

LLM Agents for Time-Series: A Survey

Anonymous ACL submission

Abstract

Recent advances in large language models (LLMs) have accelerated the development of agentic systems for time-series analysis. While traditional time-series methods typically make predictions or decisions based on a given set of evidence, many real-world applications require agentic systems that can autonomously plan workflows, reflect on intermediate results, and leverage external tools and memory. This survey presents a systematic review of LLM-based agents for time-series tasks. We adopt a problem-driven taxonomy, organizing existing systems into four problem categories: Forecasting & Reasoning, Data Augmentation & Synthesis, Anomaly Detection & Diagnosis, and Decision Support. For each category, we analyze how agent behaviors are implemented through architectural design, external tool integration, and memory mechanisms. We further discuss datasets, environments, and evaluation protocols, and outline open challenges and future directions for LLM-based time-series agents. Overall, our goal is to provide researchers with a structured view of how LLM-based agents can be designed to address different time-series problems.

1 Introduction

Time-series analysis (Hamilton, 2020) plays an important role in many real-world domains, including finance (Tsay, 2005; Taylor, 2008; Franses and Van Dijk, 2000), healthcare (Crabtree et al., 1990; Morid et al., 2023), energy systems (Deb et al., 2017), and transportation (Xu et al., 2016; Ghosh et al., 2009). Representative tasks include forecasting (Chatfield, 2000; De Gooijer and Hyndman, 2006), data augmentation (Wen et al., 2020), anomaly detection (Blázquez-García et al., 2021; Schmidl et al., 2022), and decision support. These tasks have long been addressed with statistical and classical machine learning models such as ARIMA (Shumway and Stoffer, 2017), boot-

strapping (Efron, 1992), and LOF (Breunig et al., 2000), and more recently with deep models such as RNNs (Medsker et al., 2001) and Transformers (Vaswani et al., 2017).

However, applying these methods in practice often relies heavily on expert knowledge to design analysis pipelines and interpret results. Recent advances in large language models (LLMs) (Zhao et al., 2023; Chang et al., 2024) offer an alternative by leveraging general-purpose language reasoning to assist time-series analysis, for example as modules for encoding and explaining temporal patterns. Beyond this prompt-based use, LLMs can also be deployed as *agents* (Wang et al., 2024a) that iteratively plan actions, call external tools, and maintain memory across multiple steps. Such agentic systems are particularly suitable for time-series applications because many time-series tasks require iterative refinement as new observations arrive, adaptive planning under evolving conditions, and tool-augmented integration of heterogeneous evidence rather than a fixed pipeline (Cheng et al., 2026; Tao et al., 2026; ZHANG et al., 2025). Motivated by this emerging paradigm, this survey reviews recent progress on LLM-based agent systems for time-series tasks.

Positioning. Existing LLM-for-time-series surveys (Chang et al., 2025; ZHANG et al., 2025) typically organize methods by individual LLM capabilities (e.g., planning, reasoning, memory, and tool use), rather than by how these capabilities are integrated to solve concrete time-series tasks. Conversely, domain-general LLM-agent surveys (Wang et al., 2024a; Guo et al., 2024; Li et al., 2024; Ferrag et al., 2025) seldom address challenges that are central to time-series settings, including streaming inputs, temporal dependence, semantic drift, and action–data coupling. Moreover, there is still a lack of surveys that focus specifically on LLM agents for time-series problems. To bridge this gap, we adopt a *problem-driven taxonomy* that groups meth-

ods by the tasks they target and then analyzes recurring design patterns in architecture, tool use, and memory mechanisms. A detailed comparison with closely related surveys is provided in Appendix A. **Contributions.** This survey makes four contributions: (i) We propose a problem-driven taxonomy and summarize representative systems in each family. (ii) We analyze recurring agentic patterns and cross-cutting system components (architectures, tools, memory). (iii) We connect task properties (e.g., drift, delayed feedback, operational constraints) to design choices and evaluation protocols. (iv) We curate practical resources and discuss open problems and future directions.

2 Background and Foundations

Time-Series Data and Representations. A time series consists of temporally ordered observations indexed by $\{\tau_t\}_{t=1}^T$ with $\tau_1 < \dots < \tau_T$. We denote a (possibly multivariate) sequence as

$$\mathbf{X}_{1:T} = (\mathbf{x}_1, \dots, \mathbf{x}_T), \quad \mathbf{x}_t \in \mathbb{R}^d.$$

In practice, time-series data may be irregularly sampled, noisy, or partially observed. Systems therefore operate under partial observability and typically construct representations $\mathbf{z}_t = \phi(\mathbf{X}_{1:t})$ that summarize historical context.

Time-series Modeling. Time-series models can be viewed as inference modules that map observed evidence to beliefs or predictions: statistical models (e.g., ARIMA, ARCH/GARCH, HMMs) encode explicit assumptions and interpretable uncertainty (Shumway and Stoffer, 2017; Engle, 1982; Bollerslev, 1986; Rabiner, 1989), deep models (e.g., LSTM, TCN, DeepAR) learn nonlinear mappings via gradient-based optimization (Hochreiter and Schmidhuber, 1997; Bai et al., 2018; Salinas et al., 2020), and LLM-based models (e.g., Time-LLM, LSTPrompt) represent sequences in tokenized or multimodal forms to support reasoning and tool orchestration (Jin et al., 2023; Liu et al., 2024). Despite these differences, most remain one-shot, inference-only systems that estimate future or latent quantities from past data without acting to shape outcomes.

From Inference to Agentic Systems. Static prediction pipelines are often inadequate for interactive, decision-driven time-series tasks that require iterative evidence gathering, feedback integration, and adaptive action selection, which has motivated the emergence of LLM agents. Prior surveys (Wang

et al., 2024a; ZHANG et al., 2025) characterize such systems as observe–act–update loops that combine reasoning, planning, tool use, and memory. We adopt this view and analyze how these capabilities are realized across problem settings.

3 Fundamental Design Dimensions

Before introducing the task taxonomy, we first characterize time-series agent systems along three fundamental design dimensions: architecture, tools, and memory. Together, these dimensions capture how an agent is organized, how it interacts with external resources, and how it retains information across multi-step workflows. This section provides a concise classification of these dimensions to support the task-specific analysis that follows.

3.1 Architecture

We identify the following architectural patterns for time-series agent systems; a visual summary is provided in Appendix B.

Adaptive Single-Agent. The most basic agentic form: one LLM dynamically selects tools and reasoning paths without a fixed workflow, interleaving reasoning and action as evidence arrives (Yao et al., 2023; Schick et al., 2023; Wang et al., 2024a). This distinguishes it from non-agentic LLM use.

Iterative Single-Agent. These systems use explicit self-refine/reflect loops: a single agent generates an output, critiques it, and revises across structured iterations (Madaan et al., 2023; Shinn et al., 2023; Chang et al., 2025). Compared to adaptive agents, the loop is fixed and each step keeps evaluation traces for the next round.

Lifelong / Continual Single-Agent. These systems learn continuously from streaming experience through explicit self-updating mechanisms: a single agent changes over time via training-based adaptation, rather than merely retaining memory, modifying prompts, or retrieving additional context (Parisi et al., 2019; Zheng et al., 2026).

Sequential Pipelines. Multi-stage workflows allow later-stage validation to trigger partial re-execution of earlier stages, blending structure with feedback adaptivity (Chang et al., 2025).

Committee / Debate Systems. Multiple agents generate diverse hypotheses in parallel and aggregate through debate or voting, emphasizing diversity and robustness (Liang et al., 2024).

Planner-Executor Systems. Planner-Executor structures decompose tasks into subtasks and coor-

177 dinare specialized components to solve them, im- 225
178 proving modularity and scalability (Wooldridge, 226
179 2009; Wang et al., 2024a). 227

180 3.2 Tools 228

181 Tools are callable interfaces to external programs 229
182 invoked by the LLM agent (Wang et al., 2024b). In 230
183 time-series agents, they enable access to temporal 231
184 data, specialized computation, and verifiable feed- 232
185 back that cannot be obtained reliably from para- 233
186 metric knowledge alone. Common tool families 234
187 include (i) *Database APIs*, which provide struc- 235
188 tured access to stored data; (ii) *Search & Retrieval* 236
189 *APIs*, which return relevant external information 237
190 or historical records; (iii) *Data Processing*, which 238
191 transform raw inputs into more useful represen- 239
192 tations, such as through alignment or DTW; (iv) 240
193 *Statistical & ML Models*, which produce predic- 241
194 tions, scores, or learned representations; and (v) 242
195 *Simulators, Solvers & Optimizers*, which evaluate 243
196 actions, enforce constraints, or compute solutions 244
197 in structured environments. 245
246
247

198 3.3 Memory 248

199 Memory in time-series agents helps maintain co- 249
200 herence across reasoning steps, especially in tem- 250
201 porally evolving settings (Zhang et al., 2024c). We 251
202 summarize it into three categories: (i) *Evidence* 252
203 *logs*, which record what happened during a de- 253
204 cision process, such as intermediate predictions, 254
205 retrieved evidence, tool outputs, candidate actions, 255
206 and validation results; (ii) *Pattern library*, which 256
207 stores retrievable historical cases or typical data 257
208 fragments, such as recurring market patterns, fault 258
209 signatures, or similar past situations; and (iii) *An-* 259
210 *alytical strategies*, which capture reusable expe- 260
211 rience about how to analyze a situation, such as 261
212 which tools to use, which features to focus on, or 262
213 how to interpret signals before acting. 263
264

214 4 Taxonomy 265

215 We define an *LLM agent for time series* as a sys- 266
216 tem that performs multi-step decision-making over 267
217 time-indexed evidence, where the LLM is centrally 268
218 responsible for action selection, memory updating, 269
219 or hypothesis revision under current observations. 270
220 We include methods for time-series tasks in which 271
221 the LLM selects and sequences the workflow and 272
222 adapts it based on intermediate results or external 273
223 evidence. We exclude (i) one-shot prompting, (ii) 274
224 pipelines in which the LLM serves only as a static 275

encoder or post-hoc explainer, and (iii) domain- 225
general agents that are not designed for time-series 226
constraints or temporally grounded evaluation. 227

