Healthcare TimeSeries Reasoning Benchmarks at Scale

Anonymous Author(s)

Affiliation Address email

Abstract

We introduce TimeSeriesExamAgent, a scalable and domain-agnostic framework for automatically generating and validating time series reasoning benchmarks. Existing benchmarks lack scalability, are limited to a few specific domains, while building them remains labor intensive. Automated solutions for benchmark creation have been proposed, but these typically rely on a single-step generation process without verification, leading to lower-quality exams. Our framework addresses these limitations by enabling stakeholders—such as institutions with highly confidential data—to easily construct high-quality, domain-specific benchmarks from their own private datasets. A domain expert provides a dataset, a natural language description, and a simple data-loading method. The agent then orchestrates the generation pipeline, including creating question templates, robustness verification from multiple perspectives, and iterative refinement. We demonstrate the framework on two medical datasets and evaluate multiple state-of-the-art language models on the generated benchmarks. Empirically, we demonstrate that the framework produces domain-agnostic benchmarks whose diversity matches human-generated counterparts, and our evaluation of several Large Language Models shows that accuracy remains limited, underscoring open challenges in time-series reasoning.

1 Introduction

2

3

4

6

8

9

10

11

12 13

14 15

16

17

- Many recent works have applied Large Language Models (LLMs) to time series analysis tasks such as forecasting, anomaly detection, and classification [1, 2, 3, 4, 5, 6]. More recently, attention has shifted to evaluating the reasoning capabilities of LLMs in time series tasks. These evaluations are typically framed in two ways: 1) contextualized traditional tasks such as forecasting, but with added contextual information (e.g., providing a clinical scenario before a prediction) [7, 8, 9, 10, 11], and 2) reasoning and understanding tasks that directly probe concepts in time series (e.g., "what kind of trend does the following series exhibit?") [12, 13].
- However, existing benchmarks have clear limitations. Contextualized tasks remain close to traditional metrics (e.g., mean-squared-error for forecasting) without testing deeper reasoning, while reasoning-style benchmarks often focus only on simple properties like trend or seasonality. In practice, real-world domains such as healthcare require more complex reasoning, where tasks like diagnosis naturally combine anomaly detection, classification, and domain knowledge. Curation is another challenge. Annotation or template-based benchmarks are labor-intensive, while LLM-based augmentation often lacks diversity because it simply expands existing datasets. As a result, building specialized, domain-specific benchmarks remains difficult and time-consuming.
- This challenge is especially pronounced in healthcare. Clinical datasets are both highly sensitive and highly specialized: tasks in cardiology, radiology, or genomics often require nuanced reasoning that combines statistical patterns with medical expertise and domain knowledge. Unlike open domains

where public benchmarks are available, stakeholders in healthcare cannot easily share their datasets due to patient privacy, regulatory constraints, and ethical concerns. At the same time, they need ways to rigorously evaluate whether generative models can reason about ECG signals, imaging studies, or longitudinal patient records without exposing protected health information. This creates an urgent need for customizable, automated benchmark generation tailored to private clinical datasets.

Inspired by recent agent-based approaches in other domains [14, 15], we propose TimeSeriesExamAgent, a pipeline that (1) generates domain-specific multiple-choice questions

TimeSeriesExamAgent, a pipeline that (1) generates domain-specific multiple-choice questions on time-series data, (2) scales efficiently, and (3) ensures reliable ground truth through iterative verification. We also evaluate four state-of-the-art LLMs on a benchmark of 348 automatically generated questions samples. Our results show that many models struggle on highly complex tasks that require combining quantitative analysis with domain knowledge. For brevity, we provide detailed related work in Appendix A.

9 2 TimeSeriesExamAgent

59

60

61

62

63

64

65

67

68

69

70

71

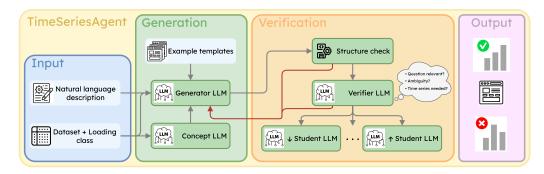


Figure 1: TimeSeriesExamAgent architecture. The user provides exam-making instructions and a custom dataset with minimal loading code. Agent outputs question templates – Python functions generated by a generator LLM and filtered through three progressive stages of verification (syntax and output format check, validation by LLM judge, capability-aligned filtering). Arrows denote data flow, red ones show direction for rejected templates.

In this section, we introduce TimeSeriesExamAgent, a multi-agent framework that combines 50 planning, generation, and verification to enable automatic benchmark construction. In this section, 51 we describe TimeSeriesExamAgent and its workflow in detail. An overview is shown in Fig. 1. 52 The Generation Agent takes as input a description of the natural language task T and a data set 53 D. The description T may include user guidelines for generation, contextual information about the 54 dataset, or other relevant instructions. For convenience, we denote each sample in D as (x_i, z_i) , 55 where $x_i \in \mathbb{R}^{n \times d}$ is a time series with n observations and d variables, and z_i is an auxiliary array 56 containing metadata or labels related to the series. The user provides a dataset class D that supports 57 basic operations such as querying the *i*-th sample.

