0.299	0.383	0.228	0.331	0.210	0.319	0.273	0.361	0.581	0.597	0.229	0.347
	Avg	0.086	0.202	0.214	0.339	0.131	0.262	0.132	0.260	0.130	0.259	0.297	0.382	0.222	0.328	0.200	0.313	0.277	0.365	0.582	0.597	0.222	0.293
Weather	12.5%	0.035	0.090	0.057	0.141	0.047	0.101	0.039	0.084	0.041	0.107	0.140	0.220	0.037	0.093	0.031	0.076	0.296	0.379	0.176	0.287	0.036	0.095
	25%	0.035	0.062	0.065	0.155	0.052	0.111	0.048	0.103	0.064	0.163	0.147	0.229	0.042	0.100	0.035	0.082	0.327	0.409	0.187	0.293	0.042	0.104
	37.5%	0.037	0.065	0.081	0.180	0.058	0.121	0.057	0.117	0.107	0.229	0.156	0.240	0.049	0.111	0.040	0.091	0.406	0.463	0.172	0.281	0.047	0.112
	50%	0.036	0.063	0.102	0.207	0.065	0.133	0.066	0.134	0.183	0.312	0.164	0.249	0.053	0.114	0.046	0.099	0.431	0.483	0.195	0.303	0.054	0.123
	Avg	0.036	0.070	0.076	0.171	0.055	0.117	0.052	0.110	0.099	0.203	0.152	0.235	0.045	0.104	0.038	0.087	0.365	0.434	0.183	0.291	0.045	0.108
Exchange	12.5%	0.004	0.038	0.043	0.151	0.040	0.137	0.023	0.106	0.112	0.224	0.235	0.407	0.345	0.483	0.280	0.411	0.017	0.090	0.025	0.112	0.019	0.098
	25%	0.004	0.039	0.104	0.240	0.043	0.143	0.037	0.134	0.184	0.288	0.412	0.551	0.443	0.561	0.253	0.383	0.030	0.122	0.037	0.136	0.034	0.130
	37.5%	0.005	0.042	0.149	0.296	0.045	0.148	0.053	0.159	0.281	0.355	0.447	0.579	0.393	0.522	0.316	0.472	0.045	0.150	0.050	0.159	0.046	0.153
	50%	0.005	0.046	0.215	0.363	0.039	0.140	0.071	0.185	0.387	0.431	0.499	0.593	0.575	0.649	0.358	0.489	0.036	0.130	0.061	0.176	0.062	0.176
	Avg	0.005	0.041	0.128	0.262	0.042	0.142	0.046	0.146	0.241	0.324	0.398	0.532	0.439	0.554	0.302	0.439	0.032	0.123	0.043	0.146	0.040	0.139
1st Count	25	32	0	0	0	0	0	0	3	1	0	0	0	0	7	2	0	0	0	0	0	0	

B.3 Imputation

Table 11 presents imputation results under different missing rates $r \in \{12.5\%, 25\%, 37.5\%, 50\%\}$. Following the experimental protocol in TimesNet (Wu et al., 2023b), we randomly mask values in the test sequences and evaluate reconstruction quality using MSE and MAE metrics.

B.4 Anomaly Detection

Table 12 reports anomaly detection results on five benchmark datasets using precision (P), recall (R), and F1-score (F1) metrics. Following standard practice (Xu et al., 2021), we apply the point-adjustment protocol where a detection is considered correct if any point within an anomalous segment is identified.

Table 12. Full results for the anomaly detection task. The P, R and F1 represent the precision, recall and F1-score (%) respectively. F1-score is the harmonic mean of precision and recall. A higher value of P, R and F1 indicates a better performance.