228 We adopt a problem-driven taxonomy that 229
230 groups systems by the time-series problems they 231
232 address, and discuss design dimensions within 233
234 each setting. This choice is user-oriented: read- 235
236 ers are often more concerned with what kinds of 237
238 designs are suitable for a given time-series prob- 239
240 lem than with design dimensions in isolation. A 241
242 problem-driven view therefore provides clearer 243
244 guidance for selecting and understanding agent de- 245
246 signs in practice. Under this taxonomy, we identify 247
248 four problem categories: *Forecasting & Reason-*
249 *ing* (forecasting, analytical QA, explanation), *Data*
250 *Augmentation & Synthesis* (data quality, new se-
251 quences/annotations), *Anomaly Detection & Diag-*
252 *nosis* (detect-then-explain), and *Decision Support*
253 (actionable recommendations under constraints, of-
254 ten via simulation/backtesting). Within each cate-
255 gory, we further distinguish several sub-problems.
256 Table 2 summarizes the categories. Figure 1 shows
257 the full hierarchical taxonomy. 258
259
260
261
262
263
264
265

266 4.1 Time-Series Forecasting & Reasoning 267

268 Time-series forecasting and reasoning share a com- 269
270 mon requirement: the agent must produce outputs 271
272 that are not only plausible in language, but also 273
274 grounded in explicit numerical evidence. In both 275
276 settings, a fixed context window or a coarse global 277
278 summary is often insufficient, because the correct 279
280 conclusion may depend on local temporal patterns, 281
282 historical analogs, reusable prototypes, or aligned 283
284 external context such as news (Zhang et al., 2025b). 285
286 As a result, effective systems usually do not treat 287
288 context as a static input. Instead, they actively con- 289
290 struct evidence beyond the context window during 291
292 inference, for example by retrieving relevant slices 293
294 on demand (Jalori et al., 2025) or querying pro- 295
296 totype memories (Jiang et al., 2025). We discuss 297
298 them separately below. 299

300 **Time-Series Forecasting.** Time-series forecast- 301
302 ing agents aim to predict future values from his- 303
304 torical observations under non-stationarity, noise, 305
306 and horizon-dependent uncertainty. In forecasting, 307
308 agent systems are mainly shaped by: (i) multi-stage 309
310 workflows, and (ii) competing hypotheses. 311

312 *Multi-stage workflows* are central in forecasting 313
314 because prediction is rarely a single-step task. A 315
316 full workflow often involves data diagnosis, prepro- 317
318 cessing, model or hyperparameter selection, and 319
320 validation, and errors in early stages can easily 321
322

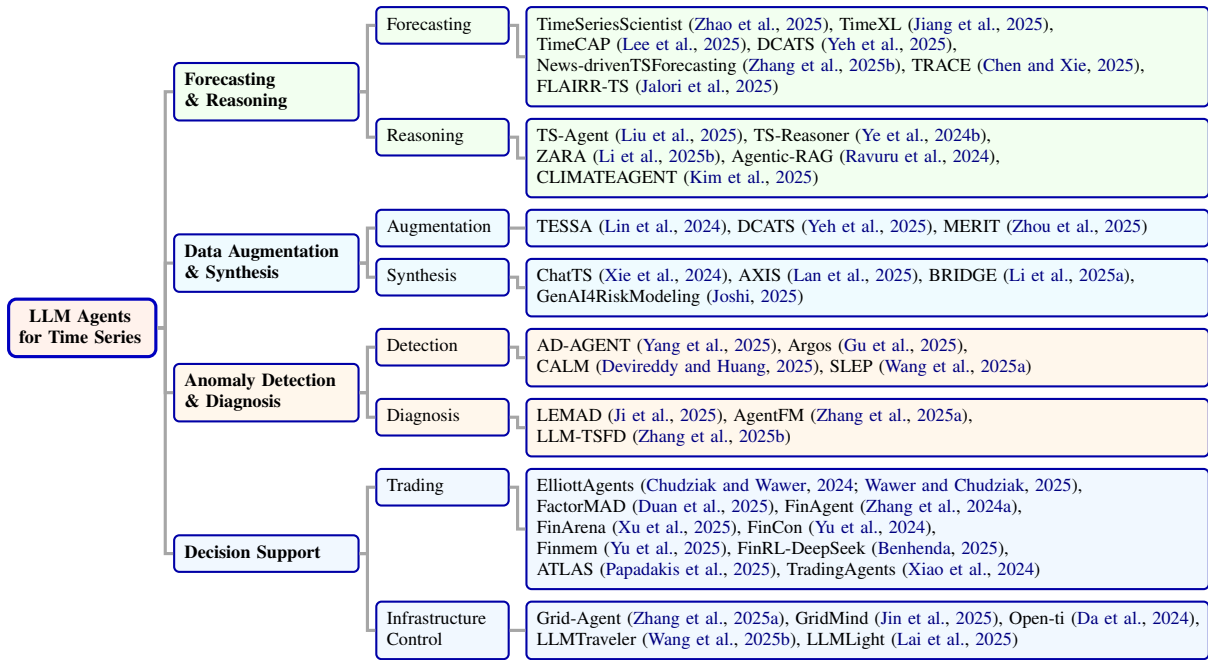


Figure 1: A taxonomy of LLM agents for time-Series.

invalidate later conclusions. A common design is therefore to ground each stage in explicit intermediate evidence rather than free-form reasoning alone. TimeSeriesScientist (Zhao et al., 2025) is a representative example: it uses statistical models and other tools to generate diagnostics, validation results, and configuration records, and stores them as evidence logs for review, rollback, and correction.

Competing hypotheses are another distinctive feature of forecasting. Different strategies may focus on different signals, temporal scales, or exogenous factors, so relying on a single reasoning path can be brittle. A common design is to use committee/debate-style architectures, where different agents represent different forecasting views and are compared through explicit error feedback. NewsTSForecasting (Zhang et al., 2025b) follows this idea by using error-based scoring and acceptance tests to control strategy updates. TRACE (Chen and Xie, 2025) similarly uses communication and multi-agent consistency refinement under sparse or missing observations, while also showing that consistency alone is not enough unless intermediate communication is checked against evidence.

Time-Series Reasoning. Time-series reasoning agents aim to derive explanations, diagnoses, or decisions from temporal evidence rather than directly predict future values. In reasoning, two task-specific considerations are especially important: (i) domain knowledge and numerical support, and (ii)

multi-step reasoning reliability.

Domain knowledge and numerical support are essential in reasoning because text-trained LLMs do not naturally encode temporal structure, domain mechanisms, or reliable quantitative operations. A common design is therefore to separate planning from computation: the main agent handles coordination, while numerical analysis is delegated to specialized sub-agents or auditable tools and operators. Agentic-RAG (Ravuru et al., 2024) illustrates the sub-agent route, while TS-Agent (Liu et al., 2025) illustrates the tool-grounded route. Reasoning systems may also require domain knowledge beyond raw observations. ZARA (Li et al., 2025b) mines discriminative features offline and stores feature-domain relations as reusable guidance, while CLIMATEAGENT (Kim et al., 2025) uses specialized data agents to handle API conventions, metadata retrieval, and parameter validation.

Multi-step reasoning reliability is similar to multi-stage forecasting workflows: reasoning tasks also involve a sequence of dependent steps, and early errors can propagate across the chain. A common design is therefore to make the process explicit, so intermediate results can be checked and revised rather than passed forward implicitly. CLIMATEAGENT (Kim et al., 2025) follows the sub-agent route and records code, retrieved data, and results as evidence logs, so later stages can build on verified outputs. TS-Reasoner (Ye et al.,

2024b) follows the operator route by compiling reasoning into an executable operator pipeline, where execution feedback can trigger plan revision and operator reselection. In both cases, evidence logs support downstream reasoning, reflection, and recovery. Verification may be handled by dedicated checking agents or by the same LLM in a critic role, as in TS-Agent (Liu et al., 2025).

Discussion: Forecasting and reasoning are closely related, as both need evidence beyond a fixed context window and can suffer from error accumulation in long workflows. These shared challenges motivate similar memory mechanisms: patterns library helps retrieve relevant precedents, while evidence logs preserve intermediate steps for reflection and revision. Their tool use is also broadly aligned, with heavy reliance on data processing and frequent calls to operators or statistical/ML models for quantitative support. The main difference lies in architecture: forecasting more often compares competing hypotheses through committee/debate designs, whereas reasoning more often relies on iterative single-agent loops to reduce long-chain errors.

4.2 Time-Series Augmentation & Synthesis

Time-series augmentation and synthesis both aim to construct additional data while preserving time-series semantics. In both settings, the main failure mode is semantic drift: the constructed data may look plausible but violate the structures that matter for downstream tasks, such as trend, periodicity, local shapes, or correlations. The key difference is semantic source: augmentation mainly relies on the target series, whereas synthesis must align control text (e.g., scenario descriptions) with domain rules. We compare them separately below.

Time-Series Augmentation. Time-series augmentation aims to construct additional data around a target sequence while preserving task-relevant semantics. It typically appears in two forms: generating numeric perturbations and generating textual annotations as an alternative representation of the same series. These two forms are shaped by different challenges: (i) *semantic preservation* is central for numeric augmentation, while (ii) *domain annotation understanding* is the main bottleneck for textual augmentation.

Semantic preservation is the core challenge in numeric augmentation. Augmented data should

remain aligned with the target series rather than merely look realistic in isolation. A common design follows one of two routes: either construct target-specific training sets from neighboring series, as in DCATS (Yeh et al., 2025), or retrieve similar sequences and then apply classic augmentations such as jittering, scaling, and time warping, as in MERIT (Zhou et al., 2025). Verification is especially important here. LLM-as-Judge based only on pretrained knowledge can be brittle, so stronger designs usually rely on downstream validation signals (Yeh et al., 2025).

Domain annotation understanding is the main challenge in textual augmentation. In series-to-text semantic augmentation, the added data is not a new numeric sequence but a textual annotation, which can be viewed as another representation or semantic view of the same series. The difficulty is that pretrained LLMs do not reliably understand domain-specific time-series annotations. A common remedy is to first learn general semantic descriptions from cross-domain annotations by decontextualizing them into domain-agnostic concepts (e.g., trend, periodicity, volatility), and then convert the general semantic annotations into domain-specific annotations via a domain agent (Lin et al., 2024).

Time-Series Synthesis. Time-series synthesis aims to generate new sequences under specified controls or constraints rather than to expand a single target series. Under this view, it mainly includes two settings: text-guided synthesis, where the system generates sequences that satisfy both scenario descriptions and domain knowledge, and domain-data construction without any text guidance. These two settings face different bottlenecks: (i) *text-series alignment* is central in the former, while (ii) *Domain constraint compliance* is the main challenge in the latter.

