OPENTSLM: TIME-SERIES LANGUAGE MODELS FOR REASONING OVER MULTIVARIATE MEDICAL TEXTAND TIME-SERIES DATA

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

028

029

031

034

039 040 041

042

043

044

045

047

048

051

052

ABSTRACT

Large language models (LLMs) have emerged as powerful tools for interpreting multimodal data (e.g., images, audio, text), often surpassing specialized models. In medicine, they hold particular promise for synthesizing large volumes of clinical information into actionable insights and patient-facing digital health applications. Yet, a major limitation remains their inability to handle time series data. To overcome this gap, we present OpenTSLM, a family of Time-Series Language Models (TSLMs) created by integrating time series as a native modality to pretrained LLMs, enabling natural-language prompting and reasoning over multiple time series of any length. We investigate two architectures that differ in how they model time series. The first, OpenTSLM-SoftPrompt, models time series implicitly by concatenating learnable time series tokens with text tokens via soft prompting. Although parameter-efficient, we hypothesize that explicit time series modeling scales better and outperforms implicit approaches. We thus introduce OpenTSLM-Flamingo, which integrates time series with text via cross-attention. We benchmark both variants with LLaMa and Gemma backbones against baselines that treat time series as text tokens or plots, across a suite of text-time-series reasoning tasks. We introduce three time-series Chain-of-Thought (CoT) datasets: HAR-CoT (human activity recognition), Sleep-CoT (sleep staging), and ECG-QA-CoT (ECG question answering). Across all, OpenTSLM models consistently outperform baselines, reaching 69.9% F1 in sleep staging and 65.4% in HAR, compared to 9.05% and 52.2% for finetuned text-only models. Notably, even 1Bparameter OpenTSLM models surpass GPT-40 (15.47 and 2.95%). OpenTSLM-Flamingo matches OpenTSLM-SoftPrompt in performance and outperforms on longer sequences, while maintaining stable memory requirements. By contrast, SoftPrompt exhibits exponential memory growth with sequence length, requiring 110 GB compared to 40 GB VRAM when training on ECG-QA with LLaMA-3B. Expert reviews by clinicians find strong reasoning capabilities and temporal understanding of raw sensor data exhibited by OpenTSLMs on ECG-QA. To facilitate further research, we provide all code, datasets, and models open-source.

1 Introduction

Medicine is inherently temporal: assessment, diagnosis, and treatment depend on how signs, symptoms, and biomarkers evolve over time Giannoula et al. (2018); Henly et al. (2011); Jørgensen et al. (2024). Clinical decision-making relies on temporal patterns—tracking vital signs, medication responses, laboratory values, and disease progression markers to guide diagnosis, prognosis, and therapeutic interventions. As time-series data from electronic health records and continuous monitoring proliferate Abernethy et al. (2022); Marra et al. (2024); Yeung et al. (2023), human-legible representations become essential for interpreting and managing this information Olex & Mcinnes (2021); Senathirajah et al. (2020); Zhou et al. (2008). Clinical summaries must translate complex temporal patterns—hemodynamic instability, biomarker trajectories, and treatment responses—into interpretable assessments that support evidence-based decision-making and care coordination.

Recent advances in multimodal large language models (LLMs) allow users to interpret complex data through natural language, synthesizing information across text, images, audio, and video Wu et al. (2023); AlSaad et al. (2024). However, reasoning over longitudinal time series data remains

055

056

057

058

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

079

081

083 084

085

087

090

092

093

095

096

097

098

099

100 101

102

103

104

105

106

107

a critical blind spot among currently supported modalities. Prior work has attempted to integrate time-series as plain text tokens Gruver et al. (2023); Kim et al. (2024); Liu et al. (2023); however results have been limited Merrill et al. (2024). Other approaches reprogram LLMs to act as feature extractors for classification heads, which then output a fixed set of classes or values, thereby losing text-generation capabilities Li et al. (2025); Nie et al. (2023); Pillai et al. (2025); Ye et al. (2025). More recently, soft prompting has been explored, concatenating learnable time-series tokens with text tokens to preserve generation Chow et al. (2024). Yet, longer series may require more tokens, increasing context length Götz et al. (2025); Nie et al. (2023) and compute due to the quadratic cost of self-attention Nie et al. (2023); Vaswani et al. (2017).

To overcome prior limitations, we propose Time-Series Language Models (TSLMs), which integrate time series as a native modality in LLMs. TSLMs provide a natural interface to complex medical data, enabling clinicians and patients to query, interpret, and reason about longitudinal health information directly through natural language. We introduce OpenTSLM, a family of TSLMs built by extending pretrained LLMs with time-series inputs. A central design question in building TSLMs is how to represent time-series signals. Prior work has primarily used soft prompting, encoding time series as learned token embeddings concatenated with text tokens. While lightweight, this captures temporal dependencies only implicitly, as additional tokens in the context, and may scale poorly to longer or multiple sequences. We hypothesize that explicit multimodal fusion via cross-attention may be more effective for modeling temporal structure. To compare both approaches, we explore two variants for OpenTSLM. The first, OpenTSLM-SoftPrompt, models time series implicitly by encoding the time series into tokens and concatenating them with text tokens via soft prompting, so the model processes both as a single sequence without distinguishing between them. The second, **OpenTSLM-Flamingo**, by contrast, models time series explicitly as a separate modality, using a cross-attention mechanism inspired by Flamingo Alayrac et al. (2022) to fuse time-series and text. We created OpenTSLM-SoftPrompt and OpenTSLM-Flamingo using Llama Touvron et al. (2023) and Gemma GemmaTeam et al. (2024) backbones. We benchmark these models against each other and against baselines including LLMs with tokenized time-series inputs Gruver et al. (2023), fine-tuned tokenized time-series models, and vision-based approaches. Unlike prior classificationbased approaches, our models are trained in text-based reasoning tasks, generating chain of thought (CoT) rationales before producing predictions. For training and evaluation, we introduce three new datasets: HAR-CoT, Sleep-CoT, and ECG-QA-CoT. To foster reproducibility and further research on TSLMs, we release OpenTSLM as an open-source framework, including models and datasets¹.

2 RELATED WORK

Table 1: Methods combining time-series data with LLMs.

				Gen.	l-Senson Raw	Oata
Name	Method	Task	Jest	Mill	r Ran	, Ek
FSHLLiu et al. (2023)	Token	CL	~	~	~	
Gruver et al. (2023)	Token	FC	~		~	
HealthLLM Kim et al. (2024)	Token	TR	~	~	~	~
Chow et al. (2024)	Soft Prom.	TR	~	~	~	~
MedualTime Ye et al. (2025)	Soft Prom.	CL			~	~
SensorLLM Li et al. (2025)	Soft Prom.	CL		~	~	~
Time2Lang Pillai et al. (2025)	Soft Prom.	CL			~	
OpenTSLM-SP (ours)	Soft Prom.	TR	~	~	~	
SensorLM Zhang et al. (2025)	Cross.Attn.	CL	~	~		
OpenTSLM-Flamingo (ours)	Cross.Attn.	TR	~	~	~	~
CL =Classification, F	C =Forecasting	TR =To	ext Rea	soning	5	

Creating Time-Series Language Models remains an open research challenge. The main barrier is the modality gap between continuous signals and discrete text representations Chow et al. (2024); Pillai et al. (2025); Zhang et al. (2025). Prior work has proposed three main strategies to bridge this gap, as summarized by Zhang et al. (2024): tokenizing time series as text (Section 2.1), applying soft prompting (Section 2.2), and using cross-attention mechanisms (Section 2.3). Table 1 provides an overview of relevant methods.

2.1 TOKENIZATION OF TIME SERIES AS TEXT INPUTS

Gruver et al. has demonstrated that LLMs can perform time series forecasting by encoding values as text tokens and predicting future values without domain-specific tuning Gruver et al. (2023).Liu et al. (2023) tokenize data from wearables and smartphones to enable LLMs to infer clinical and

¹Link to Github repo hidden during double-blind review; code is provided as a zip file; code includes data loaders that fetch our datasets from an anonymous link. More details in README inside the zip file.

wellness information through few-shot prompting. Similarly, Kim et al. (2024) propose HealthLLM, a framework for health prediction using physiological signals (e.g., heart rate, sleep) combined with user context and medical knowledge embedded in prompts.

2.2 COMBINING TEXT AND TIME SERIES TOKEN EMBEDDINGS (SOFT PROMPTING)

An alternative to manual tokenization is to encode time series into embeddings that capture time series information, using a time series encoder as presented by Nie et al. (2023). These embeddings can be input into a transformer directly or concatenated with text embeddings (softprompting) Chow et al. (2024); Nie et al. (2023); Pillai et al. (2025); Ye et al. (2025). Pillai et al. (2025) use this approach and train an encoder to produce soft prompts from time series, which are then processed by a frozen LLM for classification via a projection head; however this disables free-form text generation. Ye et al. (2025) similarly combines time series and text-token embeddings, using a classification head for prediction. Li et al. (2025) integrate sensor and text embeddings in two stages: First, they generate a caption-like summary of the time series for free-form output; Second, they classify the data via a projection head, therefore restricting free from output. Chow et al. (2024) interleave time series tokens with text tokens in the LLM input, enabling free-form text reasoning.

2.3 Cross-Attention for Time-Series Data

Few studies use cross-attention to integrate time series into LLMs. Zhang et al. (2025) apply cross-attention between a time series encoder and a text encoder, aligned with contrastive loss, to extract statistical summaries (e.g., mean, max) from a single sensor. They train a new sensor encoder, text encoder, and multimodal text decoder, rather than adapting a pretrained LLM Zhang et al. (2025).

3 METHODS

 We present two architectures for TSLMs, OpenTSLM-Soft Prompting (SP) (Section 3.2 and OpenTSLM-Flamingo (Section 3.3). To support multiple time-series inputs, we design a prompt format that interleaves sensor data with accompanying textual descriptions (e.g., "Data from Sensor X over Y days:" followed by the data representation). Figure 1 illustrates our approach.

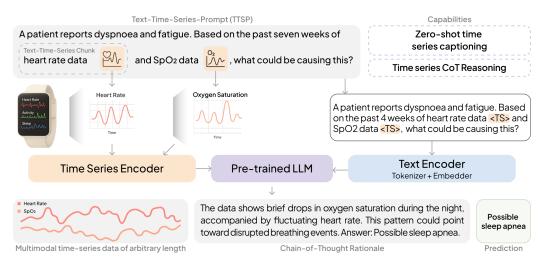


Figure 1: Overview of Text-Time-Series LLMs with support for multiple time-series inputs.

3.1 TIME-SERIES ENCODER

Both OpenTSLM architectures use a time series encoder inspired by Nie et al. (2023). It consists of a Patchencoder, followed by either a TransformerEncoder for OpenTSLM-SP or a PerceiverResampler for OpenTSLM-Flamingo (inspired by Alayrac et al. (2022); Awadalla et al. (2023)). We divide an input time series $x \in \mathbb{R}^L$ into non-overlapping patches of size p, yielding N = L/p patches. Each patch is then transformed into an embedding vector using a 1D convolution and added with a positional encoding Nie et al. (2023)

Patch Embedding:
$$\mathbf{E}_i = \text{Conv1D}(x_{i \cdot p : (i+1) \cdot p}) \in \mathbb{R}^{d_{\text{enc}}} + \mathbf{P}_i$$
 (1)

where the convolution has kernel size and stride equal to p, mapping each patch to a $d_{\rm enc}$ -dimensional embedding. P_i is the learnable positional encoding. The sequence of position-augmented embeddings is then processed by the specific Encoder (cf. Sections 3.2 and 3.3).

Preserving scale and temporal information The PatchEncoder expects inputs normalized to $x \in [-1,1]$. Since raw time series differ in scale and resolution across modalities depending on the sensor, we preserve scale and temporal context by adding the original mean, standard deviation, and time scale to the textual description. For example:

This is heart-rate data over 24 hours sampled at 50 Hz with mean=61 and std=12.