Models	SMD			MSL			SMAP			SWaT			PSM			Avg F1
Metric	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	
LSTM (1997) 1997b	78.52	65.47	71.41	78.04	86.22	81.93	91.06	57.49	70.48	78.06	91.72	84.34	69.24	99.53	81.67	77.97
Transformer 2017	83.58	76.13	79.56	71.57	87.37	78.68	89.37	57.12	69.70	68.84	96.53	80.37	62.75	96.56	76.07	76.88
LogTrans 2019	83.46	70.13	76.21	73.05	87.37	79.57	89.15	57.59	69.97	68.67	97.32	80.52	63.06	98.00	76.74	76.60
TCN 2019	84.06	79.07	81.49	75.11	82.44	78.60	86.90	59.23	70.45	76.59	95.71	85.09	54.59	99.77	70.57	77.24
Reformer 2020	82.58	69.24	75.32	85.51	83.31	84.40	90.91	57.44	70.40	72.50	96.53	82.80	59.93	95.38	73.61	77.31
Informer 2021	86.60	77.23	81.65	81.77	86.48	84.06	90.11	57.13	69.92	70.29	96.75	81.43	64.27	96.33	77.10	78.83
Autoformer 2021 88.06	82.35	85.11	77.27	80.92	79.05	90.40	58.62	71.12	89.85	95.81	92.74	99.08	88.15	93.29	84.26	
Pyraformer 2022b	85.61	80.61	83.04	83.81	85.93	84.86	92.54	57.71	71.09	87.92	96.00	91.78	71.67	96.02	82.08	82.57
Anomaly Trans. 2021	88.91	82.23	85.49	79.61	87.37	83.31	91.85	58.11	71.18	72.51	97.32	83.10	68.35	94.72	79.40	80.50
Stationary Trans22c	88.33	81.21	84.62	68.55	89.14	77.50	89.37	59.02	71.09	68.03	96.75	79.88	97.82	96.76	97.29	82.08
LSSL 2022	78.51	65.32	71.31	77.55	88.18	82.53	89.43	53.43	66.90	79.05	93.72	85.76	66.02	92.93	77.20	76.74
ETSformer 2022	87.44	79.23	83.13	85.13	84.93	85.03	92.25	55.75	69.50	90.02	80.36	84.91	99.31	85.28	91.76	82.87
DLinear 2023	83.62	71.52	77.10	84.34	85.42	84.88	92.32	55.41	69.26	80.91	95.30	87.52	98.28	89.26	93.55	82.46
TiDE 2023	76.00	63.00	68.91	84.00	60.00	70.18	88.00	50.00	64.00	98.00	63.00	76.73	93.00	92.00	92.50	74.46
PatchTST 2023b	87.26	82.14	84.62	88.34	70.96	78.70	90.64	55.46	68.82	91.10	80.94	85.72	98.84	93.47	96.08	82.79
iTransformer 2024b	78.45	65.10	71.15	86.15	62.65	72.54	90.67	52.96	66.87	99.96	65.55	79.18	95.65	94.69	95.17	76.98
MAS4TS (Ours)	86.76	83.51	85.10	90.65	74.96	82.06	89.87	55.44	68.58	91.61	94.74	93.14	91.87	98.44	95.04	84.78

C Tool Library and Adaptive Selection

This appendix provides implementation details for the tool library \mathcal{T} introduced in Section 3.4. We describe the unified invocation interface, fundamental building blocks, task-specific tool architectures, and the adaptive routing mechanism. In particular, we clarify (i) the *candidate chain space* considered by the router, (ii) the *search/selection procedure* at inference, and (iii) the *training objective* used for learning the routing policy, so that the tool-driven execution can be reproduced and audited.

C.1 Unified Tool Interface

All tools in MAS4TS expose a standardized invocation interface compatible with the routing policy π_Θ defined in the main text. The router evaluates candidate tool chains and returns the highest-confidence selection. For instance, given strong periodicity signals, it may select [TimesBlock, FFT] over [PatchEmbed, Transformer]; given sparse anchors, it may prepend [ODE_Reconstruct] to the chain.

Semantic-Adaptive Routing. The router maintains a learnable policy network g_ψ that scores each candidate tool chain based on three inputs: the task type τ , structured priors \mathcal{P} extracted by reasoning agents (e.g., detected trends, anchor density, periodicity strength), and a compressed summary of the shared memory state \mathbf{M} .

Candidate chain space. For each task τ , we maintain a registry $\mathcal{C}_\tau = \{\text{Chain}_1, \dots, \text{Chain}_{|\mathcal{C}_\tau|}\}$ of candidate tool chains. Each chain is a short sequence of tool *modules* selected from the task-specific tool set (Sections C.3–C.6) and a small set of shared preprocess/postprocess operators (e.g., normalization and masking). To keep routing efficient and avoid combinatorial explosion, we constrain the chain length to a small integer range:

$$K \in [K_{\min}, K_{\max}], \quad (13)$$

with $K_{\min} = 1$ and $K_{\max} = 3$ in our experiments.¹ Within a task, we *do not* enumerate all possible compositions; instead, \mathcal{C}_τ consists of a curated set of valid, numerically stable chains that respect module I/O schemas (Section C.7).