Generation We generate question templates instead of samples directly, as shown in Fig. 2. Templates offer two advantages: they are scalable, and their abstraction adds an extra layer of robustness. By relying on structured, rule-based generation rather than manual inputs, they reduce the chance of human errors or inconsistencies. Our generator LLM produces a predefined number of templates, each implemented as a Python function. A template contains a formatted string for the question and options, together with parameters that control how many questions to generate. For each question, the template samples a pair (x_i, z_i) from the dataset D and applies a rule-based calculation to determine the correct answer from the time series. For example, in a trend-detection template, the function computes the linear trend coefficient of x_i and selects "Yes, there is a linear trend" if the coefficient exceeds a specified threshold. In addition to such signal-derived logic, templates can also utilize the auxiliary property z_i , effectively transforming classification problems into question—answer form. For instance, if an ECG series in the dataset is labeled as exhibiting atrial fibrillation, the template can present this label as one of the multiple-choice options. Each generated sample consists of the question, its options, the correct answer, and one or more associated time series represented as

numerical values. We provide a breakdown of the Generation Agent and its prompt in Appendix B.

An example template is also provided.

Verification We observe 75 LLM-based generation frequently produces errors or 77 irrelevant outputs, motivating 78 the need for a structured veri-79 fication process. We propose a 80 multistage verification process 81 to check the accuracy and 82 relevance of each template. If a template fails at any stage, it is returned to the generation agent 85 with feedback. The generation 86 is iterative with a maximum of 87 three attempts, after which the 88 ongoing template is discarded to 89 avoid excessive context length 90 and cost from repeated failures. 91

102

103

105

106

107

108

109

110

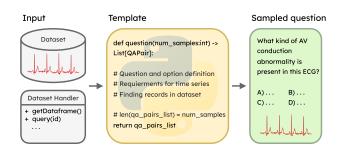


Figure 2: Question generation process: With information about dataset, TimeSeriesAgent generates question template in a form of Python functions. The created function can be called to get arbitrary number of question samples.

92 **Structure verification** We check whether the generated template can be executed successfully. We 93 execute the generated template k=5 times; if there are failures, the error message is returned as a 94 feedback.

Content verification Certain aspects of quality control are particularly well-suited for LLM-as-ajudge evaluation. For example, verifying that a question is grammatically correct, free of ambiguity
or bias, and genuinely answerable from the accompanying time series can be effectively handled by
an LLM. To this end, we use an LLM verifier to assess the validty of each template. A quantitative
score is given, and we set a threshold for rejection. If the verifier raises any rejection, its explanation
is treated similarly to a structural error and the template is regenerated. We provide the detailed
prompt in Appendix C.

Capability-Aligned Filtering To detect templates that generate overly simple or irrelevant exams, we evaluate them using a set of test-taking LLMs with varying capabilities. This approach is inspired by educational theory, particularly the expertise reversal effect [16]. A template is discarded if weaker LLMs achieve higher average accuracy than stronger models, as this typically indicates that the template is flawed or noisy rather than genuinely discriminative. Templates are retained if performance scales with model capability— or if all models perform poorly, since such questions may still capture genuine difficulty. We provide hyper-parameters in Appendix F and other design specifics in Appendix D

3 Experimental Setup, Results and Discussion

First, we generate one exam for each of the two real world datasets: PTB-XL [17], MIT-BIH [18]. In total, we have 197 samples for MIT-BIH, and 151 samples for PTB-XL. We sample 4 or 5 instances per template. Thus, the difference in the number of generated samples is a result of the template filtering mechanism above.

We select candidate models to cover a diverse range of performance levels, as indicated by the OpenVLM Leaderboard [23]. In both cases, the best results achieves Gemma-3-27b-it, which outperforms the remaining models.

To further evaluate our benchmark, we compare multiple metrics on questions generated from the dataset with those in ECG-QA [10], a template-based benchmark also built on PTB-XL. The goal is to demonstrate that our framework achieves comparable diversity without requiring manual template curation. We picked random 50 question samples from each benchmark and calculated the distances for every possible pair within the set. We used the Qwen/Qwen3-0.6B sentence transformer model to extract embeddings, as it achieved the second-best performance among all models on the Hugging Face MTEB leaderboard.

	Dataset		
Model	MIT-BIH	PTB-XL	
gpt-4o [19]	0.416	0.424	
o3-mini [20]	0.442	0.477	
Qwen2.5-VL-Instruct [21]	0.411	0.490	
Gemma-3-27b-it [22]	0.497	0.517	

Table 1: Comparative performance of four vision-language models across medical (MIT-BIH, PTB-XL) time-series datasets.

	Mean \pm Std		
Benchmark Dataset	Embedding	Normalized Levenshtein	
ECG-QA	0.207 ± 0.079	0.519 ± 0.157	
TimeSeriesExamAgent (ours)	0.301 ± 0.070	$\textbf{0.542} \pm \textbf{0.039}$	

Table 2: Question diversity comparison using embedding and normalized Levenshtein distance.