3.2 SOFT-PROMPTING ARCHITECTURE (OPENTSLM-SP)

OpenTSLM-SP has three components: (1) a time series encoder that transforms raw data into patch embeddings, (2) a projection layer mapping embeddings to the LLM hidden space, (3) a pretrained LLM, fine-tuned using LoRA adapters Hu et al. (2021) Figure 2 illustrates the architecture.

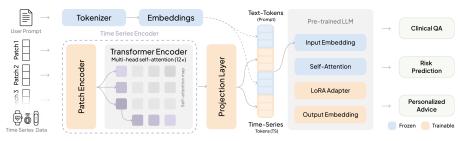


Figure 2: Architecture of OpenTSLM-SoftPrompt

Projecting Time-Series Tokens to Text Tokens We apply the patch embeddings to a transformer encoder and subsequently project the resulting tokens with an multi-layer perceptron (MLP) to align them with the embedding space of dimension $d_{\rm llm}$, corresponding to the hidden size of the language model, following Nie et al. (2023) and Chow et al. (2024).

$$\mathbf{Z} = \text{MLP}(\mathbf{TransformerEncoder}(E_{1:N})) \in \mathbb{R}^{N \times d_{\text{llm}}}$$
 (2)

where $\mathbf{Z} \in \mathbb{R}^{N \times d_{\text{llm}}}$ denotes the projected time-series tokens in the LLM embedding space.

Text-Time-Series integration via Soft Prompting We interleave any number of text and time-series tokens through a soft prompting mechanism. A typical prompt consists of (1) an initial text segment ("pre-prompt"), (2) a sequence of interleaved time-series tokens and textual descriptions, and (3) a final text segment ("post-prompt"), often a question. Formally, the model input is:

$$\mathbf{X}_{input} = [\mathbf{T}_{pre}, \mathbf{Z}_1, \mathbf{T}_{desc_1}, \mathbf{Z}_2, \mathbf{T}_{desc_2}, \dots, \mathbf{Z}_K, \mathbf{T}_{desc_K}, \mathbf{T}_{post}]$$
(3)

where T_{pre} , T_{desc_i} , and T_{post} are token embeddings of text segments, and each Z_i is a projected time-series embedding aligned with the LLM hidden space. We refer to each (Z_i, T_{desc_i}) as a text-time-series chunk. This approach implicitly integrates time series through learned tokens.

3.3 Cross-attention architecture (OpenTSLM-Flamingo)

OpenTSLM-Flamingo is inspired by the Flamingo model for vision–language tasks Alayrac et al. (2022); Awadalla et al. (2023). Following OpenFlamingo Awadalla et al. (2023), we extend pretrained LLMs with cross-attention layers to support time-series reasoning.

Architecture Overview We replace the vision encoder of Flamingo with a time series encoder and adapt the cross-attention mechanism for temporal data. The model consists of: (1) a time series patch encoder, (2) a Perceiver Resampler, (3) gated cross-attention layers integrated into the LLM, and (4) the frozen language model backbone. Figure 3 visualizes the architecture.

PerceiverResampler We use a PerceiverResampler inspired by Flamingo Awadalla et al. (2023) as Encoder for the time series patches, yielding a fixed-size latent representation:

$$\mathbf{Z}_{\text{latent}} = \text{PerceiverResampler}(\mathbf{E}_{1:N}) \in \mathbb{R}^{N_{\text{latent}} \times d_{\text{time}}},$$
 (4)

Here, d_{time} is the dimensionality of the time-series features by the perceiver, in our case (N, 1), encoding one time series with one channel at a time.

217

219

220

221

222

223 224 225

226 227

228

229

230

231

232 233

234

235

236

237

238

239

240 241

242

243

244

245 246

247 248

249

250

251

253

254

255

256

257

258

259

260

261

262

264

265

266

267

268

Text-Time-Series Gated Cross-Attention To integrate **Z**_{latent} into the LLM, we add gated crossattention layers every N (hyperparameter) transformer blocks which compute:

$$\mathbf{Q}_{\text{text}} = \mathbf{x} \mathbf{W}_{Q}, \quad \mathbf{K}_{\text{ts}} = \mathbf{Z}_{\text{latent}} \mathbf{W}_{K}, \quad \mathbf{V}_{\text{ts}} = \mathbf{Z}_{\text{latent}} \mathbf{W}_{V}$$
 (5)

GatedCrossAttention(
$$\mathbf{x}, \mathbf{Z}_{latent}$$
) = $x + \gamma \cdot softmax \left(\frac{\mathbf{Q}_{text} \mathbf{K}_{ts}^T}{\sqrt{d_k}}\right) \mathbf{V}_{ts}$. (6)

where γ_{attn} is a learnable parameter controlling the influence of the time-series, $\mathbf{x} \in \mathbb{R}^{T \times d_{\text{model}}}$, the LLM input, \mathbf{W}_O , \mathbf{W}_K , $\mathbf{W}_V \in \mathbb{R}^{d_{\text{model}} \times d_k}$ learned projection matrices, and d_k the key dimension.

Conditioning Text-tokens on Time-Series via Special Tokens The LLM processes tokens autoregressively, attending to previous inputs. Following OpenFlamingo Awadalla et al. (2023), we introduce special tokens $\langle TS \rangle$ and \langle endofchunk \rangle to indicate when time series modalities should be incorporated. Upon encountering (TS), the model conditions on the corresponding latent representation \mathbf{Z}_{latent} via gated cross-attention. A typical input prompt is

$$\mathbf{X}_{input} = [pre_prompt, \langle TS \rangle, ts_desc_1, \langle endofchunk \rangle, \langle TS \rangle, ts_desc_2, \langle endofchunk \rangle, post_prompt] \tag{7}$$

where $\langle TS \rangle$ triggers multimodal conditioning and \langle endofchunk \rangle signals the end of text describing a time series. This setup enables interleaving multiple text and time series segments Awadalla et al. (2023). The embeddings of the special tokens are learned during training.

4 **EXPERIMENTS**

In the following, we outline our training methodology and report results on multiple-choice Time Series Question Answering (TSQA) and time-series reasoning datasets. We compare OpenTSLM-SoftPrompt and OpenTSLM-Flamingo against each other and baselines in terms of performance, and report video random access memory (VRAM) requirements for training OpenTSLM. We present sample model outputs across datasets and an evaluation for ECG rationales by medical doctors.

4.1 Multi-Stage Curriculum Learning – Teaching LLMs Time Series

Following Chow et al. (2024), we adopt a two-stage curriculum to train TSLMs. In stage one (encoder warmup), we use two synthetic time-series datasets to pretrain the encoder:

- TSQA Wang et al. (2024) Multiple-choice time-series question answering on synthetic data for learning simple temporal patterns (e.g., ascending/descending trends).
- Time-Series Captioning (M4-Captions) We generate pseudo-labeled captions using ChatGPT, prompted with M4 time-series plots (see Section A.4.1).

In stage two, we introduce three new CoT time-series datasets covering human activity recognition (HAR), sleep staging, and electrocardiogram (ECG) Question Answering (QA). We generated these using GPT-40 by providing a plot and ground-truth answer for each sample, then asking the model to produce rationales leading to the correct response. Further details are provided in Section A.2.

- HAR-CoT three-axis accelerometer data combined from DaLiAc Leutheuser et al. (2013), DOMINO Arrotta et al. (2023), HHAR Stisen et al. (2015), PAMAP2 Reiss & Stricker (2012), RealWorld Sztyler & Stuckenschmidt (2016), and datasets from Shoaib et al. (2013; 2014; 2016). Sampled at 50 Hz, split into 2.56s windows, 8 activities: sitting, standing, lying, walking, running, biking, walking upstairs, walking downstairs. See Section A.2.1 for detailed description.
- Sleep-CoT Based on SleepEDF Kemp et al. (2000); Goldberger et al. (2000), using 30s electroencephalogram (EEG) segments for sleep staging. Following prior work Chow et al. (2024);

Pouliou et al. (2025), Non-rapid eye movement (REM) stages 3 and 4 are merged, yielding five classes: Wake, REM, Non-REM1, Non-REM2, Non-REM3. See Section A.2.2 for details.

• ECG-QA-CoT Based on ECG-QA Oh et al. (2023), which provides 12-lead 10s ECGs and clinical context, we excluded comparison questions, retaining 42/70 templates. This yielded 3,138 unique questions across 240k samples (see Section A.2.3).

All datasets are split into **80/10/10 train/validation/test** sets. Table 3 in Section A.1 summarizes number of samples in the datasets, number of time series and lengths.

Training objective In all stages, we frame the task as an autoregressive language modeling problem. During training and evaluation, the model is prompted to generate outputs in a structured format, consisting of a free-form rationale followed by the final prediction: `'<reasoning> Answer: <final answer>''. Formally, the loss is defined by Equation 8, where \mathbf{Z}_{ts} are the

$$\mathcal{L}_{LM} = -\sum_{t=1}^{T} \log P(y_t \mid y_{< t}, \mathbf{x}_{1:t}, \mathbf{Z}_{ts}; \Theta)$$
(8)

time-series features, and Θ the learnable weights, i.e., the TimeSeriesEncoder, MLP, and LoRA in OpenTSLM-SoftPrompt, and TimeSeriesEncoder and cross-attention in OpenTSLM-Flamingo.

4.2 BASELINES

We compare OpenTSLM against three baselines using the same open-weight LLMs, i.e., Llama-3.2(1B, 3B) and Gemma3 (270M, 1B-PT), and additionally ChatGPT-40 (gpt-4o-2024-08-06).

- 1. **Tokenized time-series**: Using the open-source code provided by Gruver et al. (2023), we tokenize time series into text inputs and report zero-shot performance on the test set.
- 2. **Tokenized finetuned**: Same as 1. (excluding GPT-4o), but finetuned with LoRA Hu et al. (2021) on the training set. We choose best model by validation loss, and report performance on test set.
- 3. **Image (Plot)**: We convert time series into plots and provide them as input to GPT-40 and Gemma-4b-pt (since the smaller Gemma 3 variants do not support image input).
- 4. **Random baseline**: For comparison, we report the expected performance of a predictor that selects labels uniformly at random, adjusted to each dataset's label distribution.

4.3 QUANTITATIVE RESULTS ON TIME-SERIES CLASSIFICATION

We present performance on the test splits of TSQA, HAR-CoT, Sleep-CoT, and ECG-QA-CoT and report macro-F1 score and accuracy in Table 2. OpenTSLM models achieve the highest performance across benchmarks, while most tokenized text-only baselines fail to produce valid outputs, not answering in the expected template but merely repeating inputs or starting to count (see Section A.3), resulting in 0.00 F1 on HAR for all models except for GPT-40 (2.95). GPT-40 yields only 2.95 F1 with text but improves substantially with plots (e.g., 10.83 on HAR, 59.24 on TSQA). Gemma3-4b similarly achieves better results TSQA and Sleep-CoT (48.77 and 6.75). Llama models achieve 2.14 and 5.65F1 on Sleep, respectively, while Gemma models again achieve 0.00, likely due to their smaller context window (32k vs. 128k). By contrast, OpenTSLM-SoftPrompt with Llama3.2-1B attains 97.50 F1 score (97.54 accuracy) on TSQA, with Llama3.2-3B at 97.37 (97.33); Flamingo variants are close (e.g., Llama3.2-1B 94.08 (94.00)), while the strongest tokenized-finetuned baseline reaches 84.54 (82.06) and GPT-40 with image inputs at 59.24 (62.10). On HAR-CoT, the strongest results are 65.44F1 (71.48 accuracy) for OpenTSLM-SoftPrompt (Llama3.2-1B) and 65.44 (71.48) for OpenTSLM-Flamingo (Gemma3-1B-pt); the best tokenized-finetuned baseline records 60.44 (66.87). On Sleep-CoT, OpenTSLM-SoftPrompt (Llama3.2-1B) achieves 69.88 (81.08), followed by OpenTSLM-SoftPrompt (Llama3.2-3B) at 54.40 (72.04) and Flamingo (Gemma3-270M) at 51.38 (68.49); tokenized-finetuned baselines remain lower (best 9.05 (24.19)). On ECG-QA-CoT, OpenTSLM-Flamingo (Llama3.2-3B) leads with 40.25 (46.25).