Search and selection at inference. Given an input context $(\mathbf{X}, \mathcal{P}, \mathbf{M})$ and a task τ , the router computes a score vector over candidate chains:

$$\mathbf{s} = g_\psi(\tau, \mathcal{P}, \text{Pool}(\mathbf{M})) \in \mathbb{R}^{|\mathcal{C}_\tau|}, \quad (14)$$

followed by a softmax to obtain $\pi_\Theta(\text{Chain}_i \mid \tau, \mathcal{P}, \mathbf{M})$ (consistent with Eq. (9) in the main text). We then select either (i) the top-1 chain (greedy) or (ii) a small top- k set for ensembling:

$$\mathcal{I}_k = \text{TopK}(\mathbf{s}, k), \quad k \in \{1, 2, 3\}. \quad (15)$$

Unless otherwise stated, we use greedy top-1 selection when the top confidence is above the threshold (0.6) and switch to top- k ensembling when it is below (Appendix C.7).

C.2 Fundamental Building Blocks

Before describing task-specific tools, we introduce the fundamental components that compose our tool library. These modules are shared across tasks and can be combined in different configurations. Here are some of the key ingredients, not all of them. And we will updated

¹The bounds are chosen to cover common patterns such as “normalize \rightarrow backbone \rightarrow head” while allowing one or two optional modules (e.g., decomposition, ODE reconstruction).

the library in the future.

C.2.1 PATCH EMBEDDING

The patch embedding module converts continuous time series into discrete tokens for transformer-based processing:

$$\mathbf{P} = \text{Linear}(\text{Unfold}(\text{Pad}(\mathbf{X}))) + \mathbf{E}_{pos}, \quad (16)$$

where $\text{Unfold}(\cdot)$ extracts overlapping patches with length p and stride s , and \mathbf{E}_{pos} denotes learnable positional embeddings. Default hyperparameters are `patch_len=16`, `stride=8`, and $d=64$. This module captures local semantic patterns that point-wise tokenization cannot represent, significantly improving temporal pattern learning.

C.2.2 MULTI-HEAD SELF-ATTENTION

Our attention mechanism follows the standard transformer formulation:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \mathbf{V}, \quad (17)$$

with multi-head extension that projects queries, keys, and values into h subspaces, enabling the model to attend to information from different representation spaces.

C.2.3 TIMESBLOCK

For tasks requiring periodic pattern capture (e.g., classification), we employ a specialized block that: (i) detects top- k periods using FFT amplitude analysis, (ii) reshapes the 1D series to 2D representation based on detected periods, (iii) applies 2D Inception-style convolution with multiple kernel sizes, and (iv) aggregates results with learned period weights:

$$\text{Period}(\mathbf{X}) = \text{FFT}^{-1}(\text{TopK}(|\text{FFT}(\mathbf{X})|, k)). \quad (18)$$

C.2.4 NORMALIZATION STRATEGIES

We provide multiple normalization strategies optimized for non-stationary time series:

- **RevIN** (Reversible Instance Normalization): Normalizes each instance independently and maintains statistics for denormalization, critical for preserving the original scale in forecasting.
- **Channel-Independent Normalization**: Normalizes each feature channel separately to avoid spurious cross-channel correlations.
- **Robust Normalization**: Uses percentiles instead of mean/std, providing resistance to outliers in anomaly detection.

C.2.5 PROJECTION HEADS

Task-specific projection heads map latent representations to output space:

- **FlattenHead**: Flattens patch features and projects to output dimension, supporting both channel-independent and shared modes.
- **MLPHead**: Two-layer MLP with configurable activation (ReLU, GELU).
- **ResidualHead**: Includes skip connections for stable gradient flow.

C.3 Forecasting Tools

Table 13 summarizes available forecasting tools. The adaptive router selects among these based on data characteristics inferred from priors \mathcal{P} .