Text-series alignment is the core challenge in text-guided synthesis. In real-world settings, paired scenario-query and time-series data are usually scarce, so LLMs cannot reliably map free-form text directly to numerical sequences. A common design is therefore to rewrite descriptions of real time series into structured text queries and then map these queries into executable generation settings. BRIDGE (Li et al., 2025a) extracts templates (e.g., length, trend, periodicity, extrema, variance), refines them, and trains a controlled diffusion generator. GenAI4RiskModeling (Joshi, 2025) similarly rewrites narratives into structured scenario-

query forms and lets the agent translate them into hyperparameter settings for GAN/VAE-based generation.

Domain constraint compliance is the main challenge in domain-data construction without text control. In this setting, the difficulty is not text alignment but ensuring that generated data still conforms to implicit domain rules and constraints. A common design is therefore to let the agent first choose a set of controllable domain attributes, then translate them into executable generation rules, and finally filter inconsistent outputs through downstream verification, keeping only aligned supervision signals (Xie et al., 2024).

Discussion: For augmentation and synthesis, the central issue is that generated outputs may look plausible while drifting away from task-relevant semantics. As a result, current systems are often organized around explicit construction–verification pipelines rather than unconstrained LLM generation. This emphasis is also reflected in tool use, which centers on data processing, augmentation or generation modules, and downstream validation signals. In existing work, memory is less central here than in forecasting or reasoning. And when it appears, it mainly appears as reusable templates or semantic patterns with domain rules.

4.3 Time-Series Anomaly Detection & Diagnosis

Time-series anomaly detection and diagnosis are closely related: detection asks whether and when abnormal patterns occur, while diagnosis asks why and where they occur and how to respond. Their main difference lies in task emphasis, and we discuss them separately below.

Time-Series Anomaly Detection. Under this survey’s definition, time-series anomaly detection focuses on producing reliable alarms, such as point-wise labels, anomalous windows, or alert events. In practice, three challenges are especially important: (i) *evidence quality*, (ii) *non-degradation* relative to the base detector, and (iii) *streaming monitoring* under continual updates.

Evidence quality matters because reliable alarms often require more than raw series values alone. Here, decision evidence refers to the information directly supporting the alarm decision, such as summary statistics, learned features, or retrieved histor-

ical snippets. A common design is to strengthen this evidence by injecting it into prompts or downstream modules. SLEP (Wang et al., 2025a) follows this direction by enriching detector inputs with additional evidence rather than relying only on the raw sequence.

Non-degradation is critical because the agent is often layered on top of an existing deployed detector, so it should improve alarm quality without underperforming the base system. A common design is to treat the deployed detector as a base model and let the agent learn complementary corrections for its typical errors. ARGOS (Gu et al., 2025) exemplifies this idea by constructing deterministic rules to correct base-detector failures and fusing rule outputs with detector outputs at inference time.

Streaming monitoring is a special but practically important setting because concept drift may require continual adaptation. The main risk is contamination control: short-lived anomalies may be absorbed as the new normal during online updates. CALM (Devireddy and Huang, 2025) addresses this with a continual design in which an LLM agent filters training data by distinguishing transient noise from sustained distribution shift, and only the latter is used for short-horizon adaptation of the forecasting-based detector.

Time-Series Diagnosis. Under this survey’s definition, time-series diagnosis focuses on explaining, localizing, and responding to abnormal events. In practice, two challenges are especially important: (i) *evidence integration*, and (ii) *limited diagnostic supervision*.

Evidence integration is important for diagnosis for reasons similar to anomaly detection, but diagnosis places more emphasis on explanation and localization. In the papers we survey, this is often handled by explicitly incorporating textual evidence, such as logs, alerts, traces, and topology context, as part of the model input or generated report. AgentFM (Zhang et al., 2025a) follows this pattern and further improves stability by retrieving labeled historical faults through RAG as few-shot references, which guides the prompt and reduces free-form drift.

Limited diagnostic supervision is another recurring bottleneck because labeled fault cases and high-quality incident narratives are often scarce. LLM-TSFD (Zhang et al., 2025b) addresses this by introducing a human-in-the-loop data preparation stage, where users specify labeling or cleaning intent in natural language and the system gener-

ates executable processing code that is iteratively refined through feedback.

Discussion: Anomaly detection and diagnosis are both concerned with building reliable and actionable monitoring pipelines from heterogeneous temporal evidence. This makes tool support especially important: both settings rely heavily on data processing, existing detectors or analytical modules, and external evidence sources. Sequential pipelines therefore remain the dominant architectural pattern, although streaming monitoring scenarios further motivate continual agent systems that can update over time. Memory serves a supporting but important role: evidence logs help keep intermediate artifacts auditable, while pattern-library memory allows retrieval of historical incidents or labeled examples to stabilize decisions.

4.4 Time-Series Decision Making

Time-series decision making differs from forecasting or diagnosis because the output is an executable action, plan, or allocation rather than a descriptive judgment. Decisions must be grounded in a predefined action space and satisfy domain constraints. The key cross-domain difference is constraint type: trading usually tolerates some risk, whereas infrastructure control (e.g., traffic or power systems) often impose hard feasibility and safety constraints. We therefore examine them separately below.

Trading. Trading agents focus on sequential financial decisions (e.g., buy, sell, hold) under evolving market conditions. Three considerations are central: (i) heterogeneous external evidence, (ii) historical experience, and (iii) risk preference.

Heterogeneous external evidence, such as news, financial statements, social media, earnings calls, visual charts, and technical indicators, often affects trading decisions, but it comes in different forms and operates at different temporal scales. Processing all of it with a single agent can easily lead to context overload and mixed signals. A common design is a planner-executor architecture, where specialized subagents or modules handle different evidence sources and their outputs are then aggregated. TradingAgents (Xiao et al., 2024), FINCON (Yu et al., 2024), and FinArena (Xu et al., 2025) all follow this pattern.

Trading decisions often rely heavily on *historical experience*. In agent systems, this is usually

supported by memory in the form of *pattern library* and *analytical strategies*, which abstract past situations into retrievable patterns or reusable decision experience. FinAgent (Zhang et al., 2024a) exemplifies this design: each stage produces a retrieval-oriented query, allowing market situations, price-driving explanations, and trading lessons to be stored separately. FINCON (Yu et al., 2024) further updates *manager-level investment beliefs*, which serve as evolving analytical strategies.

Risk preference is another special concern in trading. In some settings, the system needs to account for human preference alignment; a common design is to inject user risk preferences and feedback into prompts in a human-in-the-loop manner, so that they directly influence the final recommendation (Xu et al., 2025). In other settings, the system does not explicitly incorporate user feedback, but instead constructs internal role-based variation to induce different risk styles, again typically through prompt injection (Xiao et al., 2024; Yu et al., 2025).

Infrastructure Control. These agents mainly target control tasks in infrastructure systems, especially traffic and power systems. Compared with trading, infrastructure control is shaped mainly by two features: (i) hard constraints, and (ii) human-specified policies.

Hard constraints are central in infrastructure control because decisions usually operate in strict environments, where feasibility and safety must be validated by the environment or external tools. As a result, simulators or solvers are often essential as verification tools. A common design is a tool-grounded sequential pipeline, in which the LLM proposes candidate actions or plans, while simulators determine feasibility and effect. GridAgent (Zhang et al., 2025a) is a representative example: the LLM proposes structured mitigation plans, while power-flow solvers and validation modules determine their performance. Similarly, LLMLight (Lai et al., 2025) operates in a fixed control space and is evaluated in a traffic simulator.

Human-specified policies are another distinctive feature of infrastructure control. In some systems, humans specify policies or optimization settings, and the LLM agent mainly orchestrates downstream execution. Open-TI (Da et al., 2024) illustrates this pattern well: some tasks translate human-described policies into signal actions, while others let the user specify simulation settings or optimization techniques and have the LLM route the request to the appropriate tool chain.

Discussion: Different decision targets lead to different design principles for trading and infrastructure-control agents. Trading systems more often use planner-executor designs to decompose multimodal evidence, whereas infrastructure-control systems more often rely on sequential execution-validation pipelines. The same contrast appears in memory and tool use: trading emphasizes historical experience (pattern libraries, analytical strategies) and heterogeneous-data tooling, while infrastructure control depends more on simulators or solvers for feasibility and safety checks.

5 Resources for Implementation and Evaluation

This section summarizes key resources for implementing and evaluating LLM agents for time-series tasks. We organize them into four categories and provide representative examples (Table 3).

Datasets and Repositories. These resources provide the raw observations used for training and offline evaluation. Representative forecasting datasets include Electricity (Lai et al., 2018), METR-LA (Li et al., 2017), and ETT (Zhou et al., 2021), while common anomaly datasets include SWaT (Mathur and Tippenhauer, 2016) and SMAP/MSL (Hundman et al., 2018). Monash TSF (Godahewa et al., 2021) serves as a standardized multi-domain repository.

Benchmarks. Benchmarks define comparable tasks and metrics across methods. Forecasting benchmarks include M4/M5 (Makridakis et al., 2018, 2022) and MIRAI (Ye et al., 2024a), while anomaly suites include NAB (Lavin and Ahmad, 2015) and TSB-UAD (Paparrizos et al., 2022). TimeSeriesExam (Cai et al., 2024) extends evaluation toward reasoning.

Interactive Environments. These platforms enable closed-loop experiments with explicit state, action, and reward signals. Typical examples include Grid2Op (Marot et al., 2021) for power systems, SUMO (Behrisch et al., 2011) for traffic control, and SocioDojo (Cheng and Chin, 2024) for trading.

Toolkits. Toolkits provide reusable pipelines, APIs, and baselines for reproducible development. Common options include StatsForecast (Garza et al., 2022), NeuralForecast (Challu et al., 2023), Sktime (Löning et al., 2019), and Stable-Baselines3 (Raffin et al., 2021).

6 Conclusion and Future Work

This survey presents the first comprehensive review of LLM-based agentic systems for time-series problems. We propose a novel problem-driven taxonomy that organizes existing methods by the problems they address and highlights recurring design patterns across tasks. While this taxonomy provides a practical lens for understanding and designing time-series agents, several important challenges remain open for future research.

Trustworthiness and Validation. Although recent agent systems improve flexibility and efficiency, their trustworthiness remains limited (Jin et al., 2024; Raza et al., 2025). In many time-series settings, reliable validation methods such as simulation and backtesting are unavailable, so evaluation still relies on LLM-as-a-judge or weak heuristic checks (Tomašević et al., 2025). More reliable validation methods are therefore needed.