4.4 EVALUATION OF MEMORY USE DURING TRAINING

We evaluate peak VRAM usage during training for both OpenTSLM variants. Figure 4 summarizes peak VRAM on TSQA, HAR-CoT, SleepEDF-CoT, and ECG-QA-CoT. OpenTSLM-Flamingo shows near-constant memory across datasets: Llama-3.2-1B requires around 20–22 GB and Llama-3.2-3B around 61–72 GB; Gemma-3-270M is 5.7–7.3 GB and Gemma-3-1B-pt 15.6–18.4 GB. In contrast, OpenTSLM-SoftPrompt vary substantially with the dataset: Llama-3.2-1B requires from 4.4 GB (TSQA) up to 64.9 GB (ECG-QA-CoT), and for Llama-3.2-3B from 8.1 GB to 87.1 GB; Gemma-3-270M spans 2.4–24.1 GB and Gemma-3-1B-pt 5.1–32.7 GB.

Table 2: Performance comparison on time series question answering (TSQA) and time series reasoning (HAR-CoT, Sleep-CoT, ECG-QA-CoT) tasks between OpenTSLM models and baselines.

Method	lethod Model		TSQA HAR-CoT		R-CoT	Sleep-CoT		ECG-	QA-CoT
		F1	Acc	F1	Acc	F1	Acc	F1	Acc
Random B	aseline	33.33	33.33	11.49	12.50	17.48	20.00	16.47	20.18
p s	Llama3.2-1B	16.01	31.04	0.00^{*1}	0.00	2.14	0.65	0.00	0.00
Tokenized Time-Series	Llama3.2-3B	16.24	32.06	0.00	0.00	5.66	12.15	0.00	0.00
eni S-S	Gemma3-270M	10.52	9.58	0.00	0.00	0.00	0.00	0.00	0.00
ž,ŭ	Gemma3-1B-pt	11.76	12.92	0.00	0.00	0.00	0.00	0.00	0.00
	ChatGPT-4o	45.32	45.29	2.95	11.74	15.47	16.02	18.19	28.76
р p	Llama3.2-1B	83.74	81.40	51.28	62.71	9.05	24.19	OOM*2	OOM
uné uné	Llama3.2-3B	84.54	82.06	60.44	66.87	5.86	14.30	OOM	OOM
Tokenized Finetuned	Gemma3-270M	68.05	65.40	40.66	54.56	0.00	0.00	OOM	OOM
F 표	Gemma3-1B-pt	82.85	83.42	52.15	63.90	0.00	0.00	OOM	OOM
ge ot)	Gemma3-4B-pt	48.77	50.60	1.72	0.89	6.75	14.95	1.90	1.03
Image (Plot)	ChatGPT-40	59.24	62.10	10.83	13.90	4.82	10.75	24.95	33.30
∑ td.	Llama3.2-1B	97.50	97.54	65.44	71.48	69.88	81.08	32.84	35.49
SL.	Llama3.2-3B	97.37	97.33	64.87	67.89	54.40	72.04	33.67	36.25
든품	Gemma3-270M	40.32	26.79	1.43	0.55	7.96	5.91	1.29	1.11
OpenTSLM SoftPrompt	Gemma3-1B-pt	87.29	89.18	40.52	45.17	30.99	36.56	27.86	34.76
Σ o	Llama3.2-1B	94.08	94.00	62.93	69.27	49.33	67.31	34.62	38.14
SL	Llama3.2-3B	90.14	90.10	62.77	69.03	45.45	69.14	40.25	46.25
OpenTSLM Flamingo	Gemma3-270M	77.86	78.12	57.75	63.43	51.38	68.49	32.71	35.50
Pe FI2	Gemma3-1B-pt	92.56	92.46	65.44	71.48	43.69	60.67	35.31	37.79

Note: Gemma models have smaller context than Llama (32k vs. 128k); softprompt uses up context, performing worse. *10.00 model failed to produce "Answer: {answer}" template, often repeating input prompt (see Section A.3).*2OOM - Out of memory: 12 ECG leads of 10s tokenize to 80k tokens, requiring >100GB VRAM.

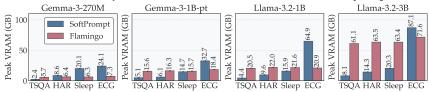


Figure 4: VRAM memory usage in training across datasets.

To further investigate memory scaling, we train models on a simulated dataset (see Section A.6.2) with random inputs of shape $(N \times L)$, where N is the number of time series processed concurrently and L the sequence length. We report max VRAM usage in Figure 5 (exact values are available in Table 10). VRAM for OpenTSLM-Flamingo effectively stays constant as N increases

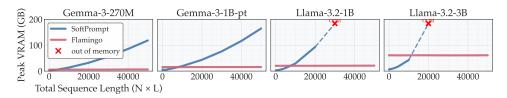


Figure 5: VRAM usage vs. total time-series size $N \times L$ (number of series × length)

from 1 to 5 and L from 10 to 10,000 (e.g., Llama-1B \approx 20.4–21.0 GB; Llama-3B \approx 60.7–61.1 GB; Gemma-270M \approx 5.7–6.4 GB; Gemma-1B \approx 15.4–15.6 GB). By contrast, SoftPrompt scales with both N and L (see Figure 5 in Section A.6.2): for Llama-1B, VRAM rises from \sim 2.6 GB at L=10, N=1 to \sim 29.5 GB at L=10,000, N=1 and exceeds memory at L=10,000, N≥3; Llama-3B shows a similar pattern (6.3 GB \rightarrow 42.7 GB at N=1, OOM by N≥3). Gemma-270M and Gemma-1B reach up to \sim 118 GB and \sim 165 GB, respectively, at L=10,000, N=5.

4.5 QUALITATIVE RESULTS AND EXPERT EVALUATION OF ECG RATIONALES

Both **OpenTSLM** variants remain text models, trained to generate rationales for classification rather than outputting only a class label. Figure 6a shows example rationales for *human activity recognition*, Figure 6b for ECG-QA, and Figure 6c for sleep staging. Fig-

379 380

381

382

383 384

385

386 387

388

389 390

391

392

393

396 397

399 400

401

402

403

404

405

406 407

408

409 410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

Prompt
Given the following accelerometer data, classify the activity. Possible activities:
blking, lying, running, sitting, standing, walking, up, walking_up, walking_down
This is the x-axis sts. This is the y-axis sts. This is the z-axis sts. /hat range does the p duration of this 12-lead ECG <TS>, ..., <TS> fall in? Answers: 1. above the normal range 2. below the normal range 3. within the normal range 3. within the normal Clinical Context: 74-year-old female patient. 12-lead ECG. clinical recording, recorded with AT-60 3. Signal quality; baseline Reason step by step, taking into account all three axis. Finish with 'Answer burst noise present, electrode artifacts present. extra beats detected, pacemaker present *Rationale
To determine the P wave duration, we must first consider the standard normal
range, which is typically 0.08 to 0.11 seconds. In this ECG, despite the
presence of baseline drift, static noise, burst noise, and electrode artifacts,
the P waves can still be evaluated. The presence of a pacemaker may • Rationale
The accelerometer data shows significant fluctuations across all three axes, indicating substantial movement. The X-axis exhibits sharp peaks and troughs suggesting rapid changes in direction, which are characteristic of dynamic activities. Similarly, the Y-axis displays pronounced oscillations, further supporting the presence of vigorous movement. The Z-axis also demonstrates consistent valiability, indicating vertical motion. The frequency and amplitude of these variations are consistent with an activity involving continuous and repetitive motion. The magnitude of changes across all axes suggests a level of intensity and movement that aligns with more dynamic activities (Answer: running). influence the rhythm, but it does not directly affect the intrinsic P wave initiatine the Injurin, but it does not allevity affect the initiatis. P wades appear consistent across the leads, and even with the noted artifacts, they do not seem prolonged. Given the patient's age and the presence of a pacemaker, it is crucial to assess whether any conduction abnormalities could affect the P wave duration. However, the morphology and duration of the P waves do not suggest any significant prolongation Therefore, considering the observed patterns and clinical context, the P wave duration falls within the normal range. Answer: within the normal range (a) Human Activity Recognition CoT (b) ECG Interpretation CoT Prompt Given the following EEG data, classify the sleep stage. The possible stages are Wake, REM, Non-REM Stage 1, Non-REM Stage 2, Non-REM Stage 3, Movemer This is the EEG recording ISS. Reason step by step. Finish with 'Answer.' Generate a detailed caption for this time-series data <TS> The time-series graph displays data points over a period, showcasing fluctuations between values of approximately 1300 and 1900. The series begins with moderate variability, experiencing a sharp increase around the 10th data The observed EEG signal displays a pattern characterized by relatively high frequency and variability, with rapid fluctuations and a lack of consistent point, reaching a peak near 1900. Following this peak, the data exhibits a rhythmic oscillations. This suggests a state of heightened neural activity, typically associated with a more alert and conscious condition. The presence of such dynamic and irregular patterns, along with the absence of the slower, more synchronized waveforms often found in deeper sleep downward trend with intermittent spikes and drops. Notably, there are significant drops around the 50th and 80th data points, where the values dip close to 1300. The latter part the series shows continued volatility with smaller peaks and troughs, indicating persistent fluctuations. The overall pattern stages, points towards a classification of wakefulness. Answer: Wake suggests a high degree of variability, with no clear long-term trend

ure 6d shows time series captioning on M4 data. To evaluate the quality of model ratio-

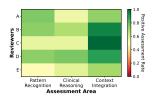
Figure 6: Example CoT rationales for HAR, Sleep Staging, ECG-QA and M4 captioning, generated with OpenTSLM-Flamingo/Llama3.2-1B. More examples are provided in Section A.4.

nales, we conducted an expert review with five cardiologists from ANONYMIZED Hospital on rationales generated by OpenTSLM-Flamingo-Llama3.2-3B (best model) for ECG-QA.

We randomly sampled two examples per template (84 total), each reviewed by at least two cardiologists. Evaluation followed a rubric derived from the American College of Cardiology/American Heart Association Clinical Competence Statement on ECGs Pangaro (1999); Committee Members et al. (2001) and based on the RIME ("Reporter–Interpreter–Manager–Educator") framework Pangaro (1999) (see A.5), assessing whether the model: (1) correctly

(c) Sleep Stage Detection CoT

80 60 60 85% 62% 85% 85% 82% Partial Reagoning Assessment Area



(a) Performance by Area

(b) Eval Distribution

Figure 7: Qualitative evaluation of CoT rationales and inter-reviewer agreement patterns.

(d) M4 Time Series Captioning

identified relevant ECG features; (2) appropriately connected them to the final answer; (3) incorporated patient context (age, artifacts, ...). Overall, the model gave a correct or partially correct ECG interpretation in 92.9% of cases, spanning ECG recognition, reasoning, and contextualization. OpenTSLM showed strongest performance in clinical context integration (85.1% positive) compared to ECG pattern recognition (65.5% positive) and clinical reasoning (62.5% positive) (Figure 7a). Assessment patterns varied notably across reviewers, with some reviewers consistently more favorable across all evaluation areas (Figure 7b). Reviewer disagreement was most common for clinical reasoning, where moderate disagreements were observed between adjacent assessment categories. Complete disagreements between positive and negative assessments were relatively rare across all areas (Figure 15 in Appendix A.5).