Table 13. Forecasting Tools Overview

Tool	Architecture	Key Components	Selection Condition
ForecastingTool	Patch-Transformer	PatchEmbed + Encoder + FlattenHead	Default; general-purpose
DecompositionTool	Trend-Seasonal	MovingAverage + Projection	Strong trend detected
NBeatsModel	Residual Stacks	Basis Expansion + Doubly Residual	Interpretability required

C.3.1 PATCH-TRANSFORMER FORECASTING TOOL

This is the primary forecasting tool, combining patch embedding with transformer encoding:

Algorithm 1 Patch-Transformer Forecasting

Require: Input $\mathbf{X} \in \mathbb{R}^{B \times L \times D}$, prediction horizon H

- 1: $\mathbf{X}_{norm}, \mu, \sigma \leftarrow \text{RevIN}(\mathbf{X})$ {Instance normalization}
- 2: $\mathbf{P} \leftarrow \text{PatchEmbed}(\mathbf{X}_{norm})$ $\{\mathbb{R}^{B \cdot D \times N_p \times d}\}$
- 3: $\mathbf{H} \leftarrow \text{TransformerEncoder}(\mathbf{P})$ {Multi-layer attention + FFN}
- 4: $\mathbf{Y}_{norm} \leftarrow \text{FlattenHead}(\mathbf{H})$ {Channel-independent projection}
- 5: $\mathbf{Y} \leftarrow \text{RevIN}^{-1}(\mathbf{Y}_{norm}, \mu, \sigma)$ {Denormalize}
- 6: **return** $\mathbf{Y} \in \mathbb{R}^{B \times H \times D}$

Channel Independence: Each feature channel is processed independently, which improves robustness when channels have heterogeneous statistical properties and enables efficient parallel computation.

C.3.2 DECOMPOSITION TOOL

When strong trend is detected (slope > 0.01), the router prepends decomposition:

$$\mathbf{X}_{\text{trend}} = \text{AvgPool}_k(\mathbf{X}), \quad \mathbf{X}_{\text{seasonal}} = \mathbf{X} - \mathbf{X}_{\text{trend}}, \quad (19)$$

where k is the moving average kernel size. The seasonal component is then processed by the main forecasting tool.

C.4 Classification Tools

Table 14. Classification Tools Overview

Tool	Architecture	Key Components	Selection Condition
ClassificationTool	TimesBlock-based	1D Conv + FFT Period + 2D Conv	Default; periodic patterns
TemporalConvNet	Dilated Convolutions	Multi-scale 1D Conv + GlobalPool	Short sequences ($L \leq 32$)
InceptionTime	Multi-branch CNN	Inception Modules + Residual	Multi-scale features needed

C.4.1 TIMESBLOCK CLASSIFICATION TOOL

The primary classification tool leverages FFT-based period discovery:

1. **Data Embedding:** 1D convolution (not linear) to capture local patterns, combined with positional encoding.
2. **TimesBlock Stack:** Multiple blocks that detect periods via FFT, reshape to 2D, apply Inception-style 2D convolution, then reshape back.
3. **Padding Mask:** Critical for variable-length sequences—zeros out padding before flattening.
4. **Full Sequence Projection:** Flattens $[\text{seq_len} \times d]$ and projects to class logits.

C.4.2 TEMPORAL CONVOLUTIONAL NETWORK TOOL

For short sequences, the router selects TCN with dilated causal convolutions:

$$(\mathbf{X} *_r f)(t) = \sum_{i=0}^{k-1} f(i) \cdot \mathbf{X}_{t-r \cdot i}, \quad (20)$$

where r is the dilation factor that increases exponentially with depth, enabling efficient long-range dependency capture with $O(\log L)$ layers.

Table 15. Imputation Tools Overview

Tool	Architecture	Key Components	Selection Condition
ImputationTool	Patch-Transformer	PatchEmbed + Encoder + Reconstruction	Default; general-purpose
SAITS	Self-Attention	Diagonal Masked Attention + ORT	Sparse missing patterns
BRITS	Bidirectional RNN	Forward + Backward GRU + Decay	Sequential missing blocks

C.5 Imputation Tools

C.5.1 PATCH-TRANSFORMER IMPUTATION TOOL

Adapts the patch-based architecture for reconstruction:

Algorithm 2 Patch-Transformer Imputation

Require: Input \mathbf{X} with missing values, binary mask \mathbf{M}

- 1: $\mathbf{X}_{masked} \leftarrow \mathbf{X} \odot (1 - \mathbf{M})$ {Zero out missing positions}
 - 2: $\mathbf{X}_{norm}, \mu, \sigma \leftarrow \text{Normalize}(\mathbf{X}_{masked})$
 - 3: $\mathbf{P} \leftarrow \text{PatchEmbed}(\mathbf{X}_{norm})$
 - 4: $\mathbf{H} \leftarrow \text{TransformerEncoder}(\mathbf{P})$
 - 5: $\hat{\mathbf{X}}_{norm} \leftarrow \text{ReconstructionHead}(\mathbf{H})$ {Same length as input}
 - 6: $\hat{\mathbf{X}} \leftarrow \text{Denormalize}(\hat{\mathbf{X}}_{norm}, \mu, \sigma)$
 - 7: $\mathbf{Y} \leftarrow \mathbf{X} \odot (1 - \mathbf{M}) + \hat{\mathbf{X}} \odot \mathbf{M}$ {Preserve known values}
 - 8: **return** \mathbf{Y}
-

C.6 Anomaly Detection Tools

Table 16. Anomaly Detection Tools Overview

Tool	Architecture	Key Components	Selection Condition
AnomalyDetectionTool	Multi-scale Patch	Multi-scale PatchEmbed + Reconstruction	Default; general anomalies
VAEAnomalyDetector	VAE	Encoder + Reparameterization + Decoder	Point anomaly focus

C.6.1 MULTI-SCALE ANOMALY DETECTION TOOL

The primary anomaly detection tool uses multi-scale reconstruction:

1. **Instance Normalization:** Uses learnable affine parameters to preserve relative anomaly magnitude (unlike RevIN which removes it).
2. **Multi-Scale Processing:** Fine scale ($p = 4$) for point anomalies; coarse scale ($p = 16$) for segment anomalies.
3. **Multi-Scale Fusion:** Concatenates reconstructions from all scales and fuses through MLP.

The anomaly score combines reconstruction error with attention-based weighting:

$$s_t = \|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2 \cdot \alpha_t, \quad (21)$$

where α_t is the attention weight highlighting suspicious regions.

C.6.2 VAE ANOMALY DETECTOR

Probabilistic detection through reconstruction likelihood:

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})), \quad (22)$$

$$s = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \beta \cdot D_{KL}(q_\phi\|p), \quad (23)$$

where the anomaly score combines reconstruction error and KL divergence.

C.7 Tool Chaining

As described in Section 3.4, MAS4TS composes multiple tools into chains where each tool’s output serves as input to the next. A tool chain is formally defined as:

$$\hat{\mathbf{Y}} = \mathcal{T}_{m_K} \circ \mathcal{T}_{m_{K-1}} \circ \cdots \circ \mathcal{T}_{m_1}(\mathbf{Z}), \quad (24)$$

where \circ denotes functional composition and $\{m_1, \dots, m_K\}$ are tool indices selected by routing policy π_Θ .

I/O schema and validity. Each tool \mathcal{T}_m exposes (i) an input schema (tensor shape, mask fields, and optional auxiliary context) and (ii) an output schema compatible with the next tool in the chain. The tool registry enforces schema consistency at composition time; invalid compositions are excluded from \mathcal{C}_τ (Section C.1). This design makes tool chains auditable and prevents silent shape/type mismatches during routing.

Default Chains. Table 17 shows the default tool chains when routing confidence exceeds 0.6. Brackets indicate conditionally included components based on priors \mathcal{P} .

Table 17. Default Tool Chains by Task

Task	Tool Chain
Forecasting	RevIN \rightarrow [Decomposition] \rightarrow PatchTransformer \rightarrow RevIN ⁻¹
Classification	Normalization \rightarrow TimesBlock \rightarrow GlobalPool \rightarrow ClassificationHead
Imputation	MaskEncode \rightarrow Normalization \rightarrow PatchTransformer \rightarrow Reconstruction
Anomaly Detection	InstanceNorm \rightarrow MultiScaleEncoder \rightarrow Reconstruction \rightarrow ScoreCompute

Ensemble Chaining. When routing confidence is low (< 0.6), multiple tool chains can be ensembled:

$$\hat{\mathbf{Y}} = \sum_{i=1}^K w_i \cdot \text{Chain}_i(\mathbf{Z}), \quad w_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)}, \quad (25)$$

where s_i is the confidence score computed by g_ψ for chain i . This provides robustness under ambiguous or out-of-distribution inputs.

Practical note on determinism. To reduce variance during evaluation, we disable stochastic sampling over chains and use deterministic top-1 selection when $c \geq 0.6$, and deterministic top- k ensembling otherwise (Section C.1). This makes routing behavior stable across runs, given a fixed random seed.