Numerical Understanding and Domain Knowledge. LLMs are general language models pretrained mainly on text, and remain limited in understanding numerical signals and specialized domain knowledge (Hung et al., 2023; Ye et al., 2024b). Many current agent designs implicitly assume that pretrained reasoning can transfer directly to time-series settings. Yet even when external tools provide quantitative results, agents still need the ability to interpret them, judge their reliability, and connect them to domain-specific reasoning. Improving this capability remains an important direction.

Online Adaptation and Continual Improvement. In real-world applications such as anomaly detection and decision making, agents often need to improve continuously after deployment. Yet most current systems remain largely static and rely mainly on memory updates, while parameter adaptation remains rare (Jaglan and Barnes, 2025; Zheng et al., 2026). A key next step is to enable stronger online adaptation in deployed agents.

Benchmarking Agent Workflows. Current studies still evaluate time-series agents mainly with standard task metrics. This leaves an important question: how should we evaluate agent workflows themselves, beyond downstream results—for example, in terms of the quality of intermediate decisions, the appropriateness of tool use, and the effectiveness of reflection? Recent work has begun to expose this gap (Weng et al., 2026; Cheng et al., 2026). Future work should therefore benchmark agent workflows more systematically.

7 Limitations

While this survey provides a comprehensive overview of LLM-based agentic systems for time-series tasks, it has several limitations:

Scope and coverage. Due to the rapid pace of advancements in LLM-based agents, some recent developments and emerging directions may not be fully captured in this survey.

Lack of quantitative comparison. The broad range of time-series tasks and heterogeneous evaluation settings make it difficult to establish a unified and fair empirical comparison across all systems.

Design guidance. The design insights are derived from recurring patterns observed in the literature rather than controlled experimental validation, and thus may not generalize to all practical settings.

Despite these limitations, we hope this survey provides a useful and structured reference for understanding and designing LLM-based time-series agents.

References

Sridhar Adepu, Venkata Reddy Palleti, Gyanendra Mishra, and Aditya Mathur. 2020. Investigation of cyber attacks on a water distribution system. In *International Conference on Applied Cryptography and Network Security*, pages 274–291. Springer.

Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Danielle C Maddix, Syama Rangapuram, David Salinas, Jasper Schulz, and 1 others. 2020. Gluonts: Probabilistic and neural time series modeling in python. *Journal of Machine Learning Research*, 21(116):1–6.

Shaojie Bai, J Zico Kolter, and Vladlen Koltun. 2018. [An empirical evaluation of generic convolutional and recurrent networks for sequence modeling](#). *CoRR*, abs/1803.01271.

Michael Behrisch, Laura Bieker, Jakob Erdmann, and Daniel Krajzewicz. 2011. Sumo—simulation of urban mobility: an overview. In *Proceedings of SIMUL 2011, the third international conference on advances in system simulation*. ThinkMind.

Mostapha Benhenda. 2025. Finrl-deepseek: Llm-infused risk-sensitive reinforcement learning for trading agents. *arXiv preprint arXiv:2502.07393*.

Aadyot Bhatnagar, Paul Kassianik, Chenghao Liu, Tian Lan, Wenzhuo Yang, Rowan Cassius, Doyen Sahoo, Devansh Arpit, Sri Subramanian, Gerald Woo, and 1 others. 2021. Merlion: A machine learning library for time series. *arXiv preprint arXiv:2109.09265*.

Ane Blázquez-García, Angel Conde, Usue Mori, and Jose A Lozano. 2021. A review on outlier/anomaly detection in time series data. *ACM computing surveys (CSUR)*, 54(3):1–33.

Tim Bollerslev. 1986. [Generalized autoregressive conditional heteroskedasticity](#). *Journal of Econometrics*, 31(3):307–327.

Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. Lof: identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 93–104.

Yifu Cai, Arjun Choudhry, Mononito Goswami, and Artur Dubrawski. 2024. Timeseriesexam: A time series understanding exam. *arXiv preprint arXiv:2410.14752*.

Yifu Cai, Xinyu Li, Mononito Goswami, Michał Wil-
iński, Gus Welter, and Artur Dubrawski. 2025. Time-seriesgym: A scalable benchmark for (time series) machine learning engineering agents. *arXiv preprint arXiv:2505.13291*.

Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza Ramirez, Max Mergenthaler Canseco, and Artur Dubrawski. 2023. Nhits: Neural hierarchical interpolation for time series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 6989–6997.

Ching Chang, Yidan Shi, Defu Cao, Wei Yang, Jeehyun Hwang, Haixin Wang, Jiacheng Pang, Wei Wang, Yan Liu, Wen-Chih Peng, and 1 others. 2025. A survey of reasoning and agentic systems in time series with large language models. *arXiv preprint arXiv:2509.11575*.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, and 1 others. 2024. A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology*, 15(3):1–45.

Chris Chatfield. 2000. *Time-series forecasting*. Chapman and Hall/CRC.

Yuxuan Chen and Haipeng Xie. 2025. Trace: Unlocking the potential of llms in time series forecasting for distributed energy resources. *IEEE Transactions on Artificial Intelligence*.

Junyan Cheng and Peter Chin. 2024. Sociodojo: Building lifelong analytical agents with real-world text and time series. In *The Twelfth International Conference on Learning Representations*.

Mingyue Cheng, Xiaoyu Tao, Qi Liu, Ze Guo, and Enhong Chen. 2026. Position: Beyond model-centric prediction – agentic time series forecasting. *arXiv preprint arXiv:2602.01776*.

787	Jaroslaw A Chudziak and Michal Wawer. 2024. Elliotta- agents: A natural language-driven multi-agent system for stock market analysis and prediction. In <i>Proceed- ings of the 38th Pacific Asia Conference on Language, Information and Computation</i> , pages 961–970.	Federico Garza, Max Mergenthaler Canseco, Cristian Challú, and Kin G Olivares. 2022. Statsforecast: Lightning fast forecasting with statistical and econo- metric models. <i>PyCon Salt Lake City, Utah, US</i> , 2022:6.	839 840 841 842 843
792	Benjamin F Crabtree, Subhash C Ray, Priscilla M Schmidt, Patrick T O’Connor, and David D Schmidt. 1990. The individual over time: time series applica- tions in health care research. <i>Journal of clinical epidemiology</i> , 43(3):241–260.	Bidisha Ghosh, Biswajit Basu, and Margaret O’Mahony. 2009. Multivariate short-term traffic flow forecasting using time-series analysis. <i>IEEE transactions on intelligent transportation systems</i> , 10(2):246–254.	844 845 846 847
797	Longchao Da, Kuanru Liou, Tiejin Chen, Xuesong Zhou, Xiangyong Luo, Yezhou Yang, and Hua Wei. 2024. Open-ti: Open traffic intelligence with aug- mented language model. <i>International Journal of Machine Learning and Cybernetics</i> , 15(10):4761– 4786.	Rakshitha Godahewa, Christoph Bergmeir, Geoffrey I Webb, Rob J Hyndman, and Pablo Montero-Manso. 2021. Monash time series forecasting archive. <i>arXiv preprint arXiv:2105.06643</i> .	848 849 850 851
803	Jan G De Gooijer and Rob J Hyndman. 2006. 25 years of time series forecasting. <i>International journal of forecasting</i> , 22(3):443–473.	Yile Gu, Yifan Xiong, Jonathan Mace, Yuting Jiang, Yigong Hu, Baris Kasikci, and Peng Cheng. 2025. Argos: agentic time-series anomaly detection with autonomous rule generation via large language mod- els (2025). <i>arXiv preprint arXiv:2501.14170</i> .	852 853 854 855 856
806	Chirag Deb, Fan Zhang, Junjing Yang, Siew Eang Lee, and Kwok Wei Shah. 2017. A review on time series forecasting techniques for building energy consump- tion. <i>Renewable and Sustainable Energy Reviews</i> , 74:902–924.	Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez, Konrad Zolna, Rishabh Agarwal, Josh S Merel, Daniel J Mankowitz, Cosmin Paduraru, and 1 others. 2020. RL unplugged: A suite of benchmarks for offline reinforcement learn- ing. <i>Advances in neural information processing sys- tems</i> , 33:7248–7259.	857 858 859 860 861 862 863
811	Ashok Devireddy and Shunping Huang. 2025. Calm: A framework for continuous, adaptive, and llm- mediated anomaly detection in time-series streams. <i>arXiv preprint arXiv:2508.21273</i> .	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi- angliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. <i>arXiv preprint arXiv:2402.01680</i> .	864 865 866 867 868
815	Yitong Duan, chuheng zhang, and Jian Li. 2025. Factor- mad: A multi-agent debate framework based on large language models for interpretable stock alpha factor mining. In <i>Proceedings of the 6th ACM International Conference on AI in Finance</i> , pages 605–613.	James D Hamilton. 2020. <i>Time series analysis</i> . Prince- ton university press.	869 870
820	Bradley Efron. 1992. Bootstrap methods: another look at the jackknife. In <i>Breakthroughs in statis- tics: Methodology and distribution</i> , pages 569–593. Springer.	Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. 2022. Adbench: Anomaly de- tection benchmark. <i>Advances in neural information processing systems</i> , 35:32142–32159.	871 872 873 874
824	Robert F Engle. 1982. Autoregressive conditional het- eroscedasticity with estimates of the variance of united kingdom inflation. <i>Econometrica</i> , 50(4):987– 1008.	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory . <i>Neural Computation</i> , 9(8):1735– 1780.	875 876 877
828	Mohamed Amine Ferrag, Norbert Tihanyi, and Mer- ouane Debbah. 2025. From llm reasoning to au- tonomous ai agents: A comprehensive review. <i>arXiv preprint arXiv:2504.19678</i> .	Kyle Hundman, Valentino Constantinou, Christopher Laporte, Ian Colwell, and Tom Soderstrom. 2018. Detecting spacecraft anomalies using lstms and non- parametric dynamic thresholding. In <i>Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining</i> , pages 387–395.	878 879 880 881 882 883
832	Philip Hans Franses and Dick Van Dijk. 2000. <i>Non- linear time series models in empirical finance</i> . Cam- bridge university press.	Chia-Chien Hung, Wiem Ben Rim, Lindsay Frost, Lars Bruckner, and Carolin Lawrence. 2023. Walking a tightrope—evaluating large language models in high- risk domains. In <i>Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation in NLP</i> , pages 99–111.	884 885 886 887 888 889
835	Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. 2020. D4rl: Datasets for deep data-driven reinforcement learning. <i>arXiv preprint arXiv:2004.07219</i> .	Aman Jaglan and Jarrod Barnes. 2025. Continual learn- ing, not training: Online adaptation for agents. <i>arXiv preprint arXiv:2511.01093</i> .	890 891 892