As shown in Table 2, benchmark generated by our framework shows a diversity comparable to one developed by humans. This indicates that the proposed framework is able to capture a wide range of expressions without relying on handcrafted templates, supporting its scalability and adaptability to other domains. We also employed G-Eval, a probabilistic LLM-as-a-judge framework [24]. An LLM is used to evaluate the relevance of each question, assigning a score between 0 and 1 to indicate how well it meets the specified criteria. Results are presented in Table 3. We provide the detailed G-Eval prompt in Appendix E.

	Mean Result				
Dataset	Specificity	Unambiguity	Domain Relevance	Answerability	
ECG-QA	0.604	0.562	0.827	0.898	
TimeSeriesExamAgent (ours)	0.922	0.932	0.989	0.992	

Table 3: Question diversity comparison using G-Eval framework.

4 Limitations and Conclusions

In this work, we present a scalable, domain-specific framework for the automatic generation of time-series benchmarks, enabling the creation of high-quality, large-scale evaluation datasets while minimizing the need for labor-intensive human annotation. A limitation of this study is that the quality of the generated exams depends on the quality and coverage of the time series dataset. Additionally, domain specialists must provide carefully crafted prompts.

For future work, we will explore human-in-the-loop improvements to template generation. In offline sessions with clinicians, we observed that exams produced with such feedback are more likely to be deemed valid. We also plan to validate exam quality by training time series—text alignment models and testing their transfer performance on other established reasoning benchmarks [8]. Finally, there is growing attention on building time series agentic frameworks [25, 26]. Enabling these frameworks to write code in order to answer our benchmark questions would provide valuable insights to the community.

References

146

147

125

126

127

128

130

[1] Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham

- Kapoor, et al. Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.
- [2] Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.
 Moment: A family of open time-series foundation models. *arXiv preprint arXiv:2402.03885*,
 2024.
- 153 [3] Defu Cao, Furong Jia, Sercan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu.
 154 Tempo: Prompt-based generative pre-trained transformer for time series forecasting. *arXiv*155 *preprint arXiv:2310.04948*, 2023.
- 156 [4] Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu
 157 Chen, Yuxuan Liang, Yuan-Fang Li, Shirui Pan, et al. Time-llm: Time series forecasting by
 158 reprogramming large language models. *arXiv preprint arXiv:2310.01728*, 2023.
- [5] Tian Zhou, Peisong Niu, Liang Sun, Rong Jin, et al. One fits all: Power general time series
 analysis by pretrained lm. *Advances in neural information processing systems*, 36:43322–43355,
 2023.
- [6] Nina Żukowska, Mononito Goswami, Michał Wiliński, Willa Potosnak, and Artur Dubrawski. Towards long-context time series foundation models. *arXiv preprint arXiv:2409.13530*, 2024.
- [7] Jialin Chen, Aosong Feng, Ziyu Zhao, Juan Garza, Gaukhar Nurbek, Cheng Qin, Ali Maatouk,
 Leandros Tassiulas, Yifeng Gao, and Rex Ying. Mtbench: A multimodal time series benchmark
 for temporal reasoning and question answering. arXiv preprint arXiv:2503.16858, 2025.
- [8] Yaxuan Kong, Yiyuan Yang, Yoontae Hwang, Wenjie Du, Stefan Zohren, Zhangyang Wang, Ming Jin, and Qingsong Wen. Time-mqa: Time series multi-task question answering with context enhancement. *arXiv preprint arXiv:2503.01875*, 2025.
- [9] Haoxin Liu, Shangqing Xu, Zhiyuan Zhao, Lingkai Kong, Harshavardhan Prabhakar Kamarthi,
 Aditya Sasanur, Megha Sharma, Jiaming Cui, Qingsong Wen, Chao Zhang, et al. Time-mmd:
 Multi-domain multimodal dataset for time series analysis. Advances in Neural Information
 Processing Systems, 37:77888–77933, 2024.
- 174 [10] Jungwoo Oh, Gyubok Lee, Seongsu Bae, Joon-myoung Kwon, and Edward Choi. Ecg-qa: A comprehensive question answering dataset combined with electrocardiogram. *Advances in Neural Information Processing Systems*, 36:66277–66288, 2023.
- 177 [11] Xu Wang, Jiaju Kang, Puyu Han, Yubao Zhao, Qian Liu, Liwenfei He, Lingqiong Zhang,
 178 Lingyun Dai, Yongcheng Wang, and Jie Tao. Ecg-expert-qa: A benchmark for evaluating
 179 medical large language models in heart disease diagnosis. arXiv preprint arXiv:2502.17475,
 180 2025.
- 181 [12] Yifu Cai, Arjun Choudhry, Mononito Goswami, and Artur Dubrawski. Timeseriesexam: A time series understanding exam. *arXiv preprint arXiv:2410.14752*, 2024.
- [13] Willa Potosnak, Cristian Challu, Mononito Goswami, Kin G. Olivares, Michał Wiliński, Nina
 Żukowska, and Artur Dubrawski. Investigating compositional reasoning in time series foundation models, 2025.
- [14] Natasha Butt, Varun Chandrasekaran, Neel Joshi, Besmira Nushi, and Vidhisha Balachandran. Benchagents: Automated benchmark creation with agent interaction. arXiv preprint arXiv:2410.22584, 2024.
- 189 [15] Gauthier Guinet, Behrooz Omidvar-Tehrani, Anoop Deoras, and Laurent Callot. Automated 190 evaluation of retrieval-augmented language models with task-specific exam generation. *arXiv* 191 *preprint arXiv:2405.13622*, 2024.
- 192 [16] Slava Kalyuga. Expertise reversal effect and its implications for learner-tailored instruction.
 193 *Educational psychology review*, 19(4):509–539, 2007.
- 194 [17] Patrick Wagner, Nils Strodthoff, Ralf-Dieter Bousseljot, Dieter Kreiseler, Fatima I Lunze, Wojciech Samek, and Tobias Schaeffter. Ptb-xl, a large publicly available electrocardiography dataset. *Scientific data*, 7(1):1–15, 2020.