D VLM Prompt Templates

This appendix documents the prompt templates used by the **Reasoner** agent for visual reasoning, as introduced in Section 3.3. Each prompt is designed to elicit structured, task-specific outputs from Vision-Language Models that can be directly consumed by downstream numerical reasoning.

D.1 Prompt Design Principles

Our prompts follow three key principles:

1. **Structured Output:** All prompts request JSON-formatted responses with predefined schema, ensuring deterministic parsing and seamless integration with downstream modules.
2. **Statistical Grounding:** Input statistics (min, max, mean, std) and semantic context are provided to anchor visual reasoning in quantitative reality, reducing out-of-range hallucinations.
3. **Task-Specific Anchors:** Different tasks require different anchor types—forecasting anchors define future trajectory shape, anomaly anchors contrast normal vs. abnormal patterns, classification anchors capture discriminative signatures, and imputation anchors guide local reconstruction.

The anchor count range $\{\text{min_anchors}\}$ – $\{\text{max_anchors}\}$ is user-configurable. We empirically suggest: forecasting (8–15), anomaly detection (8–12), classification (5–8), and imputation (5–7), where longer prediction horizons generally benefit from denser anchors.

D.2 Forecasting Prompt

System Prompt (Forecasting)

You are a senior time-series forecasting expert. Analyze the time series plot and provide predictions in strict JSON format. Always include a confidence score (0.00–1.00) based on pattern consistency.

User Prompt Template (Forecasting)

Analyze this time series plot for FORECASTING task (predict $\{\text{pred_len}\}$ future steps).
 Input sequence: $\{\text{seq_len}\}$ steps ($t=0$ to $\{\text{seq_len}-1\}$), last value= $\{\text{last_value}\}$
 Historical stats: range= $[\{\text{history_min}\}, \{\text{history_max}\}]$, std= $\{\text{history_std}\}$
 Requirements:

1. Output ONLY $\{\text{min_anchors}\}$ – $\{\text{max_anchors}\}$ key anchor points (peaks/valleys/inflections) for the PREDICTION WINDOW
2. Anchor types: 'start', 'peak', 'valley', 'inflection', 'end'
3. All values MUST stay within reasonable bounds
4. Confidence score reflects pattern continuity

Output JSON:

```
{
  "confidence": 0.85,
  "anchors": [
    {"t": 96, "v": 0.342, "type": "start"},
    {"t": 120, "v": 0.456, "type": "peak"},
    ...
  ]
}
```

The forecasting prompt explicitly constrains anchors to the prediction window (timesteps $t \geq L$) rather than summarizing historical patterns. The `last_value` field ensures continuity at the forecast boundary, while `history_std` helps calibrate reasonable fluctuation magnitudes. These anchors are subsequently fed to the Latent ODE decoder for dense trajectory reconstruction.

D.3 Classification Prompt

System Prompt (Classification)

You are a senior time-series classification expert. Analyze the time series plot and classify the sequence pattern. Provide strict JSON with confidence score.

User Prompt Template (Classification)

Analyze this time series plot for CLASSIFICATION (input length={seq_len}).
 Input stats: range=[{history_min}, {history_max}], mean={history_mean}
 Requirements:

1. Identify the dominant pattern type (periodic, trending, stationary, etc.)
2. Key anchors: {min_anchors}-{max_anchors} points that define the class signature
3. Confidence score reflects pattern uniqueness

Output JSON:

```
{
  "confidence": 0.78,
  "pattern_type": "periodic",
  "key_anchors": [
    {"t": 12, "v": 0.89, "type": "period_peak"},
    {"t": 36, "v": -0.45, "type": "period_valley"}
  ]
}
```

The prompt focuses on discriminative temporal features that characterize different classes. The `pattern_type` field provides high-level semantic guidance that can be fused with learned numerical features.

D.4 Imputation Prompt

System Prompt (Imputation)

You are a senior time-series imputation expert. Analyze the time series plot and impute MISSING VALUES in the INPUT sequence. Provide JSON with confidence score.