893	Gunjan Jalori, Preetika Verma, and Sercan	<i>The 41st international ACM SIGIR conference on re-</i>	949
894	Ö Ar	<i>search & development in information retrieval</i> , pages	950
895	ik. 2025. Flairr-ts—forecasting llm-agents with iter-	95–104.	951
896	ative refinement and retrieval for time series. <i>arXiv</i>		
897	<i>preprint arXiv:2508.19279</i> .		
898	Xin Ji, Le Zhang, Wenya Zhang, Fang Peng, Yi-	Siqi Lai, Zhao Xu, Weijia Zhang, Hao Liu, and Hui	952
899	fan Mao, Xingchuang Liao, and Kui Zhang. 2025.	Xiong. 2025. Llmight: Large language models as	953
900	Lemad: Llm-empowered multi-agent system for	traffic signal control agents. In <i>Proceedings of the</i>	954
901	anomaly detection in power grid services. <i>Electron-</i>	<i>31st ACM SIGKDD Conference on Knowledge Dis-</i>	955
902	<i>ics</i> , 14(15):3008.	<i>covery and Data Mining V. 1</i> , pages 2335–2346.	956
903	Yushan Jiang, Zijie Pan, Xikun Zhang, Sahil Garg, An-	Tian Lan, Hao Duong Le, Jinbo Li, Wenjun He, Meng	957
904	derson Schneider, Yuriy Nevmyvaka, and Dongjin	Wang, Chenghao Liu, and Chen Zhang. 2025. Axis:	958
905	Song. 2024. Empowering time series analysis with	Explainable time series anomaly detection with large	959
906	large language models: A survey. <i>arXiv preprint</i>	language models. <i>arXiv preprint arXiv:2509.24378</i> .	960
907	<i>arXiv:2402.03182</i> .		
908	Yushan Jiang, Wenchao Yu, Geon Lee, Dongjin Song,	Alexander Lavin and Subutai Ahmad. 2015. Evaluating	961
909	Kijung Shin, Wei Cheng, Yanchi Liu, and Haifeng	real-time anomaly detection algorithms—the numenta	962
910	Chen. 2025. Timexl: Explainable multi-modal time	anomaly benchmark. In <i>2015 IEEE 14th interna-</i>	963
911	series prediction with llm-in-the-loop. In <i>The Thirty-</i>	<i>national conference on machine learning and applica-</i>	964
912	<i>ninth Annual Conference on Neural Information Pro-</i>	<i>tions (ICMLA)</i> , pages 38–44. IEEE.	965
913	<i>cessing Systems</i> .		
914	Hongwei Jin, Kibaek Kim, and Jonghwan Kwon. 2025.	Geon Lee, Wenchao Yu, Kijung Shin, Wei Cheng, and	966
915	Gridmind: Llms-powered agents for power system	Haifeng Chen. 2025. Timecap: Learning to con-	967
916	analysis and operations. In <i>Proceedings of the SC'25</i>	textualize, augment, and predict time series events	968
917	<i>Workshops of the International Conference for High</i>	with large language model agents. In <i>Proceedings</i>	969
918	<i>Performance Computing, Networking, Storage and</i>	<i>of the AAAI Conference on Artificial Intelligence</i> ,	970
919	<i>Analysis</i> , pages 560–568.	volume 39, pages 18082–18090.	971
920	Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu,	Hao Li, Yuhao Huang, Chang Xu, Viktor Schlegel,	972
921	James Y. Zhang, Xiaoming Shi, Pin-Yu Chen, Yux-	Renhe Jiang, Riza Batista-Navarro, Goran Nenadic,	973
922	uan Liang, Yifan Li, Shirui Pan, and Qingsong Wen.	and Jiang Bian. 2025a. Bridge: Bootstrapping text	974
923	2023. Time-llm: Time series forecasting by repro-	to control time-series generation via multi-agent it-	975
924	gramming large language models . <i>arXiv preprint</i>	erative optimization and diffusion modelling. <i>arXiv</i>	976
925	<i>arXiv:2310.01728</i> . Accepted at ICLR 2024.	<i>preprint arXiv:2503.02445</i> .	977
926	Ming Jin, Yifan Zhang, Wei Chen, Kexin Zhang, Yux-	Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang.	978
927	uan Liang, Bin Yang, Jindong Wang, Shirui Pan,	2024. A survey on llm-based multi-agent sys-	979
928	and Qingsong Wen. 2024. Position: What can large	tems: workflow, infrastructure, and challenges. <i>Vici-</i>	980
929	language models tell us about time series analysis.	<i>nagearth</i> , 1(1):9.	981
930	In <i>Forty-first International Conference on Machine</i>	Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu.	982
931	<i>Learning</i> .	2017. Diffusion convolutional recurrent neural net-	983
932	Satyadhar Joshi. 2025. Using gen ai agents with gae and	work: Data-driven traffic forecasting. <i>arXiv preprint</i>	984
933	vae to enhance resilience of us markets. <i>Available at</i>	<i>arXiv:1707.01926</i> .	985
934	<i>SSRN 5123068</i> .		
935	Hyeonjae Kim, Chenyue Li, Wen Deng, Mengxi Jin,	Zechen Li, Baiyu Chen, Hao Xue, and Flora D Salim.	986
936	Wen Huang, Mengqian Lu, and Binhang Yuan. 2025.	2025b. Zara: Zero-shot motion time-series analysis	987
937	Climateagent: Multi-agent orchestration for com-	via knowledge and retrieval driven llm agents. <i>arXiv</i>	988
938	plex climate data science workflows. <i>arXiv preprint</i>	<i>preprint arXiv:2508.04038</i> .	989
939	<i>arXiv:2511.20109</i> .		
940	Kenneth R Knapp, Michael C Kruk, David H Levin-	Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang,	990
941	son, Howard J Diamond, and Charles J Neumann.	Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and	991
942	2010. The international best track archive for cli-	Zhaopeng Tu. 2024. Encouraging divergent thinking	992
943	mate stewardship (ibtracs) unifying tropical cyclone	in large language models through multi-agent debate .	993
944	data. <i>Bulletin of the American Meteorological Soci-</i>	<i>Preprint</i> , arXiv:2305.19118.	994
945	<i>ety</i> , 91(3):363–376.		
946	Guokun Lai, Wei-Cheng Chang, Yiming Yang, and	Minhua Lin, Zhengzhang Chen, Yanchi Liu, Xujiang	995
947	Hanxiao Liu. 2018. Modeling long-and short-term	Zhao, Zongyu Wu, Junxiang Wang, Xiang Zhang,	996
948	temporal patterns with deep neural networks. In	Suhang Wang, and Haifeng Chen. 2024. Decod-	997
		ing time series with llms: A multi-agent frame-	998
		work for cross-domain annotation. <i>arXiv preprint</i>	999
		<i>arXiv:2410.17462</i> .	1000
		Haoxin Liu, Zhiyuan Zhao, Jindong Wang, Harshavard-	1001
		han Kamarthi, and B. Aditya Prakash. 2024. Lst-	1002
		prompt: Large language models as zero-shot time	1003

1004	series forecasters by long-short-term prompting. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 7832–7840, Bangkok, Thailand. Association for Computational Linguistics.	1058
1005		1059
1006		1060
1007		1061
1008	Penghang Liu, Elizabeth Fons, Svitlana Vyetrenko, Daniel Borrajo, Vamsi Potluru, and Manuela Veloso. 2025. Ts-agent: A time series reasoning agent with iterative statistical insight gathering. <i>arXiv preprint arXiv:2510.07432</i> .	1062
1009		1063
1010		1064
1011		1065
1012		1066
1013	Xu Liu, Yutong Xia, Yuxuan Liang, Junfeng Hu, Yiwei Wang, Lei Bai, Chao Huang, Zhenguang Liu, Bryan Hooi, and Roger Zimmermann. 2023. Largest: A benchmark dataset for large-scale traffic forecasting. <i>Advances in Neural Information Processing Systems</i> , 36:75354–75371.	1067
1014		1068
1015		1069
1016		1070
1017		1071
1018		1072
1019	Markus Löning, Anthony Bagnall, Sajaysurya Ganesh, Viktor Kazakov, Jason Lines, and Franz J Király. 2019. sktime: A unified interface for machine learning with time series. <i>arXiv preprint arXiv:1909.07872</i> .	1073
1020		1074
1021		1075
1022		1076
1023		1077
1024	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Swabha Swayamdipta, Yiming Yang, and Hannaneh Hajishirzi. 2023. Self-refine: Iterative refinement with self-feedback. <i>Preprint</i> , arXiv:2303.17651.	1078
1025		1079
1026		1080
1027		1081
1028		1082
1029		1083
1030	Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. 2018. The m4 competition: Results, findings, conclusion and way forward. <i>International Journal of forecasting</i> , 34(4):802–808.	1084
1031		1085
1032		1086
1033		1087
1034	Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. 2022. M5 accuracy competition: Results, findings, and conclusions. <i>International journal of forecasting</i> , 38(4):1346–1364.	1088
1035		1089
1036		1090
1037		1091
1038	Antoine Marot, Benjamin Donnot, Gabriel Dulac-Arnold, Adrian Kelly, Aidan O’ Sullivan, Jan Viebahn, Mariette Awad, Isabelle Guyon, Patrick Panciatici, and Camilo Romero. 2021. Learning to run a power network challenge: a retrospective analysis. In <i>NeurIPS 2020 competition and demonstration track</i> , pages 112–132. PMLR.	1092
1039		1093
1040		1094
1041		1095
1042		1096
1043		1097
1044		1098
1045	Aditya P Mathur and Nils Ole Tippenhauer. 2016. Swat: A water treatment testbed for research and training on ics security. In <i>2016 international workshop on cyber-physical systems for smart water networks (CySWater)</i> , pages 31–36. IEEE.	1099
1046		1100
1047		1101
1048		1102
1049		1103
1050	Larry R Medsker, Lakhmi Jain, and 1 others. 2001. Recurrent neural networks. <i>Design and applications</i> , 5(64-67):2.	1104
1051		1105
1052		1106
1053	Mohammad Amin Morid, Olivia R Liu Sheng, and Joseph Dunbar. 2023. Time series prediction using deep learning methods in healthcare. <i>ACM Transactions on Management Information Systems</i> , 14(1):1–29.	1107
1054		1108
1055		1109
1056		1110
1057		1111
	Charidimos Papadakis, Angeliki Dimitriou, Giorgos Filandrianos, Maria Lymperaioi, Konstantinos Thomas, and Giorgos Stamou. 2025. Atlas: Adaptive trading with llm agents through dynamic prompt optimization and multi-agent coordination. <i>arXiv preprint arXiv:2510.15949</i> .	1112
		1113
	John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S Tsay, Themis Palpanas, and Michael J Franklin. 2022. Tsb-uad: An end-to-end benchmark suite for univariate time-series anomaly detection. <i>Proc. VLDB Endow.</i> , 15(8):1697–1711.	1114
		1115
	German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. 2019. Continual lifelong learning with neural networks: A review. <i>Neural Networks</i> , 113:54–71.	1116
		1117
	Lawrence R Rabiner. 1989. A tutorial on hidden markov models and selected applications in speech recognition. <i>Proceedings of the IEEE</i> , 77(2):257–286.	1118
		1119
	Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. 2021. Stable-baselines3: Reliable reinforcement learning implementations. <i>Journal of machine learning research</i> , 22(268):1–8.	1120
		1121
	Chidaksh Ravuru, Sagar Srinivas Sakhinana, and Venkataramana Runkana. 2024. Agentic retrieval-augmented generation for time series analysis. <i>arXiv preprint arXiv:2408.14484</i> .	1122
		1123
	Shaina Raza, Ranjan Sapkota, Manoj Karkee, and Christos Emmanouilidis. 2025. Trism for agentic ai: A review of trust, risk, and security management in llm-based agentic multi-agent systems. <i>arXiv preprint arXiv:2506.04133</i> .	1124
		1125
	Hansheng Ren, Bixiong Xu, Yujing Wang, Chao Yi, Congrui Huang, Xiaoyu Kou, Tony Xing, Mao Yang, Jie Tong, and Qi Zhang. 2019. Time-series anomaly detection service at microsoft. In <i>Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining</i> , pages 3009–3017.	1126
		1127
	David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. 2020. Deepar: Probabilistic forecasting with autoregressive recurrent networks. <i>International Journal of Forecasting</i> , 36(3):1181–1191.	1128
		1129
	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. <i>Preprint</i> , arXiv:2302.04761.	1130
		1131
	Sebastian Schmidl, Phillip Wenig, and Thorsten Papenbrock. 2022. Anomaly detection in time series: a comprehensive evaluation. <i>Proceedings of the VLDB Endowment</i> , 15(9):1779–1797.	1132
		1133
	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. <i>Preprint</i> , arXiv:2303.11366.	1134
		1135