- [18] George B Moody and Roger G Mark. The impact of the mit-bih arrhythmia database. *IEEE* engineering in medicine and biology magazine, 20(3):45–50, 2001.
- [19] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,
 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv
 preprint arXiv:2410.21276, 2024.
- 202 [20] OpenAI. Openai o3-mini system card. https://openai.com/index/ 203 o3-mini-system-card/, January 2025. Accessed: 2025-08-22.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,
 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang
 Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen
 Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report.
 arXiv preprint arXiv:2502.13923, 2025.
- [22] Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona
 Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma
 3 technical report. arXiv preprint arXiv:2503.19786, 2025.
- [23] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong,
 Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating
 large multi-modality models. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 11198–11201, 2024.
- 216 [24] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval:
 217 Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*,
 218 2023.
- 219 [25] Yifu Cai, Xinyu Li, Mononito Goswami, Michał Wiliński, Gus Welter, and Artur Dubrawski.
 220 Timeseriesgym: A scalable benchmark for (time series) machine learning engineering agents.
 221 arXiv preprint arXiv:2505.13291, 2025.
- 222 [26] Wen Ye, Wei Yang, Defu Cao, Yizhou Zhang, Lumingyuan Tang, Jie Cai, and Yan Liu. Domain-223 oriented time series inference agents for reasoning and automated analysis. *arXiv preprint* 224 *arXiv:*2410.04047, 2024.
- Yilin Wang, Peixuan Lei, Jie Song, Yuzhe Hao, Tao Chen, Yuxuan Zhang, Lei Jia, Yuanxiang
 Li, and Zhongyu Wei. Itformer: Bridging time series and natural language for multi-modal qa
 with large-scale multitask dataset. arXiv preprint arXiv:2506.20093, 2025.
- [28] Wanying Wang, Zeyu Ma, Pengfei Liu, and Mingang Chen. Testagent: A framework for domain adaptive evaluation of llms via dynamic benchmark construction and exploratory interaction.
 arXiv preprint arXiv:2410.11507, 2024.
- [29] Mike A Merrill, Mingtian Tan, Vinayak Gupta, Tom Hartvigsen, and Tim Althoff. Language
 models still struggle to zero-shot reason about time series. arXiv preprint arXiv:2404.11757,
 2024.
- [30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
 Advances in neural information processing systems, 35:24824–24837, 2022.
- [31] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented
 generation for knowledge-intensive nlp tasks. Advances in neural information processing
 systems, 33:9459–9474, 2020.
- [32] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- [33] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer,2020.

- 246 [34] Jesse Mu, Xiang Lisa Li, and Noah Goodman. Learning to compress prompts with gist tokens, 247 2024.
- Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang, Zhiliang Tian, and Liang Ding. Recursively
 summarizing enables long-term dialogue memory in large language models, 2025.
- Zio [36] Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi
 Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. arXiv preprint arXiv:2407.16741, 2024.

A Related Work

Time series benchmarks The task of creating domain-specific time series reasoning benchmarks is challenging. Existing benchmarks are either domain-agnostic, or limited to a specific domains with high quality datasets. For example, TimeSeriesExam [12] introduced over 700 multiple-choice questions to evaluate five general reasoning skills, but its questions primarily assess signal properties (e.g. trend, cyclicity, stationarity) and lack the contextual depth needed for real-world applications. Domain-specific benchmarks address this gap but have limited scope and poor extensibility, since their curation often relies on templates. For instance, ECG-QA [10] and ECG-Expert-QA [11] focus on ECG interpretation, while EngineMT-QA [27] targets industrial settings. Automatic benchmark generation offers a scalable alternative but raises concerns about quality and diversity of generated questions. Without extensive verification, LLM-generated questions often require heavy manual curation [8, 9], which is both difficult and time-consuming—undermining the main advantage of automation.