5 DISCUSSION

All OpenTSLM models consistently outperform baselines. Text-only models often fail to follow the answer template and thus perform at or below chance (Section 4.1). Finetuned baselines improve substantially on HAR-CoT (60.44% F1 vs. 0% for Llama-3.2-1B) but only slightly on Sleep-CoT (9.05 vs. 2.14). ECG-QA finetuning was infeasible due to high VRAM demands (80k tokens require >100GB per sample). OpenTSLM-SoftPrompt performs best on shorter sequences (Sleep-CoT, TSQA) but becomes impractical as VRAM requirements grow with sequence length (>180GB in simulations with 10,000-length series). With softprompting, smaller models like Gemma-3 270M and 1B quickly exhaust their context and underperform. In contrast, OpenTSLM-Flamingo sustains stable memory across sequence lengths and series (up to 60GB for Llama-3.2-3B with five 10,000-length series). This allows even tiny models, such as Gemma-270M, to deliver strong results, highlighting the efficiency of cross-attention for treating time series as a native modality.

Practical implications. Our results show that even frontier LLMs like GPT-40 are poorly suited for time-series reasoning and that time series must be treated as a distinct modality. With OpenTSLM, even small models like Gemma3 270M outperform GPT-40 (~200B parameters Abacha et al. (2025)) at a fraction of the compute and cost, enabling efficient on-device or mobile deployment. OpenTSLM-SoftPrompt is preferable for short time series, requiring only a few additional parameters for finetuning, but scales exponentially in memory with sequence length, making it impractical for long or multi-series inputs. In contrast, OpenTSLM-Flamingo maintains nearly constant memory across longer or multiple series, performs better on complex datasets, and should therefore be considered the general-purpose option for TSLMs. Perhaps the greatest advantage of TSLMs is the interface they provide for contextualizing results. In ECG-QA, OpenTSLM correctly identified the relevant ECG features in most cases, with missing context only 7.1% of the time. The model demonstrated particularly strong clinical context integration (85.1% positive assessments), thereby offering clinicians and researchers a transparent window into the model's reasoning. As trust is important in medicine, this transparency underscores the value of applying LLMs to time series.

Comparison with prior work. Our approach differs from prior work in several ways. First, we introduce time series as a new modality for LLMs, unlike Sivarajkumar & Wang (2023) and Kim et al. (2024), which tokenize time series. Second, we frame tasks as joint text–time-series reasoning, training models to generate rationales that integrate temporal information. This contrasts with MedualTime Ye et al. (2025) and Time2Lang Pillai et al. (2025), which reprogrammed LLMs with fixed classification or forecasting heads, removing language generation capabilities. Notably, OpenTSLM achieves 40.25 F1 on ECG-QA-CoT, producing rationales across 3,138 questions and 42 templates with diverse answer options. By comparison, Ye et al. report 76 F1 on PTB-XL (underlying dataset of ECG-QA) with only four classes and a fixed classification head Ye et al. (2025). Third, unlike SensorLM Zhang et al. (2025), which is trained from scratch, our models build on pretrained open-weight LLMs, retaining pretrained knowledge. Fourth, while prior work used soft prompting Chow et al. (2024) to model time series implicitly, we find it scales poorly, whereas our OpenTSLM-Flamingo models them explicitly via cross-attention, scaling to long sequences.

Limitations. We acknowledge several limitations. First, our method of encoding time series may not be optimal, as we rely on including mean and standard deviation in accompanying texts to preserve temporal scale. Second, we generated CoT datasets using GPT-40 on plots, which we have shown to perform poorly on these plots alone. Curated datasets likely lead to better rationales. Third, framing tasks as natural language generation does not ensure that the model prioritizes the correct label, underscoring the need for loss functions that explicitly enforce correct answers. Fourth, we did not conduct ablation studies; for example, although OpenTSLM-Flamingo introduces gated cross-attention layers between every two transformer blocks, comparable performance might be achievable with fewer. Finally, while we report strong results on individual datasets, we have not yet demonstrated generalization to unseen data, an essential step toward general TSLMs.

6 CONCLUSION

Our results show that both OpenTSLM variants enable small-scale LLMs to outperform much larger text-only models on time-series tasks, demonstrating that lightweight, domain-adapted architectures can achieve strong performance without massive model scales. With OpenTSLM, we extend openweight pretrained LLMs to process time series retaining knowledge while adapting them to temporal domains. This work may lay the foundation for general-purpose TSLMs capable of handling diverse time-series datasets. Although our focus is healthcare, the ability to reason over longitudinal data has broad relevance in domains such as finance, supply chain management, and industrial monitoring.

REPRODUCIBILITY STATEMENT

All source code associated with this work is publicly available. All external datasets used are open source, and any datasets generated by us have also been released as open source. We additionally release all trained model weights. We also provide the notebooks annotated by clinical doctors for rationale generation on the ECG-QA dataset. These resources ensure full reproducibility of our results.

USE OF LARGE LANGUAGE MODELS

Large Language Models (LLMs) were partially used for text editing, in limited instances, to improve the grammar and clarity of the original text. LLMs were additionally used for reviewing parts of the source code to identify critical errors or bugs. No LLMs were used for data analysis, experimental design, or drawing scientific conclusions.

REFERENCES

- Asma Ben Abacha, Wen wai Yim, Yujuan Fu, Zhaoyi Sun, Meliha Yetisgen, Fei Xia, and Thomas Lin. Medec: A benchmark for medical error detection and correction in clinical notes, 2025. URL https://arxiv.org/abs/2412.19260.
- Amy Abernethy, Laura Adams, Meredith Barrett, Christine Bechtel, Patricia Brennan, Atul Butte, Judith Faulkner, Elaine Fontaine, Stephen Friedhoff, John Halamka, Michael Howell, Kevin Johnson, Peter Long, Deven McGraw, Redonda Miller, Peter Lee, Jonathan Perlin, Donald Rucker, Lewis Sandy, and Kristen Valdes. The promise of digital health: Then, now, and the future. *NAM Perspectives*, 6, 06 2022. doi: 10.31478/202206e.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. URL https://arxiv.org/abs/2204.14198.
- Rawan AlSaad, Alaa Abd-alrazaq, Sabri Boughorbel, Arfan Ahmed, Max-Antoine Renault, Rafat Damseh, and Javaid Sheikh. Multimodal large language models in health care: Applications, challenges, and future outlook. *J Med Internet Res*, 26:e59505, Sep 2024. ISSN 1438-8871. doi: 10.2196/59505. URL https://doi.org/10.2196/59505.
- Luca Arrotta, Gabriele Civitarese, Riccardo Presotto, and Claudio Bettini. Domino: A dataset for context-aware human activity recognition using mobile devices. In 2023 24th IEEE International Conference on Mobile Data Management (MDM), pp. 346–351, 2023. doi: 10.1109/MDM58254. 2023.00063.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-source framework for training large autoregressive vision-language models, 2023. URL https://arxiv.org/abs/2308.01390.
- Winnie Chow, Lauren Gardiner, Haraldur T. Hallgrímsson, Maxwell A. Xu, and Shirley You Ren. Towards time series reasoning with llms, 2024. URL https://arxiv.org/abs/2409.11376.
- Committee Members, Alan H. Kadish, Alfred E. Buxton, Harold L. Kennedy, Bradley P. Knight, Jay W. Mason, Claudio D. Schuger, Cynthia M. Tracy, William L. Winters, Alan W. Boone, Michael Elnicki, John W. Hirshfeld, Beverly H. Lorell, George P. Rodgers, Cynthia M. Tracy, and Howard H. Weitz. Acc/aha clinical competence statement on electrocardiography and ambulatory electrocardiography. *Circulation*, 104(25):3169–3178, 2001. doi: 10.1161/circ.104.25.3169. URL https://www.ahajournals.org/doi/abs/10.1161/circ.104.25.3169.
- GemmaTeam, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex

Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. URL https://arxiv.org/abs/2403.08295.

- Alexia Giannoula, Alba Gutiérrez-Sacristán, Àlex Bravo, Ferran Sanz, and Laura I Furlong. Identifying temporal patterns in patient disease trajectories using dynamic time warping: A population-based study. *Scientific Reports*, 8, 03 2018. doi: 10.1038/s41598-018-22578-1.
- Ary Goldberger, Luís Amaral, L. Glass, Shlomo Havlin, J. Hausdorg, Plamen Ivanov, R. Mark, J. Mietus, G. Moody, Chung-Kang Peng, H. Stanley, and Physiotoolkit Physiobank. Components of a new research resource for complex physiologic signals. *PhysioNet*, 101, 01 2000.
- Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew Gordon Wilson. Large language models are zero-shot time series forecasters. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- Leon Götz, Marcel Kollovieh, Stephan Günnemann, and Leo Schwinn. Byte pair encoding for efficient time series forecasting, 05 2025.
- Susan Henly, Jean Wyman, and Mary Findorff. Health and illness over time the trajectory perspective in nursing science. *Nursing research*, 60:S5–14, 05 2011. doi: 10.1097/NNR. 0b013e318216dfd3.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL https://arxiv.org/abs/2106.09685.
- Isabella Jørgensen, Amalie Haue, Davide Placido, Jessica Hjaltelin, and Søren Brunak. Disease trajectories from healthcare data: Methodologies, key results, and future perspectives. *Annual review of biomedical data science*, 7:251–276, 08 2024. doi: 10.1146/annurev-biodatasci-110123-041001.
- Bob Kemp, Aeilko H. Zwinderman, Bert Tuk, Hilbert A. C. Kamphuisen, and Josefien J. L. Oberye. Analysis of a sleep-dependent neuronal feedback loop: the slow-wave microcontinuity of the eeg. *IEEE Transactions on Biomedical Engineering*, 47:1185–1194, 2000. URL https://api.semanticscholar.org/CorpusID:837298.
- Yubin Kim, Xuhai Xu, Daniel McDuff, Cynthia Breazeal, and Hae Won Park. Health-Ilm: Large language models for health prediction via wearable sensor data. In Tom Pollard, Edward Choi, Pankhuri Singhal, Michael Hughes, Elena Sizikova, Bobak Mortazavi, Irene Chen, Fei Wang, Tasmie Sarker, Matthew McDermott, and Marzyeh Ghassemi (eds.), *Proceedings of the fifth Conference on Health, Inference, and Learning*, volume 248 of *Proceedings of Machine Learning Research*, pp. 522–539. PMLR, 27–28 Jun 2024. URL https://proceedings.mlr.press/v248/kim24b.html.
- Heike Leutheuser, Dominik Schuldhaus, and Bjoern Eskofier. Hierarchical, multi-sensor based classification of daily life activities: Comparison with state-of-the-art algorithms using a benchmark dataset. *PloS one*, 8:e75196, 10 2013. doi: 10.1371/journal.pone.0075196.