User Prompt Template (Imputation)

Analyze this time series plot for IMPUTATION (input length={seq_len}).
 Missing values: marked at positions {missing_positions}
 Input stats: range=[{history_min}, {history_max}], mean={history_mean}
 Requirements:

1. Impute ONLY marked missing positions
2. Specify reasoning: 'interpolation' or 'extrapolation'
3. Key anchors: {min_anchors}-{max_anchors} critical points guiding imputation

Output JSON:

```
{
  "confidence": 0.88,
  "imputed_values": [
    {"t": 25, "v": 0.34, "reason": "interpolation"},
  ],
  "key_anchors": [
    {"t": 20, "v": 0.28, "type": "observed"},
  ]
}
```


Visual Reasoning over Time Series via Multi-Agent Systems

The imputation prompt explicitly marks missing positions to focus the VLM’s attention on relevant gaps. The `reason` field distinguishes interpolation (bounded by observations on both sides) from extrapolation (extending from one-sided context), informing the Latent ODE about reconstruction confidence. Key anchors at observed boundaries provide hard constraints that ensure consistency with known data points.

D.5 Anomaly Detection Prompt

System Prompt (Anomaly Detection)

You are a senior time-series anomaly detection expert. Analyze the time series plot and identify anomalous time steps in the INPUT sequence. Provide strict JSON with confidence scores.

User Prompt Template (Anomaly Detection)

Analyze this time series plot for ANOMALY DETECTION (input length={seq.len}).
Input stats: range=[{history_min}, {history_max}], mean={history_mean},
std={history_std}
Requirements:

1. Focus ONLY on INPUT SEQUENCE (t=0 to {seq.len-1})
2. Anomaly scores: 0.0-1.0 (higher = more anomalous)
3. For consecutive anomalies, report ONLY the most significant one
4. Key anchors: {min.anchors}-{max.anchors} points including normal patterns as reference

Output JSON:

```
{
  "confidence": 0.82,
  "anomaly_scores": [
    {"t": 45, "score": 0.91, "reason": "spike"},
    {"t": 78, "score": 0.85, "reason": "level_shift"}
  ],
  "key_anchors": [
    {"t": 10, "v": 0.23, "type": "normal"},
    {"t": 45, "v": 1.82, "type": "anomaly"}
  ]
}
```

The anomaly detection prompt requests both anomaly scores and reference anchors. This dual-output structure enables the reconstruction-based detection pipeline to distinguish anomalous patterns from normal temporal dynamics. The `reason` field (spike, level_shift, etc.) provides interpretability for downstream analysis. Requiring only the most significant anomaly among consecutive ones prevents redundant detections.

E Justification of MAS

Although forecasting, imputation, anomaly detection, and classification differ in outputs and objectives, they share a common workflow: (i) **perceive** salient temporal structures, (ii) **reason** about global shape constraints and local irregularities, and (iii) **execute** decisions under validity constraints (e.g., consistency with observations, boundedness, and smoothness). This aligns naturally with an agentic view and motivates modeling time-series analysis as analysis–reasoning–execution pipeline.

A single monolithic model typically entangles these steps into one latent representation with a fixed head, which makes it hard to (a) adapt the execution pathway to heterogeneous series characteristics and task requirements, and (b) incorporate domain tools (e.g., decomposition, robust statistics, constraint projection) in an auditable and controllable manner. A multi-agent design provides a principled way to modularize the workflow into specialized roles while still enabling coordination via a controlled shared-memory interface. Specifically, MAS4TS follows an **Analyzer–Reasoner–Executor (ARE)** paradigm: the *Analyzer* produces grounded statistical priors and a plot-based interface, the *Reasoner* converts perceived structures into explicit anchors and trajectory constraints, and the *Executor* selects, composes, and verifies task-specific tool chains.

This separation supports adaptive computation (different tool compositions for different regimes), improves interpretability (anchors, constraints, and verifier checks are inspectable), and reduces training burden by focusing learning on coordination-critical components (e.g., routing and memory gating) rather than retraining task-specific backbones. In this sense, the benefit of MAS is not merely parallelization, but the ability to unify diverse time-series tasks via a shared reasoning-and-execution protocol with explicit intermediate artifacts.

F Limitations and Future Work

MAS4TS inherits the capabilities and failure modes of its underlying foundation models (especially the VLM used for anchor extraction), which can affect robustness under distribution shift. In addition, for complex multivariate and long-context scenarios, the agent workflow may introduce redundant context and communication overhead, potentially reducing the efficiency gains from tool-driven execution.

We will continue to extend MAS4TS as an evolving time-series agent system by (i) enriching the tool library and improving tool routing for harder real-world constraints, and (ii) exploring agent **skills** and foundation-model reasoning strategies to better handle complex contexts with less redundancy and stronger generalization.