Title	Multi-Domain	Curation	# Samples	Skill type		
		Fully Automatic		P	R	PS
Time-MQA [8]	✓	Х	200,000	/	/	✓
TimeSeriesExam [12]	×	X	763	1	/	Х
Time-MMD [9]	✓	×	17,113	1	Х	Х
MT-Bench [7]	✓	✓	22,000	1	/	Х
ECG-QA [10]	×	×	414,348	1	/	Х
TimeSeriesExamAgent (ours)	✓	✓	600+	✓	✓	✓

Table 4: Overview of time series and multimodal datasets with curation and skill types (P – Prediction, R – Reasoning, PS – Practical skills (tasks beyond classification and reasoning such as performing calculations and applying formulas). TimeSeriesExamAgent is universal – tailored to user's needs and with advanced automatic verifications.

Agents for benchmark creation An AI agent is an autonomous system that can observe its environment, reason about possible actions, and act toward achieving a goal. In LLM-based settings, the language model often provides the reasoning or planning layer that guides the agent's decisions. Recent work has shown success in using agents for automatic benchmark creation. Most solutions adopt a multi-agent pipeline with planning, generation, validation, and evaluation modules [14]. For instance, [28] integrates exploratory evaluation using reinforcement learning, while [14] takes a natural language task description as input. However, most of these approaches are not tailored to time series and struggle to generate questions conditioned on numeric data. One recent solution does incorporate time series but is limited to single-step design and lacks extensive verification [29].

B Generation Agent Workflow

We rely on two stages of generation for the templates: planning and generating, inspired by the chain-of-thought (CoT) prompting[30].

Generation planning To provide a relevant and diverse set of templates, we rely on a comprehensive list of domain-specific concepts. There are several ways our pipeline generates a list of concepts:

- 1. LLM generation: User guidelines and dataset descriptions are provided as input to an LLM, which proposes the concepts.
- 2. Web Search: We provide the option for generator LLM obtain concepts through web search.
- 3. Retrieval Augmented Generation: As an option, the user could also provide a relevant file from which the LLM reads and generates concepts[31].

Template generation As input to our generator, the following components are provided:

- User-provided guidelines: a document containing the user's goal or specific requirements,
- Dataset description: a list of columns and example values with ranges from the dataset, with a short usage example,
- List of concepts: generated in previous step. For each template, our pipeline will choose a concept at random to ensure diversity.
- Example templates[Optional]: user-provided few-shot examples presenting required structural elements [32].

B.1 Generation Prompt

280

281

282

283

284

286

287

288

289

290

291

292

293

314

```
Here is the goal of the exam questions:
294
    {user_info_text}
295
296
    Here are sample concepts on which you can base your question generation:
297
    {concept_conversation}
298
299
    Use the concept numbered {concept_no} from the list to guide the design of
300
        your question template.
301
302
303
    Here is the description of the dataset you will use to generate the
304
        question:
    {dataset_describe}
305
306
    In your template, use the provided 'user_dataset' object. Use its 'query(
307
        index)' method to load relevant time series data.
308
309
    Do not select time series randomly. First, formulate the question, and then
310
         choose a time series that fits its logic and reasoning needs.
311
312
    Generate one function-based question template now.
313
```