- Zechen Li, Shohreh Deldari, Linyao Chen, Hao Xue, and Flora D. Salim. Sensorllm: Human-intuitive alignment of multivariate sensor data with llms for activity recognition, 2025. URL https://arxiv.org/abs/2410.10624.
 - Xin Liu, Daniel McDuff, Geza Kovacs, Isaac Galatzer-Levy, Jacob Sunshine, Jiening Zhan, Ming-Zher Poh, Shun Liao, Paolo Di Achille, and Shwetak Patel. Large language models are few-shot health learners, 2023. URL https://arxiv.org/abs/2305.15525.
 - Caroline Marra, Tim Chico, April Alexandrow, Will Dixon, Norman Briffa, Erin Rainaldi, Max Little, Kristin Size, Athanasios Tsanas, Joseph Franklin, Ritu Kapur, Helen Grice, Anwar Gariban, Joy Ellery, Cathie Sudlow, Amy Abernethy, and Andrew Morris. Addressing the challenges of integrating digital health technologies to measure patient-centred outcomes in clinical registries. *The Lancet Digital Health*, 7, 12 2024. doi: 10.1016/S2589-7500(24)00223-1.
 - Mike Merrill, Mingtian Tan, Vinayak Gupta, Thomas Hartvigsen, and Tim Althoff. Language models still struggle to zero-shot reason about time series. pp. 3512–3533, 01 2024. doi: 10.18653/v1/2024.findings-emnlp.201.
 - Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=Jbdc0vTOcol.
 - 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.
 - Amy Olex and Bridget Mcinnes. Review of temporal reasoning in the clinical domain for timeline extraction: Where we are and where we need to be. *Journal of Biomedical Informatics*, 118: 103784, 04 2021. doi: 10.1016/j.jbi.2021.103784.
 - Louis Pangaro. A new vocabulary and other innovations for improving descriptive in-training evaluations. *Academic medicine: journal of the Association of American Medical Colleges*, 74(11): 1203–7, 1999. doi: 10.1097/00001888-199911000-00012.
 - Arvind Pillai, Dimitris Spathis, Subigya Nepal, Amanda C Collins, Daniel M Mackin, Michael V Heinz, Tess Z Griffin, Nicholas C Jacobson, and Andrew Campbell. Time2lang: Bridging timeseries foundation models and large language models for health sensing beyond prompting, 2025. URL https://arxiv.org/abs/2502.07608.
 - Areti Pouliou, Vasileios Papageorgiou, Georgios Petmezas, Diogo Pessoa, Rui Pedro Paiva, N. Maglaveras, and George Tsaklidis. A new approach for sleep stage identification combining hidden markov models and eeg signal processing. *Journal of Medical and Biological Engineering*, 45, 02 2025. doi: 10.1007/s40846-025-00928-5.
 - Attila Reiss and Didier Stricker. Creating and benchmarking a new dataset for physical activity monitoring. 06 2012. doi: 10.1145/2413097.2413148.
 - Yalini Senathirajah, David Kaufman, Kenrick Cato, Elizabeth Borycki, Jaime Fawcett, and Andre Kushniruk. Characterizing and visualizing display and task fragmentation in the electronic health record: methodological approaches (preprint). *JMIR Human Factors*, 7, 05 2020. doi: 10.2196/18484.
 - Muhammad Shoaib, Hans Scholten, and Paul Havinga. Towards physical activity recognition using smartphone sensors. pp. 80–87, 12 2013. doi: 10.1109/UIC-ATC.2013.43.
 - Muhammad Shoaib, Stephan Bosch, Ozlem Durmaz Incel, Hans Scholten, and Paul J. M. Havinga. Fusion of smartphone motion sensors for physical activity recognition. *Sensors*, 14(6):10146–10176, 2014. ISSN 1424-8220. doi: 10.3390/s140610146. URL https://www.mdpi.com/1424-8220/14/6/10146.

- Muhammad Shoaib, Stephan Bosch, Ozlem Durmaz Incel, Hans Scholten, and Paul J. M. Havinga.
 Complex human activity recognition using smartphone and wrist-worn motion sensors. *Sensors*, 16(4), 2016. ISSN 1424-8220. doi: 10.3390/s16040426. URL https://www.mdpi.com/1424-8220/16/4/426.
 - Sonish Sivarajkumar and Yanshan Wang. Healthprompt: A zero-shot learning paradigm for clinical natural language processing. *AMIA* ... *Annual Symposium proceedings*. *AMIA Symposium*, 2022: 972–981, 04 2023.
 - Allan Stisen, Henrik Blunck, Sourav Bhattacharya, Thor Siiger Prentow, Mikkel Baun Kjærgaard, Anind Dey, Tobias Sonne, and Mads Møller Jensen. Smart devices are different: Assessing and mitigatingmobile sensing heterogeneities for activity recognition. In *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*, SenSys '15, pp. 127–140, New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450336314. doi: 10.1145/2809695.2809718. URL https://doi.org/10.1145/2809695.2809718.
 - Timo Sztyler and Heiner Stuckenschmidt. On-body localization of wearable devices: An investigation of position-aware activity recognition. In 2016 IEEE International Conference on Pervasive Computing and Communications (PerCom), pp. 1–9, 2016. doi: 10.1109/PERCOM.2016. 7456521.
 - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL https://arxiv.org/abs/2302.13971.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. 06 2017. doi: 10.48550/arXiv.1706.03762.
 - Patrick Wagner, Nils Strodthoff, Ralf-Dieter Bousseljot, Dieter Kreiseler, Fatima Lunze, Wojciech Samek, and Tobias Schaeffter. Ptb-xl, a large publicly available electrocardiography dataset. *Scientific Data*, 7:154, 05 2020. doi: 10.1038/s41597-020-0495-6.
 - Chengsen Wang, Qi Qi, Jingyu Wang, Haifeng Sun, Zirui Zhuang, Jinming Wu, Lei Zhang, and Jianxin Liao. Chattime: A unified multimodal time series foundation model bridging numerical and textual data, 2024. URL https://arxiv.org/abs/2412.11376.
 - Jiayang Wu, Wensheng Gan, Zefeng Chen, Wan Shicheng, and Philip Yu. Multimodal large language models: A survey. 11 2023. doi: 10.1109/BigData59044.2023.10386743.
 - Jiexia Ye, Weiqi Zhang, Ziyue Li, Jia Li, Meng Zhao, and Fugee Tsung. Medualtime: A dual-adapter language model for medical time series-text multimodal learning, 2025. URL https://arxiv.org/abs/2406.06620.
 - Andy Wai Kan Yeung, Ali Torkamani, Atul Butte, Benjamin Glicksberg, Björn Schuller, Blanca Rodriguez, Daniel Ting, David Bates, Eva Schaden, Hanchuan Peng, Harald Willschke, Jeroen van der Laak, Josip Car, Kazem Rahimi, Leo Celi, Maciej Banach, M. Kletecka-Pulker, Oliver Kimberger, Roland Eils, and Atanas Atanasov. The promise of digital healthcare technologies. *Frontiers in Public Health*, 11, 09 2023. doi: 10.3389/fpubh.2023.1196596.
 - Xiyuan Zhang, Ranak Roy Chowdhury, Rajesh K. Gupta, and Jingbo Shang. Large language models for time series: a survey. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, IJCAI '24, 2024. ISBN 978-1-956792-04-1. doi: 10.24963/ijcai.2024/921. URL https://doi.org/10.24963/ijcai.2024/921.
 - Yuwei Zhang, Kumar Ayush, Siyuan Qiao, A. Ali Heydari, Girish Narayanswamy, Maxwell A. Xu, Ahmed A. Metwally, Shawn Xu, Jake Garrison, Xuhai Xu, Tim Althoff, Yun Liu, Pushmeet Kohli, Jiening Zhan, Mark Malhotra, Shwetak Patel, Cecilia Mascolo, Xin Liu, Daniel McDuff, and Yuzhe Yang. Sensorlm: Learning the language of wearable sensors, 2025. URL https://arxiv.org/abs/2506.09108.
 - Li Zhou, Simon Parsons, and George Hripcsak. The evaluation of a temporal reasoning system in processing clinical discharge summaries. *Journal of the American Medical Informatics Association: JAMIA*, 15:99–106, 01 2008. doi: 10.1197/jamia.M2467.

APPENDIX Α

702

703

704 705

706

714

715 716

717 718

719

720

721

722

723

724

725

726

727

728

729

730 731

732

733

734

735

736 737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

A.1 TRAINING DETAILS

Table Table 3 provides an overview of the datasets used during training. All data was split into ratios

	Dataset	#Samples (Train/Val/Test)	Num series	Length	Frequency
Stage 1	TSQA*1 M4-Captions	38,400 / 4,800 / 4,800 80,000 / 10,000 / 10,000	1 1	Hours to Years 64-512 points	Not specified Not specified
2	HAR-CoT	68,542 / 8,718 / 8,222	3	2.56s	50Hz
Stage	Sleep-CoT	7,434 / 930 / 930	1	30s	100Hz
St	ECG-QA-CoT	159,313 / 31,137 / 6,019	12	10s	100Hz

Table 3: *1TSQA Wang et al. (2024) Overview of datasets used in Stage 1 (pretraining tasks) and Stage 2 (task-specific CoT reasoning). Datasets are split in 80/10/10 ration.

of 80/10/10 for train/val/test sets.

A.1.1 TRAINING CONFIGURATION

The models were trained with the following configuration:

• Optimizer: AdamW • Learning Rates:

OpenTSLM-SP:

* Time series encoder: 2×10^{-4}

* LoRA: 2×10^{-4}

* Projector: 1×10^{-4} - OpenTSLM-Flamingo:

* Encoder: 2×10^{-4}

* Cross-attention layers: 2×10^{-4}

• Scheduler: Linear learning rate schedule with warmup

• Warmup: 10% of total training steps

• Gradient Clipping: ℓ_2 -norm capped at 1.0

• Weight Decay: 0.01

• Training Length: Up to 200 epochs with early stopping (patience = 5 epochs)

Learning rate choices were informed by Chow et al. (2024).

A.2 GENERATION OF MULTIVARIATE TIME SERIES COT DATASETS

This section provides detailed descriptions of the CoT datasets generated for our study: Human Activity Recognition (HAR-CoT), Sleep Stage Classification (SleepEDF-CoT), and Electrocardiogram Question Answering (ECG-QA-CoT).

Our objective was to enable TSLMs not only to classify time series but also to generate explicit reasoning chains. Since few datasets include CoT text, we generated our own multivariate time series CoT datasets using widely adopted benchmarks in HAR, sleep staging, and ECG-QA, following a similiar approach as proposed by Chow et al. (2024).

For each dataset, we generated rationales with GPT-40 by providing a plot of the data along with the correct label, and prompting the model to produce a rationale leading to that label. The exact prompts are described in Sections A.2.1, A.2.2, and A.2.3. We carefully engineered the prompts and manually reviewed a subset of samples to ensure the generated rationales were consistent and sensible. When plotting, original data was used without normalization. If multiple time series were present in a sample (e.g., three in HAR or twelve in ECG), all were plotted as separate subplots but combined into a single figure.

• **GPT-4o snapshot:** gpt-4o-2024-08-06

• Temperature: 0.3

• Seed: 42

The following subsections describe dataset-specific methodologies, data processing, prompts, answer selection, and final class distributions.

A.2.1 HUMAN ACTIVITY RECOGNITON (HAR) COT

We merged multiple HAR datasets spanning DaLiAc Leutheuser et al. (2013), DOMINO Arrotta et al. (2023), HHAR Stisen et al. (2015), PAMAP2 Reiss & Stricker (2012), RealWorld Sztyler & Stuckenschmidt (2016), and datastes from Shoaib et al. (2013; 2014; 2016). We retain only those activity classes present in all datasets. The final dataset includes eight activity classes: sitting, walking, standing, running, walking up stairs, walking down stairs, lying, and biking. Data is 2 second window of 3 axis acceleration data with 12 class labels.

Data Processing The dataset was processed to create 2.56-second windows of triaxial accelerometer data (X, Y, Z axes). Each sample was visualized as a multi-panel plot showing the acceleration signals across all three axes over the time window.

Prompt for CoT generation We generated CoT rationales by prompting the model with a correct and dissimilar label. The following prompt template was used for HAR-CoT generation:

```
You are shown a time-series plot of accelerometer over a 2.56 second window. This data corresponds to one of two possible activities:
```

This data corresponds to one of two possible activities: [CORRECT_ACTIVITY]

772 [CORRECT_ACTIVITY]
773 [DISSIMILAR_ACTIVITY]

Your task is to classify the activity based on analysis of the data.

Instructions:

- Begin by analyzing the time series without assuming a specific label.
- Think step-by-step about what the observed patterns suggest regarding movement intensity and behavior.
- Write your rationale as a single, natural paragraph, do not use bullet points, numbered steps, or section headings.
- Do not refer back to the plot or to the act of visual analysis in your rationale; the plot is only for reference but you should reason about the time-series data.
- Do **not** assume any answer at the beginning, analyze as if you do not yet know which class is correct.
- Do $\star\star$ not $\star\star$ mention either class label until the final sentence.
- Make sure that your last word is the answer. You MUST end your response with "Answer: [CORRECT_ACTIVITY]":

Answer Selection Strategy For each sample, we implemented a dissimilarity-based answer selection strategy. Given a correct activity label, we selected the most dissimilar activity from a predefined mapping:

- Sitting: walking, running, biking, walking up, walking down
- Walking: sitting, lying, standing, biking, running
- Standing: walking, running, biking, walking up, walking down
- Running: sitting, lying, standing, biking, walking
- Walking up: sitting, lying, standing, biking, running
- Walking down: sitting, lying, standing, biking, running
- Lying: walking, running, biking, walking up, walking down
- Biking: sitting, lying, standing, walking, running

This strategy ensured that the binary classification tasks were challenging and required genuine analysis of movement patterns rather than simple pattern recognition.