1114	Robert H Shumway and David S Stoffer. 2017. Arima models. In <i>Time series analysis and its applications: with R examples</i> , pages 75–163. Springer.	1167
1115		1168
1116		
1117	Xiaoyu Tao and 1 others. 2026. Cast-r1: Learning tool-augmented sequential decision policies for time series forecasting. <i>arXiv preprint arXiv:2602.13802</i> .	1169
1118		1170
1119		1171
1120	Stephen J Taylor. 2008. <i>Modelling financial time series</i> . world scientific.	1172
1121		1173
1122	Aleksandar Tomašević, Darja Cvetković, Sara Major, Slobodan Maletić, Miroslav Anđelković, Ana Vranić, Boris Stupovski, Dušan Vudragović, Aleksandar Bojgojević, and Marija Mitrović Dankulov. 2025. Towards operational validation of llm-agent social simulations: A replicated study of a reddit-like technology forum. <i>arXiv preprint arXiv:2508.21740</i> .	1174
1123		1175
1124		1176
1125		
1126		
1127		
1128		
1129	Ruey S Tsay. 2005. <i>Analysis of financial time series</i> . John wiley & sons.	1177
1130		1178
1131	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	1179
1132		1180
1133		1181
1134		
1135		
1136	Bingrui Wang, Yuan Zhou, Leijiao Ge, and Sun-Yuan Kung. 2025a. Large-model-based smart agent for time series anomaly detection in power systems. <i>Expert Systems with Applications</i> , page 128917.	1182
1137		1183
1138		1184
1139		1185
1140	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024a. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):186345.	1186
1141		1187
1142		1188
1143		1189
1144		1190
1145	Leizhen Wang, Peibo Duan, Zhengbing He, Cheng Lyu, Xin Chen, Nan Zheng, Li Yao, and Zhenliang Ma. 2025b. Agentic large language models for day-to-day route choices. <i>Transportation Research Part C: Emerging Technologies</i> , 180:105307.	1191
1146		1192
1147		1193
1148		1194
1149		1195
1150	Zhiruo Wang, Zhoujun Cheng, Hao Zhu, Daniel Fried, and Graham Neubig. 2024b. What are tools anyway? a survey from the language model perspective. <i>arXiv preprint arXiv:2403.15452</i> .	1196
1151		1197
1152		1198
1153		
1154	Michał Wawer and Jarosław A Chudziak. 2025. Integrating traditional technical analysis with ai: A multi-agent llm-based approach to stock market forecasting. <i>arXiv preprint arXiv:2506.16813</i> .	1199
1155		1200
1156		1201
1157		1202
1158	Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. 2020. Time series data augmentation for deep learning: A survey. <i>arXiv preprint arXiv:2002.12478</i> .	1203
1159		1204
1160		1205
1161		1206
1162	Muyan Weng, Defu Cao, Wei Yang, Yashaswi Sharma, and Yan Liu. 2026. Temporalbench: A benchmark for evaluating llm-based agents on contextual and event-informed time series tasks. <i>arXiv preprint arXiv:2602.13272</i> .	1207
1163		1208
1164		1209
1165		1210
1166		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222

1223	Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen, Jordan Suchow, Zhenyu Cui, Rong Liu, and 1 others. 2024. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. <i>Advances in Neural Information Processing Systems</i> , 37:137010–137045.	Haokun Zhao, Xiang Zhang, Jiaqi Wei, Yiwei Xu, Yuting He, Siqi Sun, and Chenyu You. 2025. Timeseries-scientist: A general-purpose ai agent for time series analysis. <i>arXiv preprint arXiv:2510.01538</i> .	1279
1224			1280
1225			1281
1226			1282
1227			
1228		Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> , 1(2).	1283
1229			1284
			1285
1230	Zhiyuan Zeng, Jiashuo Liu, Siyuan Chen, Tianci He, Yali Liao, Yixiao Tian, Jinpeng Wang, Zaiyuan Wang, Yang Yang, Lingyue Yin, and 1 others. 2025. Futurex: An advanced live benchmark for llm agents in future prediction. <i>arXiv preprint arXiv:2508.11987</i> .		1286
1231			1287
1232			
1233		Junhao Zheng, Chengming Shi, Xidi Cai, Qiuke Li, Duzhen Zhang, Chenxing Li, Dong Yu, and Qianli Ma. 2026. Lifelong learning of large language model based agents: A roadmap. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> .	1288
1234			1289
			1290
1235	Lingzhe Zhang, Yunpeng Zhai, Tong Jia, Xiaosong Huang, Chiming Duan, and Ying Li. 2025a. Agentfm: Role-aware failure management for distributed databases with llm-driven multi-agents. In <i>Proceedings of the 33rd ACM International Conference on the Foundations of Software Engineering</i> , pages 525–529.		1291
1236			1292
1237		Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 2021. Informer: Beyond efficient transformer for long sequence time-series forecasting. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 35, pages 11106–11115.	1293
1238			1294
1239			1295
1240			1296
1241			1297
1242	Qi Zhang, Chao Xu, Jie Li, Yicheng Sun, Jinsong Bao, and Dan Zhang. 2025b. Llm-tsfd: An industrial time series human-in-the-loop fault diagnosis method based on a large language model. <i>Expert Systems with Applications</i> , 264:125861.		1298
1243		Shu Zhou, Yunyang Xuan, Yuxuan Ao, Xin Wang, Tao Fan, and Hao Wang. 2025. Merit: Multi-agent collaboration for unsupervised time series representation learning. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 24011–24028.	1299
1244			1300
1245			1301
1246			1302
1247	REGINA ZHANG, SIQI GOH, XINGSHENG CHEN, ZONGRU LI, HONGGANG WEN, JIALE GUO, MENGLIN YANG, HONGZHI YIN, QIANG YANG, SIU-MING YIU, and 1 others. 2025. From prompts to agents: A comprehensive survey of llm-driven time series analysis. <i>arXiv preprint</i> .		1303
1248			
1249			
1250			
1251			
1252			
1253	Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, and 1 others. 2024a. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. In <i>Proceedings of the 30th acm sigkdd conference on knowledge discovery and data mining</i> , pages 4314–4325.		
1254			
1255			
1256			
1257			
1258			
1259			
1260	Xiyuan Zhang, Ranak Roy Chowdhury, Rajesh K Gupta, and Jingbo Shang. 2024b. Large language models for time series: A survey. <i>arXiv preprint arXiv:2402.01801</i> .		
1261			
1262			
1263			
1264	Yan Zhang, Ahmad Mohammad Saber, Amr Youssef, and Deepa Kundur. 2025a. Grid-agent: An llm-powered multi-agent system for power grid control. <i>arXiv preprint arXiv:2508.05702</i> .		
1265			
1266			
1267			
1268	Yuxuan Zhang, Yangyang Feng, Daifeng Li, Kexin Zhang, Junlan Chen, and Bowen Deng. 2025b. Can competition enhance the proficiency of agents powered by large language models in the realm of news-driven time series forecasting? <i>arXiv preprint arXiv:2504.10210</i> .		
1269			
1270			
1271			
1272			
1273			
1274	Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2024c. A survey on the memory mechanism of large language model based agents. <i>Preprint</i> , arXiv:2404.13501.		
1275			
1276			
1277			
1278			

A Survey Comparison

This section positions our survey against closely related surveys on LLMs for time-series analysis. Table 1 compares them along five dimensions:

(i) *TS-specific*, which indicates whether the survey is explicitly designed for time-series problems rather than general LLM or agent settings;

(ii) *LLM Agents*, which indicates whether the survey focuses on agentic systems in which the LLM plays an active role in multi-step decision-making, where “Partial” means that agentic systems are included but are not the main focus;

(iii) The four task columns, which indicate whether the survey clearly covers each major problem family in our taxonomy—*Forecasting & Reasoning*, *Augmentation & Synthesis*, *Anomaly Detection & Diagnosis*, and *Decision Support*—where ✓ denotes clear coverage, ~ partial coverage, and × that the topic is not a main focus;

(iv) *Taxonomy*, which describes the main organizing principle of the survey;

(v) *Design Guidance*, which indicates whether the survey provides practical guidance for selecting or understanding system designs in concrete settings. Under *Taxonomy*, “Problem-driven” groups methods by time-series problems, “Capability-driven” by agent capabilities such as planning, memory, or tool use, “Method-driven” by methodological categories such as prompting, fine-tuning, tokenization, or model integration, and “Other” covers organizing principles that do not fit these categories cleanly. Under *Design Guidance*, “Limited” means that guidance is provided only indirectly or at a high level.