B.2 Example of Question Template

```
def question_hypertrophic_cardiomyopathy(num_samples, verbose=False):
315
        hyperparameters = {
316
            "min_probability_threshold": 75.0,
317
            "target_abnormalities": ["SEHYP", "LVH", "VCLVH", "RVH"],
318
            "normal_codes": ["NORM"],
319
            "max_attempts": 1000,
320
        }
321
322
        question = "Based on the morphological characteristics and voltage
323
        patterns observed in this ECG recording, what is the most likely
324
        structural cardiac finding?"
325
326
        options = [
327
```

```
"Septal hypertrophy with voltage criteria consistent with
328
        hypertrophic cardiomyopathy",
329
             "Left ventricular hypertrophy with strain pattern indicating
330
        pressure overload",
331
            "Right ventricular hypertrophy suggesting pulmonary hypertension",
332
            "Normal cardiac structure with physiological variant morphology",
333
334
335
        def has_high_probability_abnormality(scp_codes_dict, target_codes,
336
        threshold):
337
            for code in target_codes:
338
                 if code in scp_codes_dict and scp_codes_dict[code] >= threshold:
339
340
                     return code
            return None
342
343
        def is_normal_ecg(scp_codes_dict, normal_codes, threshold):
344
            for code in normal_codes:
345
                 if code in scp_codes_dict and scp_codes_dict[code] >= threshold:
346
347
                     return True
348
            return False
349
350
        qa_pairs = []
351
        attempts = 0
352
        df = user_dataset.get_dataframe()
353
354
        while len(qa_pairs) < num_samples and attempts < hyperparameters["</pre>
355
        max_attempts"]:
356
            attempts += 1
357
            if verbose:
358
                 print(f"[Hypertrophic Cardiomyopathy] Generating question {len(
359
        qa_pairs)+1}/{num_samples}")
360
361
            # Sample a random record
362
            sample_row = df.sample(1).iloc[0]
363
            ecg_id = sample_row['ecg_id']
            scp_codes = sample_row['scp_codes']
365
366
            if not isinstance(scp_codes, dict):
367
                 continue
368
369
            # Check for target abnormalities with high probability
370
371
            detected_abnormality = has_high_probability_abnormality(
                 scp_codes, hyperparameters["target_abnormalities"],
372
373
        hyperparameters["min_probability_threshold"]
374
375
            is_normal = is_normal_ecg(scp_codes, hyperparameters["normal_codes
376
377
        "], hyperparameters["min_probability_threshold"])
378
            if detected_abnormality is None and not is_normal:
379
                 continue
380
381
            try:
382
                 ts = user_dataset.query(ecg_id)
383
                 if ts is None or ts.shape != (12, 1000):
384
                     continue
385
386
```

```
except Exception as e:
387
                 if verbose:
388
                     print(f"Error loading ECG {ecg_id}: {e}")
389
                 continue
390
391
            # Determine correct answer based on detected abnormality
392
393
            if detected_abnormality == "SEHYP":
                 correct_answer = options[0]
394
                 answer_type = "septal_hypertrophy"
395
            elif detected_abnormality == "LVH" or detected_abnormality == "
396
        VCLVH":
397
                 correct_answer = options[1]
398
                 answer_type = "left_ventricular_hypertrophy"
399
            elif detected_abnormality == "RVH":
400
                 correct_answer = options[2]
401
                 answer_type = "right_ventricular_hypertrophy"
402
            elif is_normal:
403
                 correct_answer = options[3]
404
                 answer_type = "normal"
405
            else:
406
407
                 continue
408
            qa_pairs.append({
409
                 "question": question,
410
                 "options": options,
411
                 "answer": correct_answer,
412
                 "answer_type": answer_type,
413
                 "ecg_id": ecg_id,
414
                 "ts": ts.tolist(),
                 "scp_codes": scp_codes,
416
                 "detected_abnormality": detected_abnormality,
417
                 "relevant_concepts": ["Hypertrophic Cardiomyopathy", "
418
        Structural Heart Disease", "ECG Morphology", "Voltage Criteria"],
419
                 "domain": "cardiology",
420
                 "detractor_types": ["Similar structural abnormalities", "Normal
421
         variants"],
422
                 "question_type": "multiple_choice",
423
424
                 "format_hint": "Please answer the question and provide the
        correct option letter, e.g., [A], [B], [C], [D], and option content at
425
        the end of your answer. All information needed to answer the question
426
        is given. If you are unsure, please provide your best guess.",
427
            })
428
429
        return qa_pairs
430
```

B.3 Example of Natural Language Description

```
I want to create time series exam testing model understanding of ecg
432
433
        signals.
434
   To load the data, use the provided user_dataset object.
435
436
   The PTB-XL ECG dataset is a large dataset of 21799 clinical 12-lead ECGs
437
        from 18869 patients
438
   of 10 second length. The raw waveform data was annotated by up to two
439
440
        cardiologists, who assigned
   potentially multiple ECG statements to each record. The in total 71
441
       different ECG statements conform
442
```

```
to the SCP-ECG standard and cover diagnostic, form, and rhythm statements.
443
        The dataset is complemented by extensive metadata on demographics,
444
        infarction characteristics, likelihoods for diagnostic ECG statements
445
        as well as annotated signal properties.
446
447
    You can focus some of the questions around those scp codes:
448
449
   NDT
            non-diagnostic T abnormalities
450
   \mathtt{NST}_-
            non-specific ST changes
451
   DTG
            digitalis-effect
452
   LNGQT
            long QT-interval
453
   NORM
            normal ECG
454
   IMI
            inferior myocardial infarction
455
    ASMI
            anteroseptal myocardial infarction
456
   LVH
            left ventricular hypertrophy
457
   LAFB
            left anterior fascicular block
458
459
    ISC_
            non-specific ischemic
   IRBBB
            incomplete right bundle branch block
460
   1AVB
            first degree AV block
461
   IVCD
            non-specific intraventricular conduction disturbance (block)
462
            ischemic in anterolateral leads
   ISCAL
463
            complete right bundle branch block
   CRBBB
   CLBBB
            complete left bundle branch block
465
            inferolateral myocardial infarction
   TT.MT
466
   LAO/LAE left atrial overload/enlargement
467
   AMI
            anterior myocardial infarction
468
   ALMI
            anterolateral myocardial infarction
469
   ISCIN
            ischemic in inferior leads
470
    INJAS
            subendocardial injury in anteroseptal leads
471
   LMI
            lateral myocardial infarction
472
    ISCIL
            ischemic in inferolateral leads
473
   LPFB
            left posterior fascicular block
474
    ISCAS
            ischemic in anteroseptal leads
475
   INJAL
            subendocardial injury in anterolateral leads
476
   ISCLA
            ischemic in lateral leads
477
   RVH
            right ventricular hypertrophy
478
   ANEUR
            ST-T changes compatible with ventricular aneurysm
480
   RAO/RAE right atrial overload/enlargement
            electrolytic disturbance or drug (former EDIS)
481
            Wolf-Parkinson-White syndrome
482
   ILBBB
            incomplete left bundle branch block
483
   TPI.MT
            inferoposterolateral myocardial infarction
484
   ISCAN
            ischemic in anterior leads
485
   IPMI
            inferoposterior myocardial infarction
486
    SEHYP
            septal hypertrophy
487
488
    INJIN
            subendocardial injury in inferior leads
    INJLA
            subendocardial injury in lateral leads
489
   PMI
            posterior myocardial infarction
490
            third degree AV block
    3AVB
491
492
   INJIL
            subendocardial injury in inferolateral leads
493
  2AVB
            second degree AV block
   ABQRS
494
            abnormal QRS
   PVC
            ventricular premature complex
495
   STD
            non-specific ST depression
496
   VCLVH
            voltage criteria (QRS) for left ventricular hypertrophy
497
   QWAVE
            Q waves present
498
            low amplitude T-waves
   T.OWT
499
   {
m NT}_-
            non-specific T-wave changes
500
   PAC
            atrial premature complex
501
```

```
prolonged PR interval
   LPR
502
   INVT
            inverted T-waves
503
   LVOLT
            low QRS voltages in the frontal and horizontal leads
504
            high QRS voltage
   HVOLT
505
   TAB_
            T-wave abnormality
506
            non-specific ST elevation
   STE_
   PRC(S)
           premature complex(es)
   SR
            sinus rhythm
   AFIB
            atrial fibrillation
510
  STACH
            sinus tachycardia
511
  SARRH
            sinus arrhythmia
512
   SBRAD
            sinus bradycardia
513
514 PACE
            normal functioning artificial pacemaker
   SVARR
            supraventricular arrhythmia
516
   BIGU
            bigeminal pattern (unknown origin, SV or Ventricular)
            atrial flutter
517
   AFLT
   SVTAC
            supraventricular tachycardia
518
   PSVT
            paroxysmal supraventricular tachycardia
519
   TRIGU
            trigeminal pattern (unknown origin, SV or Ventricular)
520
```