Label distribution

A.2.2 SLEEP STAGE CLASSIFICATION CHAIN-OF-THOUGHT (SLEEPEDF-COT)

The SleepEDF-CoT dataset was generated from the Sleep-EDF database, which contains polysomnography recordings with expert-annotated sleep stage labels. The dataset includes five sleep stages: Wake (W), Non-REM stage 1 (N1), Non-REM stage 2 (N2), Non-REM stage 3 (N3), and REM sleep (REM).

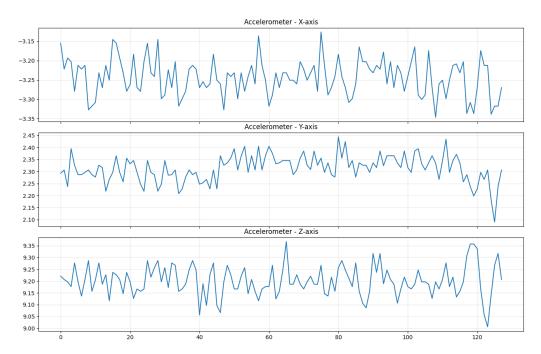


Figure 8: Sample HAR signal input to GPT-40 for rationale generation

Table 4: Per-class sample distribution for HAR-CoT train, validation, and test sets

Class	Train (n=68542)	Val (n=8718)	Test (n=8222)
Biking	4037 (5.9%)	435 (5.0%)	473 (5.8%)
Lying	4305 (6.3%)	682 (7.8%)	444 (5.4%)
Running	8101 (11.8%)	948 (10.9%)	1057 (12.9%)
Sitting	18997 (27.7%)	2315 (26.6%)	2342 (28.5%)
Standing	11001 (16.1%)	1449 (16.6%)	1264 (15.4%)
Walking	12675 (18.5%)	1611 (18.5%)	1508 (18.3%)
Walking Down	4514 (6.6%)	710 (8.1%)	542 (6.6%)
Walking Up	4912 (7.2%)	568 (6.5%)	592 (7.2%)

Data Processing The dataset was processed to create 30-second windows of EEG data from the Fpz-Cz channel. Each sample was visualized as a single-channel EEG plot showing brain activity patterns characteristic of different sleep stages.

Prompt for CoT generation We generated CoT rationales by prompting the model with a correct and dissimilar label. The following prompt template was used for SleepEDF-CoT generation:

You are presented with a time-series plot showing EEG data collected over a 30-second interval. This signal corresponds to one of two possible sleep stages:

- [SLEEP_STAGE_1]
- [SLEEP_STAGE_2]

Your task is to determine the correct sleep stage based solely on the observed patterns in the time series.

Instructions:

- Analyze the data objectively without presuming a particular label.
- Reason carefully and methodically about what the signal patterns $\operatorname{suggest}$
 - regarding sleep stage.

- Write your reasoning as a single, coherent paragraph. Do not use bullet points, lists, or section headers.
- ${\hspace{0.25cm}\text{-}\hspace{0.25cm}}$ Do not reference the plot, visuals, or the process of viewing the data in your explanation; focus only on the characteristics of the time series
- ${\hspace{0.25cm}\text{-}\hspace{0.25cm}}$ Do not mention or speculate about either class during the rationale, only reveal the correct class at the very end.
- Never state that you are uncertain or unable to classify the data. You must always provide a rationale and a final answer.
- Your final sentence must conclude with: "Answer: [CORRECT_SLEEP_STAGE]"

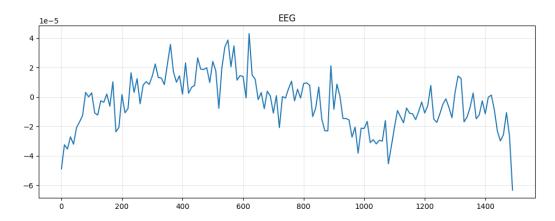


Figure 9: Sample EEG signal input to GPT-40 for sleep stage rationale generation

Answer Selection Strategy For sleep stage classification, we implemented a dissimilarity-based strategy that pairs physiologically distinct sleep stages:

- Wake (W): N3, N4, REM
- N1: W, N3, N4
- N2: W, REM
- N3: W, REM
- N4: W, REM
- REM: N2, N3, N4

This approach ensured that the binary classification tasks required understanding of fundamental differences in brain activity patterns between sleep stages.

Label distribution SleepEDF dataset

Table 5: Per-class sample distribution for train, validation, and test sets (Sleep stages)

Label	Train (n=7434)	Val (n=930)	Test (n=930)
Non-REM 1	410 (5.5%)	52 (5.6%)	51 (5.5%)
Non-REM 2	2057 (27.7%)	257 (27.6%)	257 (27.6%)
Non-REM 3	357 (4.8%)	45 (4.8%)	45 (4.8%)
Non-REM 4	299 (4.0%)	37 (4.0%)	38 (4.1%)
REM	944 (12.7%)	118 (12.7%)	118 (12.7%)
Wake	3367 (45.3%)	421 (45.3%)	421 (45.3%)

A.2.3 ELECTROCARDIOGRAM QUESTION ANSWERING CHAIN-OF-THOUGHT (ECG-QA-COT)

The ECG-QA-CoT dataset was generated from the PTB-XL Wagner et al. (2020) database combined with the ECG-QA Oh et al. (2023) question templates. This dataset contains 12-lead ECG recordings with clinical questions covering various aspects of cardiac analysis, including rhythm analysis, morphology assessment, and diagnostic classification.

Data Processing The dataset was processed to create complete 12-lead ECG recordings (I, II, III, aVR, aVL, aVF, V1, V2, V3, V4, V5, V6) sampled at 100 Hz. Each ECG was visualized as a multi-panel plot showing all 12 leads simultaneously, enabling comprehensive cardiac analysis.

Prompt for CoT generation The following prompt template was used for ECG-QA-CoT generation:

```
You are presented with a complete 12-lead ECG recording showing all standard leads (I, II, III, aVR, aVL, aVF, V1, V2, V3, V4, V5, V6).
```

Clinical Context: [CLINICAL_CONTEXT]

934 Question: [QUESTION]

This question has one of two possible answers:

- [ANSWER_OPTION_1]
- [ANSWER_OPTION_2]

Your task is to analyze the ECG and determine the correct answer based on the observed cardiac patterns. You may include the clinical context in your analysis if it helps you determine the correct answer.

Instructions:

- Analyze the ECG systematically without presuming a particular answer.
- Consider rhythm, rate, morphology, intervals, and any abnormalities you observe across all 12 leads.
- Think step-by-step about what the ECG patterns indicate regarding the clinical question above.
- Write your reasoning as a single, coherent paragraph. Do not use bullet points, lists, or section headers.
- Do not reference the visual aspects of viewing the ECG plot; focus on the cardiac characteristics and clinical significance.
- Do not mention or assume either answer option during your rationale, only reveal the correct answer at the very end.
- NEVER state uncertainty or inability to determine the answer. You MUST always provide clinical reasoning and a definitive answer.
 - Your final sentence must conclude with: "Answer: [CORRECT_ANSWER]"

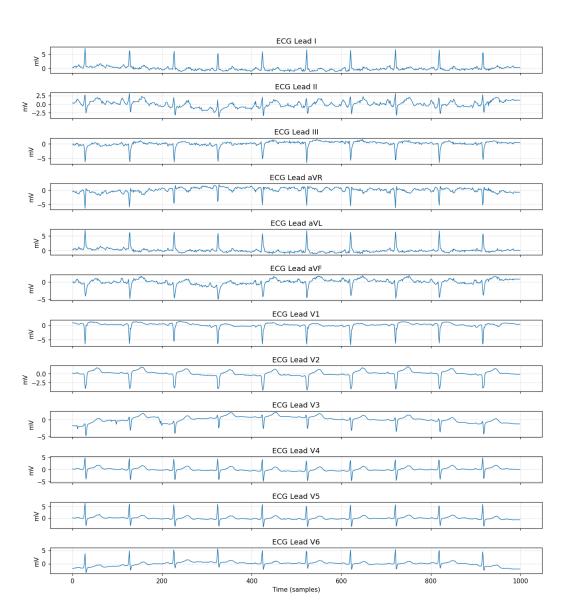


Figure 10: Sample ECG signal input to GPT-40 for rationale generation

Table 6: Per-template sample distribution for ECG-QA CoT train, validation, and test sets

Template ID	Train (n=159,306)	Val (n=31,137)	Test (n=41,093)
Template 1	17,089 (10.7%)	2,924 (9.4%)	3,467 (8.4%)
Template 2	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 3	240 (0.2%)	48 (0.2%)	48 (0.1%)
Template 4	20,861 (13.1%)	3,782 (12.1%)	4,096 (10.0%)
Template 5	20,104 (12.6%)	3,599 (11.6%)	3,905 (9.5%)
Template 6	5,356 (3.4%)	1,022 (3.3%)	1,085 (2.6%)
Template 7	1,137 (0.7%)	221 (0.7%)	224 (0.5%)
Template 8	4,371 (2.7%)	747 (2.4%)	1,466 (3.6%)
Template 9	3,563 (2.2%)	610 (2.0%)	1,200 (2.9%)
Template 10	894 (0.6%)	311 (1.0%)	377 (0.9%)
Template 11	2,861 (1.8%)	533 (1.7%)	964 (2.3%)
Template 12	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 13	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 14	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 15	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 16	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 17	19,952 (12.5%)	3,013 (9.7%)	4,416 (10.7%)
Template 18	9,580 (6.0%)	2,178 (7.0%)	3,806 (9.3%)
Template 19	4,122 (2.6%)	698 (2.2%)	1,395 (3.4%)
Template 20	1,200 (0.8%)	228 (0.7%)	237 (0.6%)
Template 21	180 (0.1%)	36 (0.1%)	36 (0.1%)
Template 22	400 (0.3%)	131 (0.4%)	167 (0.4%)
Template 23	744 (0.5%)	126 (0.4%)	168 (0.4%)
Template 24	90 (0.1%)	18 (0.1%)	18 (0.0%)
Template 25	399 (0.3%)	160 (0.5%)	178 (0.4%)
Template 26	10,585 (6.6%)	1,894 (6.1%)	2,193 (5.3%)
Template 27	1,038 (0.7%)	180 (0.6%)	210 (0.5%)
Template 28	3,600 (2.3%)	720 (2.3%)	720 (1.8%)
Template 29	300 (0.2%)	60 (0.2%)	60 (0.1%)
Template 30	224 (0.1%)	36 (0.1%)	43 (0.1%)
Template 31	1,235 (0.8%)	198 (0.6%)	274 (0.7%)
Template 32	697 (0.4%)	246 (0.8%)	313 (0.8%)
Template 33	6,102 (3.8%)	2,189 (7.0%)	2,775 (6.8%)
Template 34	2,411 (1.5%)	494 (1.6%)	872 (2.1%)
Template 35	246 (0.2%)	18 (0.1%)	50 (0.1%)
Template 36	900 (0.6%)	176 (0.6%)	180 (0.4%)
Template 37	108 (0.1%)	21 (0.1%)	22 (0.1%)
Template 38	523 (0.3%)	192 (0.6%)	241 (0.6%)
Template 39	5,100 (3.2%)	1,019 (3.3%)	1,020 (2.5%)
Template 40	480 (0.3%)	104 (0.3%)	104 (0.3%)
Template 41	1,700 (1.1%)	819 (2.6%)	849 (2.1%)
Template 42	9,114 (5.7%)	2,026 (6.5%)	3,554 (8.6%)