B Architecture Overview

Figure 2 provides a visual summary of the architectural patterns discussed in Section 3. It organizes time-series agent systems into single-agent and multi-agent settings and summarizes the main architectural patterns in our taxonomy, including adaptive, iterative, and continual single-agent designs, as well as sequential pipeline, planner-executor, and committee/debate systems.

C Table 2: Method Summary

Table 2 summarizes the representative methods covered in this survey. The methods are organized primarily by *Problem Type*, following the problem-driven taxonomy adopted in this work. For each method, the table reports its publication

Year and *Venue*, and summarizes its classification along three agent-oriented dimensions: *Architecture*, *Tools*, and *Memory*. Detailed discussions are provided in Section 4.

D Table 3: Resources Summary

Table 3 provides a supplementary summary of the resources discussed in Section 5. The resources are organized by *Problem Type*, following the problem-driven taxonomy adopted in this work, and include benchmarks, datasets, environments, and toolkits. For each listed resource, the table reports the following practical attributes: (i) *Category*, the category of resource, such as benchmark, dataset, environment, toolkit, or dataset repository; (ii) *Problem Type*, the primary time-series problem or application domain the resource is associated with; (iii) *Interactive*, whether the resource supports interactive or closed-loop evaluation, where agent actions can affect subsequent observations or outcomes; (iv) *Timesteps*, the approximate temporal length of the data, when such information is available; (v) *Series*, the approximate number of variables or channels in the resource, when applicable; and (vi) *Last Update*, the most recent reported release, update year, or maintenance status. This table is intended as a compact reference to the practical characteristics of commonly used resources.

E Typical Design Patterns Across Time-Series Tasks

Figure 3 provides a supplementary visual summary of the typical design patterns discussed across different time-series tasks in Section 4. Each panel shows a common design pattern rather than the exact design of a specific system.

In the figure, (i) dashed boxes denote agent modules together with their tools, while solid boxes group entities with similar roles; (ii) solid arrows show the main flow, while dashed arrows indicate conditional or fallback paths; and (iii) background colors are used to group panels from the same higher-level problem type.

Table 1: Comparison with related surveys. ✓: clearly covered; ~: partially covered; ×: not covered.

Survey	Year	TS-specific	LLM Agents	Forecasting & Reasoning	Augmentation & Synthesis	Anomaly Detection & Diagnosis	Decision Support	Taxonomy	Design Guidance
(Jiang et al., 2024)	2024	Yes	No	✓	~	✓	~	Method-driven	Limited
(Zhang et al., 2024b)	2024	Yes	No	✓	✓	✓	×	Method-driven	Limited
(Chang et al., 2025)	2025	Yes	Partial	✓	✓	✓	✓	Other	Limited
(ZHANG et al., 2025)	2025	Yes	Partial	✓	✓	✓	~	Capability-driven	Limited
Ours	2026	Yes	Yes	✓	✓	✓	✓	Problem-driven	Yes

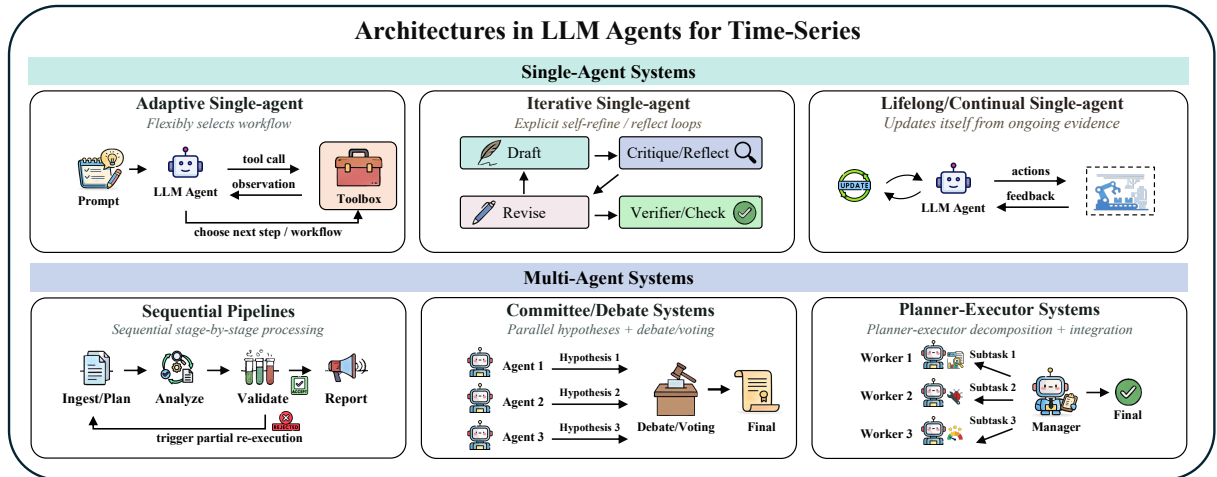


Figure 2: An illustration of the main architectural patterns for time-series agent systems.

Table 2: Summary of representative LLM-based agentic methods for time-series tasks. Each method is grouped by **Problem Type** and further characterized by its publication year, venue, architecture, tools, and memory.

Method	Year	Venue	Problem Type	Architecture	Tools	Memory
TimeSeriesScientist (Zhao et al., 2025)	2024	arXiv	Forecasting	Multi (Sequential Pipeline)	Statistical & ML Models	Evidence Logs
TimeXL (Jiang et al., 2025)	2025	arXiv	Forecasting	Multi (Sequential Pipeline)	None	Evidence Logs, Pattern Library
TimeCAP (Lee et al., 2025)	2025	AAAI	Forecasting	Multi (Sequential Pipeline)	Data Processing	Pattern Library
NewsTSForecasting (Zhang et al., 2025b)	2025	arXiv	Forecasting	Multi (Hybrid)	Search & Retrieval APIs	Evidence Logs, Analytical Strategies
TRACE (Chen and Xie, 2025)	2025	TAI	Forecasting	Multi (Committee / Debate)	None	Evidence Logs
FLAIRR-TS (Jalori et al., 2025)	2025	arXiv	Forecasting	Multi (Sequential Pipeline)	None	Evidence Logs, Analytical Strategies
TS-Agent (Liu et al., 2025)	2024	arXiv	Reasoning	Single (Iterative)	Data Processing, Statistical & ML Models	Evidence Logs
Agentic-RAG (Ravuru et al., 2024)	2024	arXiv	Reasoning	Multi (Planner–Executor)	None	Pattern Library
TS-Reasoner (Ye et al., 2024b)	2024	arXiv	Reasoning	Single (Iterative)	Database APIs, Data Processing, Statistical	Evidence Logs
ZARA (Li et al., 2025b)	2025	arXiv	Reasoning	Multi (Sequential Pipeline)	Database APIs; Data Processing; Statistical & ML Models	Pattern Library, Analytical Strategies
CLIMATEAGENT (Kim et al., 2025)	2025	arXiv	Reasoning	Multi (Planner–Executor)	Database APIs, Data Processing, Simulators & Solvers & Optimizers	Evidence Logs
TESSA (Lin et al., 2024)	2024	arXiv	Augmentation	Multi (Sequential Pipeline)	Data Processing	None
DCATS (Yeh et al., 2025)	2025	arXiv	Augmentation	Single (Iterative)	Statistical & ML Models	Evidence Logs
MERIT (Zhou et al., 2025)	2025	ACL	Augmentation	Multi (Sequential Pipeline)	Data Processing, Statistical & ML Models	None
ChatTS (Xie et al., 2024)	2024	ICSP	Synthesis	Multi (Sequential Pipeline)	None	None
AXIS (Lan et al., 2025)	2025	arXiv	Synthesis	Multi (Sequential Pipeline)	None	None
GenAI4RiskModeling (Joshi, 2025)	2025	SSRN	Synthesis	Multi (Sequential Pipeline)	Database APIs, Data Processing, Statistical & ML Models	None
BRIDGE (Li et al., 2025a)	2025	arXiv	Synthesis	Multi (Hybrid)	Search APIs, Statistical & ML Models	Pattern Library
AD-AGENT (Yang et al., 2025)	2025	arXiv	Detection	Multi (Sequential Pipeline)	Data Processing, Statistical & ML Models	Evidence Logs, Analytical Strategies
Argos (Gu et al., 2025)	2025	arXiv	Detection	Multi (Sequential Pipeline)	Data Processing, Statistical & ML Models	Evidence Logs
SLEP	2025	SSRN	Detection	Single (Adaptive)	Data Processing	Evidence Logs, Analytical Strategies
CALM	2025	arXiv	Detection	Single (Lifelong / Continual)	Data Processing, Statistical & ML Models	None
LEMAD (Zhang et al., 2025b)	2025	Electron.	Diagnosis	Multi (Sequential Pipeline)	Data Processing, Statistical & ML Models	Evidence Logs
AgentFM (Zhang et al., 2025a)	2025	FSE	Diagnosis	Multi (Planner–Executor)	Database APIs, Data Processing, Statistical & ML Models	None
LLM-TSFD (Zhang et al., 2025b)	2025	ESWA	Diagnosis	Single (Iterative)	Database APIs, Data Processing, Search & Retrieval APIs, Statistical & ML Models	Pattern Library
ElliottAgents (Chudziak and Wawer, 2024; Wawer and Chudziak, 2025)	2024	PACLIC	Trading	Multi (Planner–Executor)	Data Processing, Statistical & ML Models	Pattern Library
TradingAgents (Xiao et al., 2024)	2024	arXiv	Trading	Multi (Hybrid)	Data Processing, Search & Retrieval APIs	Evidence Logs
FinCon (Yu et al., 2024)	2024	NeurIPS	Trading	Multi (Planner–Executor)	Database APIs, Data Processing, Search & Retrieval APIs	Evidence Logs, Pattern Library, Analytical Strategies
FinAgent (Zhang et al., 2024a)	2024	WWW	Trading	Multi (Sequential Pipeline)	Data Processing	Pattern Library, Analytical Strategies
FactorMAD (Duan et al., 2025)	2024	arXiv	Trading	Multi (Committee / Debate)	Data Processing, Statistical & ML Models	Patterns Library
FinMem (Yu et al., 2025)	2025	AAAI	Trading	Single (Adaptive)	Database APIs, Data Processing	Evidence Logs, Patterns Library
FinArena (Xu et al., 2025)	2025	arXiv	Trading	Multi (Planner–Executor)	Data Processing, Search & Retrieval APIs	None
ATLAS (Papadakis et al., 2025)	2025	arXiv	Trading	Multi (Planner–Executor)	Search & Retrieval APIs, Simulators & Solvers & Optimizers	Analytical Strategies
FinRL-DeepSeek (Benhenda, 2025)	2025	arXiv	Trading	Single (Adaptive)	Statistical & ML Models	None
Open-TI (Da et al., 2024)	2024	arXiv	Infrastructure Control	Multi (Sequential Pipeline)	Database APIs, Data Processing, Statistical & ML Models	None
Grid-Agent (Zhang et al., 2025a)	2024	arXiv	Infrastructure Control	Multi (Sequential Pipeline)	Simulators & Solvers & Optimizers	Evidence Logs
GridMind (Jin et al., 2025)	2025	SC Work-shops	Infrastructure Control	Multi (Planner–Executor)	Simulators & Solvers & Optimizers	Evidence Logs
LLMTraveler (Wang et al., 2025b)	2025	TRC	Infrastructure Control	Single (Iterative)	None	Evidence Logs
LLMLight (Lai et al., 2025)	2025	KDD	Infrastructure Control	Single (Iterative)	Simulators & Solvers & Optimizers	None