521 B.4 Examples of Generated Questions

ECG Question Example

Q: Analyze the P-wave morphology and amplitude characteristics in this recording. What atrial abnormality is present?

- A. RAO/RAE: Right atrial overload/enlargement with prominent P-waves
- B. LAO/LAE: Left atrial overload/enlargement with bifid P-waves
- C. Normal P-wave morphology with no atrial abnormalities
- D. Absent P-waves indicating atrial fibrillation

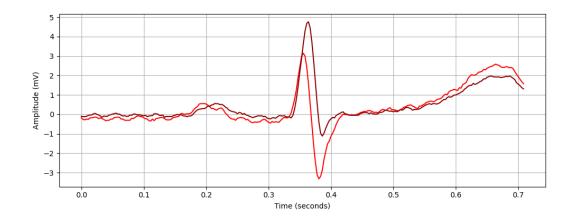




Q: Examine the ECG at shown. Which of the following statements best describes the T-wave morphology in this single-beat ECG?

- A. The T-wave is upright (positive), suggesting normal ventricular repolarization.
- B. The T-wave is inverted (negative), which may indicate myocardial ischemia or ventricular
- C. The T-wave is biphasic (partly positive, partly negative), which may indicate regional repolarization abnormalities.
- D. The T-wave is flattened, which may indicate electrolyte disturbances such as hypokalemia.

523



C LLM Verifier

For each template, we use an LLM to evaluate the generated question. Specifically, we ask:

- Is the question relevant to the given concept?
- Does answering the question require the provided time series?
- Are the question and answer free from ambiguity and bias?

C.1 Validation Prompt

```
You are an expert validator of question templates involving reasoning over
530
    {exam_type} time series data.
531
532
```

You are given an exam question template:

```
{exam_template}
```

534 535

539 540

533

526

527

528

529

Your task is to validate the question template using the following criteria:

- 1. Is the question relevant to {exam_type} time series analysis? 537
- 2. Would you need the time series itself to answer the question? 538
 - 3. Are there no ambiguity in the question or its answer?

If the answer to all is YES or MOSTLY YES, return only the number 1. 541

- If the answer to either is NO, return your objections.
- Return 1 (do not include any additional text then) or describe your objections.

D Other Design Specifics

Detractors In addition, the mechanism of plausible but incorrect answer choices was implemented. The LLM is prompted to reflect on possible mistakes that the test taker might make while solving the exam. Using this knowledge, misleading, incorrect option choices can be generated.

Context Condensation A common issue we encountered in the framework was context window overflow during exam regeneration. To mitigate this, we applied context condensation, which 549 reduces the number of tokens while preserving essential information. In our setup, the agent 550 generates templates in a conversational manner. The process begins with a generation prompt, 551 followed by a message containing the generated exam. If errors occur or the exam is rejected during 552 verification, the feedback and regenerated exams are appended to the conversation. Several context 553 554 condensation techniques exist, such as windowing [33] and context compression [34]. We adopt a 555 summarization-based method [35, 36], which has shown strong results in prior work and fits our use 556 case. Specifically, we summarize non-recent pairs of failing exams and error messages into short descriptions that highlight the issues encountered. These summaries provide the LLM with concise 557 feedback, supporting the generation of higher-quality templates. 558

559 E G-Eval

We evaluated a set of generated questions under the G-Eval framework. We used the following criteria:

```
562 1. SPECIFICITY
```

564

565 566

568

577

580

581

584

586

587

588

589

590

592 593

596

563 Evaluate the specificity of the generated ECG multiple-choice question.

A good question should target a single phenomenon.