Label distribution

Per-Template Label Distribution Summary

Template ID	Train Labels	Val Labels	Test Labels
Template 1	, ,	no: 1995, yes: 796, not	
	not sure: 978	sure: 133	sure: 261
Template 2	no: 200, yes: 100	no: 40, yes: 20	no: 40, yes: 20

1080	Template 3	st/t change: 60, myocar-	st/t change: 12, myocar-	st/t change: 12, myocar-
1081		dial infarction: 60, none:	dial infarction: 12, none:	dial infarction: 12, none:
1082		60, hypertrophy: 60,	12, hypertrophy: 12,	12, hypertrophy: 12,
1083		conduction disturbance:	conduction disturbance:	conduction disturbance:
1084	Tampleta 4	60 none: 6300, myocardial	12	12
1085	Template 4	infarction in anteroseptal	none: 1258, left ventricular hypertrophy: 110,	none: 1260, myocardial infarction in anterosep-
1086 1087		leads: 618, left anterior	myocardial infarction in	tal leads: 122, myocar-
1088		fascicular block: 593,	anteroseptal leads: 109,	dial infarction in infe-
1089		myocardial infarction in	left anterior fascicular	rior leads: 118, left
1090		inferior leads: 586, first	block: 107, first degree	ventricular hypertrophy:
1091		degree av block: 585	av block: 107	117, left anterior fascic-
1092	Tamplata 5	nama, 6200 myanadial	nana, 1240 laft antarian	ular block: 117
1093	Template 5	none: 6300, myocardial infarction in anteroseptal	none: 1248, left anterior fascicular block: 105,	none: 1260, myocardial infarction in anteroseptal
1094		leads: 578, left anterior	first degree av block:	leads: 117, left anterior
1095		fascicular block: 565,	103, myocardial in-	fascicular block: 116,
1096		first degree av block:	farction in anteroseptal	non-specific intraven-
1097		558, non-specific intra-	leads: 99, left ventricular	tricular conduction
1098		ventricular conduction	hypertrophy: 95	disturbance (block):
1099		disturbance (block): 522		112, first degree av block: 109
1100	Template 6	none: 1530, non-	none: 306, non-specific	none: 306, ventricular
1101	remplate o	diagnostic t abnormal-	st depression: 57, non-	premature complex: 64,
1102		ities: 306, ventricular	diagnostic t abnormali-	non-specific st depres-
1103		premature complex:	ties: 56, ventricular pre-	sion: 63, non-diagnostic
1104		300, non-specific st	mature complex: 55,	t abnormalities: 60,
1105		changes: 295, non-	voltage criteria (qrs) for	atrial premature com-
1106 1107		specific st depression: 294	left ventricular hypertro- phy: 52	plex: 60
1108	Template 7	none: 360, bigeminal	none: 72, sinus rhythm:	none: 72, bigeminal
1109	r	pattern (unknown origin,	19, bigeminal pat-	pattern (unknown ori-
1110		supraventricular, or ven-	tern (unknown origin,	gin, supraventricular, or
1111		tricular): 105, atrial flut-	supraventricular, or	ventricular): 21, sinus
1112		ter: 99, sinus rhythm:	ventricular): 19, atrial	rhythm: 19, atrial fibril-
1113		98, atrial fibrillation: 98	flutter: 18, atrial fibrillation: 17	lation: 18, sinus tachy- cardia: 18
1114	Template 8	myocardial infarction in	myocardial infarction	myocardial infarction in
1115	1	anteroseptal leads: 1050,	in inferior leads: 130,	anteroseptal leads: 304,
1116		myocardial infarction in	left ventricular hyper-	left ventricular hypertro-
1117		inferior leads: 830, left	trophy: 129, myocardial	phy: 282, myocardial in-
1118		ventricular hypertrophy:	infarction in anteroseptal leads: 127, left anterior	farction in inferior leads:
1119		791, left anterior fasci- cular block: 705, non-	fascicular block: 114,	259, left anterior fasci- cular block: 236, non-
1120		specific ischemic: 512	none: 100	specific ischemic: 177
1121	Template 9	myocardial infarction in	left anterior fascicular	left anterior fascicular
1122 1123	-	anteroseptal leads: 635,	block: 111, none: 100,	block: 206, myocardial
1124		left anterior fascicular	non-diagnostic t abnor-	infarction in anterosep-
1125		block: 592, non-specific	malities: 79, myocardial	tal leads: 194, non-
1126		ischemic: 459, left ventricular hypertrophy:	infarction in anteroseptal leads: 74, incom-	specific ischemic: 155, left ventricular hypertro-
1127		432, first degree av	plete right bundle branch	phy: 149, non-specific
1128		block: 399	block: 70	intraventricular conduc-
1129				tion disturbance (block):
1130				127
1131				

1134	Template 10	none: 200, sinus rhythm:	sinus rhythm: 56, none:	none: 100, sinus rhythm:
1135	1	135, atrial fibrillation:	56, atrial fibrillation: 51,	56, sinus tachycardia:
1136		118, sinus tachycardia:	sinus tachycardia: 51, si-	52, atrial fibrillation: 52,
1137		108, sinus bradycardia:	nus arrhythmia: 42	sinus bradycardia: 51
1138	TD 1 . 11	107	100	
1139	Template 11	non-specific st de-	none: 100, non-	non-specific st de-
1140		pression: 692, non-	diagnostic t abnormali-	pression: 194, non-
1141		diagnostic t abnormalities: 570, ventricular	ties: 99, non-specific st depression: 81, ventric-	diagnostic t abnormalities: 182, ventricular
1142		premature complex:	ular premature complex:	premature complex:
1143		414, low amplitude	64, abnormal qrs: 64	142, voltage criteria
1144		t-wave: 334, voltage	on, achierman que, on	(qrs) for left ventricular
1145		criteria (qrs) for left		hypertrophy: 123, q
1146		ventricular hypertrophy:		waves present: 105
1147		329		_
1148	Template 12	no: 200, yes: 100	no: 40, yes: 20	no: 40, yes: 20
1149	Template 13	no: 200, yes: 100	no: 40, yes: 20	no: 40, yes: 20
1150	Template 14	no: 200, yes: 100	no: 40, yes: 20	no: 40, yes: 20
1151	Template 15	no: 200, yes: 100	no: 40, yes: 20	no: 40, yes: 20
1152	Template 16	no: 200, yes: 100	no: 40, yes: 20	no: 40, yes: 20
1153	Template 17 Template 18	no: 14455, yes: 5497 none: 2400, non-specific	no: 2270, yes: 743	no: 3150, yes: 1266
1154	rempiate 18	st depression: 1848,	none: 1150, non-specific st depression: 378, volt-	none: 1200, voltage criteria (qrs) for left
1155		voltage criteria (qrs)	age criteria (qrs) for left	ventricular hypertrophy:
1156		for left ventricular	ventricular hypertrophy:	675, non-specific st
1157		hypertrophy: 1510,	216, q waves present:	depression: 645, non-
1158		non-diagnostic t ab-	114, non-diagnostic t ab-	diagnostic t abnormal-
1159		normalities: 1385, low	normalities: 107	ities: 473, non-specific
1160		amplitude t-wave: 1138		t-wave changes: 308
1161	Template 19	none: 1695, lead I: 1509,	none: 415, lead I: 165,	none: 655, lead I: 438,
1162		lead V6: 1453, lead V5:	lead V6: 154, lead V5:	lead V6: 431, lead V5:
1163	T1-4- 20	1322, lead aVL: 1242	153, lead aVL: 138	399, lead aVL: 392
1164	Template 20	no: 800, yes: 400	no: 160, yes: 68	no: 160, yes: 77
1165	Template 21	none: 60, left axis deviation: 30, right axis	none: 12, left axis deviation: 6, right axis devia-	none: 12, left axis deviation: 6, right axis devia-
1166		deviation: 30, extreme	tion: 6, extreme axis de-	tion: 6, extreme axis de-
1167		axis deviation: 30, nor-	viation: 6, normal heart	viation: 6, normal heart
1168		mal heart axis: 30	axis: 6	axis: 6
1169	Template 22	left axis deviation: 100,	left axis deviation: 50,	left axis deviation: 50,
1170		right axis deviation: 100,	normal heart axis: 50,	right axis deviation: 50,
1171		extreme axis deviation:	right axis deviation: 23,	normal heart axis: 50,
1172		100, normal heart axis:	extreme axis deviation: 8	extreme axis deviation:
1173	T1-4- 22	100	05 21	17
1174	Template 24	no: 545, yes: 199	no: 95, yes: 31	no: 120, yes: 48
1175	Template 24	none: 30, early stage of myocardial infarction:	none: 6, early stage of myocardial infarction: 4,	none: 6, early stage of myocardial infarction: 4,
1176		20, middle stage of my-	middle stage of myocar-	middle stage of myocar-
1177		ocardial infarction: 20,	dial infarction: 4, old	dial infarction: 4, old
1178		old stage of myocardial	stage of myocardial in-	stage of myocardial in-
1179		infarction: 20	farction: 4	farction: 4
1180	Template 25	none of myocardial in-	none of myocardial in-	none of myocardial in-
1181		farction: 100, unknown	farction: 50, unknown	farction: 50, unknown
1182		stage of myocardial in-	stage of myocardial in-	stage of myocardial in-
1183		farction: 100, middle	farction: 50, middle	farction: 50, middle
1184		stage of myocardial in-	stage of myocardial in-	stage of myocardial in-
1185		farction: 100, early stage	farction: 49, early stage	farction: 50, early stage
1186		of myocardial infarction:	of myocardial infarction:	of myocardial infarction:
1187		70, old stage of myocardial infarction: 29	6, old stage of myocar- dial infarction: 5	19, old stage of myocardial infarction: 9
		diai iiiiaiCii0ii, 29	mai imaicuon. J	uiai iiiiaicii0ii. 9