Table 3: Summary of representative resources for time-series tasks. Each resource is grouped by **Problem Type** and further characterized by its type, interactivity, temporal scale, number of series, and update status.

Resource	Category	Problem Type	Interactive	Timesteps	Series	Last Update
<i>Forecasting</i>						
M4 (Makridakis et al., 2018)	Benchmark	Forecasting	No	Up to 10k	~100k	2020
M5 (Makridakis et al., 2022)	Benchmark	Forecasting	No	1.9k	42k	2022
MIRAI (Ye et al., 2024a)	Benchmark	Forecasting	No	~ 1M	59k	2025
FutureX (Zeng et al., 2025)	Benchmark	Forecasting	No	~ 1k	195	2025
Monash TSF (Godahewa et al., 2021)	Dataset Repo	Multi-Domain Forecasting	No	Up to 527k	Up to 145k	2021
IHEPC	Dataset	Energy Forecasting	No	2M	9	2011
Electricity (Lai et al., 2018)	Dataset	Energy Forecasting	No	26k	321	2015
METR-LA (Li et al., 2017)	Dataset	Traffic Forecasting	No	34k	207	2018
PEMS-BAY (Li et al., 2017)	Dataset	Traffic Forecasting	No	52k	325	2018
ETT (Zhou et al., 2021)	Dataset	Energy Forecasting	No	Up to 69k	7	2021
Exchange (Lai et al., 2018)	Dataset	Finance Forecasting	No	7.6k	8	2021
Traffic (Lai et al., 2018)	Dataset	Traffic Forecasting	No	17k	862	2023
Weather	Dataset	Weather Forecasting	No	53k	21	2023
LargeST (Liu et al., 2023)	Dataset	Traffic Forecasting	No	526k	Up to 8.6k	2023
ILI	Dataset	Health Forecasting	No	Unclear	7	Ongoing
StatsForecast (Garza et al., 2022)	Toolkit	Forecasting	No	N/A	N/A	2025
NeuralForecast (Challu et al., 2023)	Toolkit	Forecasting	No	N/A	N/A	2026
<i>Reasoning</i>						
TimeSeriesExam (Cai et al., 2024)	Benchmark	Reasoning	No	N/A	>700 Tasks	2024
IBTrACS (Knapp et al., 2010)	Dataset	Reasoning	No	N/A	Unclear	2025
<i>Anomaly Detection</i>						
TSB-UAD (Paparrizos et al., 2022)	Benchmark Suite	Anomaly Detection	No	Varies	236 Datasets	2022
ADBench (Han et al., 2022)	Benchmark Suite	Anomaly Detection	No	Varies	57 Datasets	2022
NAB (Lavin and Ahmad, 2015)	Benchmark	Anomaly Detection	No	Up to 22k	58	2015
UCR-AD (Wu and Keogh, 2021)	Benchmark	Anomaly Detection	No	Unclear	250	2021
SWaT (Mathur and Tippenhauer, 2016)	Dataset	Anomaly Detection	No	Unclear	Unclear	2016
SMAP (Hundman et al., 2018)	Dataset	Anomaly Detection	No	430k	55	2018
MSL (Hundman et al., 2018)	Dataset	Anomaly Detection	No	67k	27	2018
WADI (Adepu et al., 2020)	Dataset	Anomaly Detection	No	Unclear	103	2019
KPI (Ren et al., 2019)	Dataset	Anomaly Detection	No	Unclear	781	2019
ADTK	Toolkit	Anomaly Detection	No	N/A	N/A	2020
Luminol	Toolkit	Anomaly Detection	No	N/A	N/A	2023
<i>Decision-support</i>						
RL Unplugged (Gulcehre et al., 2020)	Dataset	Offline RL	No	N/A	N/A	2020
D4RL (Fu et al., 2020)	Dataset	Offline RL	No	N/A	N/A	2021
Grid2Op (Marot et al., 2021)	Environment	Power System Control	Yes	N/A	N/A	2021
SocioDojo (Cheng and Chin, 2024)	Environment	Trading	Yes	N/A	N/A	2024
CityLearn	Environment	Power System Control	Yes	N/A	N/A	2025
SUMO (Behrisch et al., 2011)	Environment	Traffic Control	Yes	N/A	N/A	Ongoing
Stable-Baselines3 (Raffin et al., 2021)	Toolkit	Reliable RL	No	N/A	N/A	2025
<i>Others</i>						
TimeSeriesGym (Cai et al., 2025)	Benchmark	Multi-Task	No	N/A	N/A	2025
Sktime (Löning et al., 2019)	Toolkit	Multi-Task	No	N/A	N/A	2025
GluonTS (Alexandrov et al., 2020)	Toolkit	Multi-Task	No	N/A	N/A	2025
Merlion (Bhatnagar et al., 2021)	Toolkit	Multi-Task	No	N/A	N/A	2024

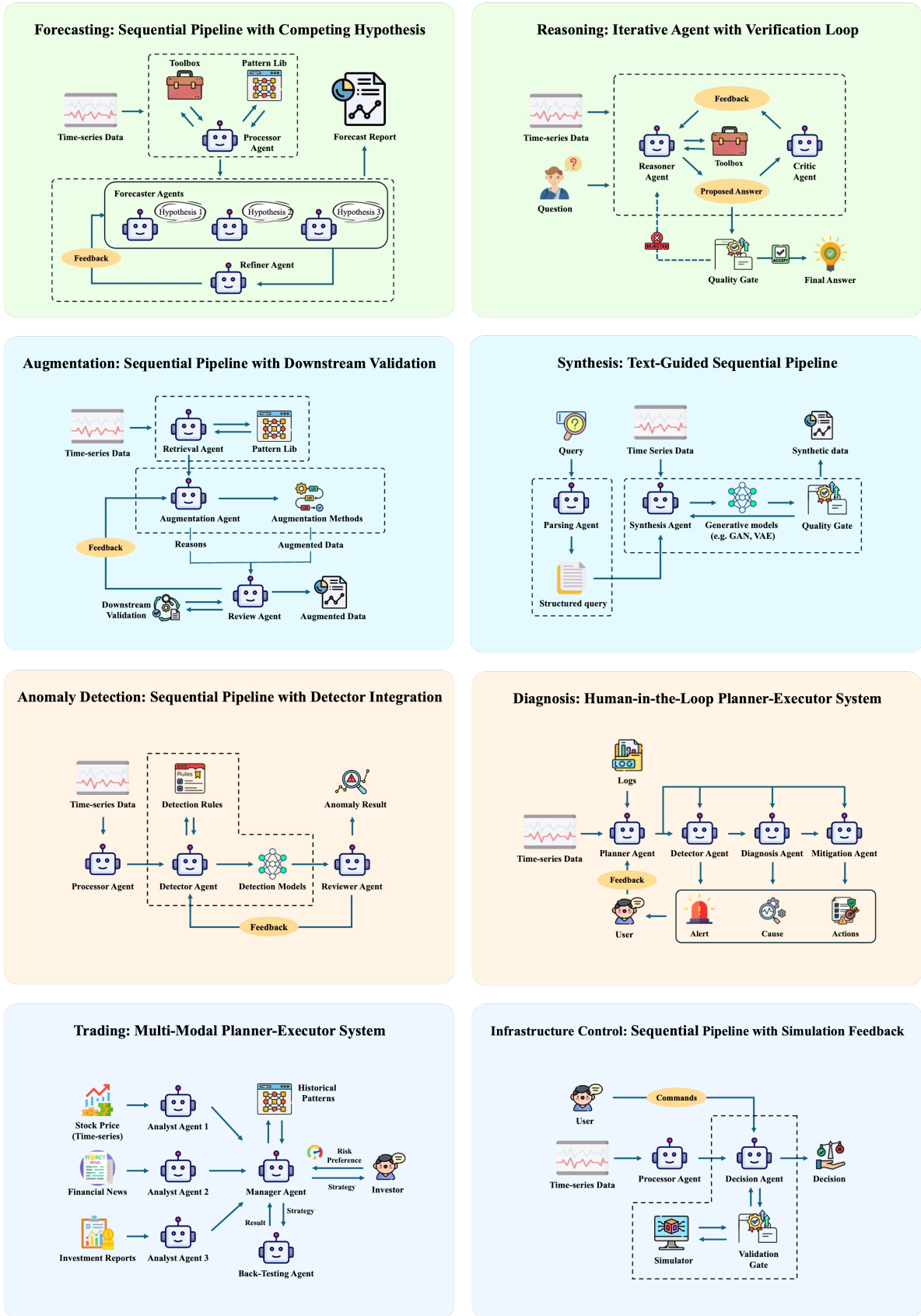


Figure 3: Typical design patterns of time-series agent systems across representative tasks. From top left to bottom right, the panels correspond to forecasting, reasoning, augmentation, synthesis, anomaly detection, diagnosis, trading, and infrastructure control.