567 Evaluation steps:

- 1. Read the question and all answer options.
- 2. Determine if the question targets a single, clearly defined ECG finding or clinical interpretation.
- 3. Assess the ratio of unique medical terms to general words.
- 572 4. Penalize if:
- The question is overly broad or open-ended (e.g., "Is this ECG normal ?").
- The wording leaves diagnostic interpretation unclear.
- The question covers multiple unrelated phenomena.

578 Score highest if the question has one precise focus (e.g., "Is there ST elevation in lead V3?").

1. UNAMBIGUITY

882 Evaluate the unambiguity of the generated ECG multiple-choice question.

583 A question and the answers should not have multiple interpretations.

585 Evaluation steps:

- 1. Read the question and all answer options.
- 2. Determine if the question can be objectively assessed.
- 3. Check if the answers are clear and unambiguous.
- 4. Penalize if:
 - The question uses subjective terms (e.g., "Does this look strange?").
- The answers are open to multiple interpretations.
 - The question cannot be objectively answered.

Score highest if the question is clear and objective (e.g., "Is there tachycardia?"),

2. DOMAIN RELEVANCE

- 598 Evaluate the domain relevance of the generated ECG multiple-choice question.
- 600 Does the question actually pertain to ECGs and medicine?

602 Evaluation steps:

- 603 1. Read the question and all answer options.
- 604 2. Identify medical and ECG-specific terminology.
- 605 3. Determine if the question is relevant to ECG interpretation and medical diagnosis.
- 607 4. Penalize if:
- The question contains non-medical terms (e.g., "Is the line pretty?").
- The question is not related to ECG interpretation.
 - The question lacks medical context.

610 611

614

619

623

626

599

601

612 Score highest if the question contains relevant medical terms

613 (e.g., "QRS," "arrhythmia," "P wave") and pertains to ECG interpretation.

615 3. ANSWERABILITY

616 Evaluate the answerability of the generated ECG multiple-choice question.

Even without an answer provided, the question should be answerable based on the data (ECG).

620 Evaluation steps:

- 621 1. Read the question and all answer options.
- 622 2. Determine if the question can be answered by analyzing ECG waveform data.
- 624 3. Assess whether the question requires time series analysis or could be 625 answered without it.
 - 4. Penalize if:
- The question asks about non-ECG factors (e.g., "Was the patient nervous?").
- The question can be answered without analyzing the ECG time series data.
 - The question is too general and doesn't require specific ECG analysis.

631 632

 $_{633}$ Score highest if the question requires specific ECG time series analysis $_{634}$ (e.g., "Is there atrial fibrillation?").

635 Give fewer points if the question can be answered without time series data.

536 F Hyperparameters

In this section, we list all the hyperparameter used for our agentic workflow.

- 1. Generator LLM: the LLM use to generate concepts and the corresponding template. We used claude-sonnet-4-20250514 (initial generation with reasoning_effort="medium").
- 2. Concept LLM: the LLM use to generate concepts. We used gpt-4o-2024-08-06.
- 3. Verifier LLM: the LLM use to verify templates. We used gpt-4o-2024-08-06.
- 4. Student LLMs: the student LLMs we use to check the exam differentiability. Currently we have two student LLMs: stronger: gpt-4o-2024-08-06 and weaker: gpt-4o-mini.
- 5. Exam type: We are generating the data connected to specific domain. We used "ECG".
- 6. Few-shot examples: 3 templates prepared beforehand were used to present the desired structure. For each generation, they were randomly sampled from set of 9.

647 G Evaluation Protocol

- All used models were accessed by API with LiteLLM Python library. The following API providers were used with default parameters:
 - Closed source models OpenAI API, Anthropic API
 - Open source models Hugging Face Inference Providers API
- During the evaluation, the images of the plots were encoded with base64 encoding and provided to the models. Plots were created with DPI = 50. We used setup without context condensation.

654 H Expert evaluation

650

651

662

663

664

665

666

667

- We also presented samples of our work to the clinicians. During the first meeting, we received feedback that our work was interesting but required further improvements. Specifically, the plots needed to be both stretched, annotated and all 12 leads needed to be included. In addition, the language and jargon we used could be confusing for specialists. Overall, the clinicians considered 3 out of 5 questions to be answerable and flagged 2 out of 5 as problematic.
- At the following meeting, after incorporating experts' comments into the prompt, we asked a specialist to provide their opinion on improvements on the exams, using the following criteria:
 - Correctness: The answer must be unequivocally accurate according to current medical knowledge and guidelines.
 - No Ambiguity: Only one answer should be valid; distractors must be plausible but clearly incorrect.
 - Precision of Wording: Both questions and answers should be clear, concise, and medically
 accurate, avoiding vague phrasing.
 - Relevance: The question should be engaging and meaningful to the specialist.
- Of the 8 questions evaluated, 7 were rated positively across all four criteria: correctness, lack of ambiguity, precision of wording, and relevance.