1188				
	Template 26	no: 7335, yes: 3250	no: 1335, yes: 559	no: 1470, yes: 723
1189 1190	Template 27	no: 715, yes: 323	no: 120, yes: 60	no: 145, yes: 65
1191	Template 28	no: 2400, yes: 1200	no: 480, yes: 240	no: 480, yes: 240
1192	Template 29 Template 30	no: 200, yes: 100 none: 60, baseline drift:	no: 40, yes: 20 none: 12, baseline drift:	no: 40, yes: 20
1193	Template 50	58, static noise: 56,	10, static noise: 10, burst	none: 12, static noise: 11, baseline drift: 10,
1194		burst noise: 50, elec-	noise: 10	burst noise: 10, elec-
1195		trodes problems: 44	noise. 10	trodes problems: 7
1196	Template 31	static noise: 448, none:	static noise: 95, none:	static noise: 99, none:
1197	· r	430, baseline drift: 333,	72, burst noise: 47, base-	88, burst noise: 80, base-
1198		burst noise: 309, elec-	line drift: 45	line drift: 71, electrodes
1199		trodes problems: 17		problems: 1
	Template 32	baseline drift: 252, static	none: 100, static noise:	baseline drift: 112, static
1200		noise: 241, none: 200,	83, baseline drift: 78,	noise: 109, none: 100,
1201		burst noise: 174, elec-	burst noise: 22	burst noise: 58, elec-
1202	T1.4 22	trodes problems: 23	1200	trodes problems: 5
1203	Template 33	none: 2400, static noise:	none: 1200, static noise:	none: 1200, static noise:
1204		1824, baseline drift: 1729, burst noise: 823,	675, baseline drift: 358, burst noise: 79	744, baseline drift: 712, burst noise: 283, elec-
1205		electrodes problems: 27	burst noise. 19	trodes problems: 6
1206	Template 34	lead III: 972, lead II:	none: 215, lead III: 182,	lead III: 339, lead II:
1207	rempiace o .	904, lead I: 864, lead	lead II: 175, lead I: 169,	327, lead I: 320, lead
1208		aVR: 844, lead aVL: 779	lead aVR: 165	aVR: 305, lead aVL: 270
1209	Template 35	no: 200, yes: 46	no: 15, yes: 3	no: 40, yes: 10
1210	Template 36	no: 600, yes: 300	no: 120, yes: 56	no: 120, yes: 60
1211	Template 37	supraventricular ex-	supraventricular ex-	supraventricular ex-
1212		trasystoles: 38, ventric-	trasystoles: 7, ex-	trasystoles: 8, extrasys-
1213		ular extrasystoles: 30,	trasystoles: 6, none: 6,	toles: 6, ventricular
1214		none: 30, extrasystoles:	ventricular extrasystoles:	extrasystoles: 6, none: 6
1215	Template 38	28	5	nana: 100 supravantria
1216	Template 36	none: 200, supraventricular extrasystoles: 125,	none: 100, extrasystoles: 55, supraventricular ex-	none: 100, supraventricular extrasystoles: 57,
1217		ventricular extrasystoles:	trasystoles: 27, ventricu-	extrasystoles: 54, ven-
1218 1219		115, extrasystoles: 108	lar extrasystoles: 16	tricular extrasystoles: 38
1219	Template 39	no: 3400, yes: 1700	no: 680, yes: 339	no: 680, yes: 340
1221	Template 40	none: 160, within the	none: 36, within the nor-	none: 36, within the nor-
		normal range: 110,	mal range: 24, above the	mal range: 24, above the
1222		above the normal range:	normal range: 24, below	normal range: 24, below
1223 1224		110, below the normal	the normal range: 20	the normal range: 20
1224	T1.441	range: 100	M.C. d	24.2
	Template 41	within the normal range:	within the normal range:	within the normal range:
1226		600, above the normal range: 600, below the	300, above the normal range: 300, below the	300, above the normal range: 300, below the
1227		normal range: 500	normal range: 219	normal range: 249
1228	Template 42	qt interval: 4393, rr in-	rr interval: 902, qt in-	rr interval: 1730, qt in-
1229		terval: 4336, qt cor-	terval: 880, qt corrected:	terval: 1672, p duration:
1230		rected: 4262, p duration:	879, p duration: 872, qrs	1614, qt corrected: 1592,
1231		4093, qrs duration: 4010	duration: 779	qrs duration: 1486
1232		<u> </u>		-

A.3 Example of Baselines failing to produce meaningful output

As shown in Table 2 in Appendix 4.3, some text-only models achieve 0% F1 score on the CoT datasets. This is because they fail to answer in the " $\langle rationale \rangle$ Answer: $\langle answer \rangle$ " template (see Appendix 4.1). We present some examples of such outputs in the following.

A.3.1 LLAMA3.2-3B BASELINE OUTPUT ON HAR-COT

1239 INPUT PROMPT (TRUNCATED)

1233

1234

1235

1236

1237 1238

1241

You are given accelerometer data in all three dimensions. Your task is to classify the activity based on analysis of the data.

```
1242
1243
      Instructions:
1244
       - Begin by analyzing the time series without assuming a specific label.
1245
       - Think step-by-step about what the observed patterns suggest regarding
      movement intensity and behavior.
1246
       - Write your rationale as a single, natural paragraph, do not use bullet
1247
       points, numbered steps, or section headings.
1248
       - Do **not** mention any class label until the final sentence.
1249
1250
      The following is the accelerometer data on the x-axis, it has mean
       -3.2434 and std 0.0474:\n1 8 6 6 ,4 4 9 ,1 0 5 7 ,8 5 5 , -7 6 2 ,6 5 2
1251
       ,450,652, -1773, -1571, -1369,248, -560,652,
       -1\ 5\ 6\ ,2\ 0\ 6\ 8\ ,1\ 8\ 6\ 6\ ,1\ 0\ 5\ 6\ ,2\ 4\ 8\ ,\ -7\ 6\ 2\ ,\ -3\ 9\ 8\ ,1\ 2\ 5\ 9\ ,\ -5
1253
        6 \ 0 \ , \ -7 \ 6 \ 3 \ , 8 \ 5 \ 5 \ , 1 \ 8 \ 6 \ 5 \ , 2 \ 4 \ 8 \ , 4 \ 6 \ , 2 \ 0 \ 6 \ 8 \ , \ -1 \ 1 \ 6 \ 6 \ , \ -9 \ 6 \ 4 \ , 4 
1254
       1 0 , -5 6 0 ,8 5 5 ...
      The following is the accelerometer data on the y-axis, it has mean 2.3132
1255
       and std 0.0550:\n -3 7 5 , -1 2 4 , -1 3 7 5 ,1 4 8 2 ,2 3 2 , -4 8 1 ,
1256
       -4 8 2 , -3 0 3 , -1 2 4 , -4 8 1 , -6 6 0 ,2 3 2 ,5 3 , -1 7 3 2 , -8 3
1257
       9, -303,946, -303, -1017,767,339,589, -303, -1
1258
      1 9 6 , -1 7 3 2 ,5 8 9 , -3 0 3 , -4 8 1 , -1 7 3 2 , -1 1 9 6 ,5 8 9 , -5 1 7 , -4 8 1 , -1 2 4 , -1 9 1 1 , -1 5 5 3 , -6 6 0 , -1 2 4 , ...
1259
      The following is the accelerometer data on the z-axis, it has mean 9.2017
1260
       and std 0.0639:\n2 9 8 ,7 8 , -7 7 , -3 9 0 ,1 1 7 5 , -7 7 , -1 0 1 8
1261
       ,7 9 ,1 3 3 3 , -7 0 4 ,4 7 ,1 1 7 6 , -2 3 4 ,3 9 2 , -1 3 3 1 ,5 4 9 ,3
1262
       9 2 ,7 9 , -8 6 1 ,5 4 9 , -7 7 , -1 1 7 4 , -5 4 7 , -7 0 4 , -5 4 7 ,1
1263
       3 3 2 ,2 3 6 ,8 6 2 ,1 3 3 2 , -7 7 , ...
1264
1265
       Possible activity labels are:
      biking, lying, running, sitting, standing, walking, walking_down,
1266
       walking_up.
```

- Please now write your rationale. Make sure that your last word is the answer. You MUST end your response with "Answer:"

GENERATED ANSWER

1267 1268

1269

1270

12711272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1285

1286

1287

1288

1289

1290

1291 1292

1295

The accelerometer data on the x-axis has a mean of -3.2434 and a standard deviation of 0.0474. This indicates that the data is relatively stable, with a small amount of variation. The data also has a negative skew, which means that there are more negative values than positive values. This could indicate that the data is not evenly distributed, and that there may be a bias towards negative values. The accelerometer data on the y-axis has a mean of 2.3132 and a standard deviation of 0.0550. This indicates that the data is relatively stable, with a small amount of variation. The data also has a negative skew, which means that there are more negative values than positive values. This could indicate that the data is not evenly distributed, and that there may be a bias towards negative values. Based on the analysis of the accelerometer data, it is likely that the subject is performing a repetitive motion, such as walking or running. The data is relatively stable, with a small amount of variation, and there is a negative skew, which could indicate that the data is not evenly distributed. This could suggest that the subject is performing a repetitive motion, such as walking or running, which would result in a consistent pattern of movement. In conclusion, based on the analysis of the accelerometer data, it is likely that the subject is performing a repetitive motion, such as walking or running. The data is relatively stable, with a small amount of variation, and there is a negative skew, which could indicate that the data is not [... repeated multiple times]

A.3.2 GEMMA3-1B-PT BASELINE OUTPUT ON SLEEP-COT

1293 1294 INPUT PROMPT (TRUNCATED)

You are given a 30-second EEG time series segment. Your task is to classify the sleep stage based on analysis of the data.

```
1296
1297
      Instructions:
1298
       - Analyze the data objectively without presuming a particular label.
1299
       - Reason carefully and methodically about what the signal patterns
1300
       regarding sleep stage.
1301
       - Write your reasoning as a single, coherent paragraph.
1302
       - Only reveal the correct class at the very end.
1303
       - Never state that you are uncertain or unable to classify the data.
1304
        You must always provide a rationale and a final answer.
1305
       The following is the EEG time series, it has mean 0.0000 and std 0.0000:\
1306
       n -3 2 0 , -1 8 0 ,5 5 , -2 1 1 , -6 5 , -4 6 6 , -2 6 9 , -4 4 1 , -7 5 9 , -8 2 9 , -1 0 6 5 , -8 2 9 , -1 0 7 7 , -8 2 9 , -8 9 3 , -7 2 7 , -1
1307
1308
       0 2 6 , -7 7 2 , -9 0 5 , -1 0 4 5 , -6 4 4 , -7 4 0 , -5 3 0 , -8 6 1 ,
       -7 5 3 , -1 4 6 6 , -1 0 2 6 , -9 1 8 , -9 1 8 , -7 5 9 , -3 3 2 , -2 9
1309
       4 \ , \ -2 \ 2 \ 4 \ , \ -4 \ 9 \ 8 \ , \ -1 \ 9 \ 9 \ , \ -2 \ 0 \ 5 \ , \ -3 \ 0 \ 1 \ , \ -6 \ 0 \ 0 \ , \ -5 \ 3 \ 6 \ , \ -6 \ 0 \ 0
1310
       , -4 5 3 , -4 2 2 , -2 3 1 , -7 1 , -2 8 1 , -4 0 9 , -1 4 1 ,1 7 , -7 1 , -2 6 2 ,5 5 , -2 8 8 , -2 3 1 , -4 9 8 ,5 5 , -1 , -1 4 ,5 3 3 ,1 0 0
1311
1312
       , -2 1 1 , -2 1 8 , -2 6 9 ,1 6 4 ,2 9 ,1 9 5 ,1 2 5 , -5 9 , -3 5 8 , -1
1313
       6 7 ...
1314
       Possible sleep stages are:
1315
       Wake, Non-REM stage 1, Non-REM stage 2, Non-REM stage 3, REM sleep,
1316
      Movement
1317
1318
       - Please now write your rationale. Make sure that your last word is the
1319
       answer. You MUST end your response with "Answer:"
1320
1321
       GENERATED ANSWER (TRUNCATED)
1322
       - Please do not use bullet points, lists, or section headers.
1323
       - Please do not use the word "sleep" in your response.
1324
       - Please do not use the word "sleep" in your response.
       - Please do not use the word "sleep" in your response.
       (repeated until out of new tokens)
1326
1327
       A.3.3 GEMMA3-270M BASELINE OUTPUT ON HAR-COT
1328
      INPUT PROMPT (TRUNCATED)
1329
1330
      You are given accelerometer data in all three dimensions. Your task is to
1331
       classify the activity based on analysis of the data.
1332
1333
       Instructions:
1334
       - Begin by analyzing the time series without assuming a specific label.
       - Think step-by-step about what the observed patterns suggest regarding
1335
       movement intensity and behavior.
1336
       - Write your rationale as a single, natural paragraph, do not use bullet
1337
       points, numbered steps, or section headings.
       - Do **not** mention any class label until the final sentence.
1339
       The following is the accelerometer data on the x-axis, it has mean
1340
       -1.9818 and std 1.8034:\n1 2 7 7 ,9 8 5 ,1 2 1 3 ,1 2 5 1 ,1 3 5 1 ,1 8 7
1341
       2 ,1 6 1 2 ,6 9 8 ,4 4 3 ,6 2 9 ,4 3 8 ,6 1 3 ,9 3 2 ,9 2 7 ,1 0 3 2 ,9
1342
       2 1 ,9 3 7 ,6 7 7 ,5 4 4 ,6 5 6 ,5 3 9 ,9 2 7 ,8 9 5 ,9 6 4 ,1 0 7 5 ,1 0
1343
       4 9 ,8 5 2 ,9 3 2 ,1 5 9 6 ,1 9 5 2 ,1 8 8 3 ,1 4 1 0 ,3 7 4 , ...
       The following is the accelerometer data on the y-axis, it has mean 5.8203
       and std 4.7959:\n7 1 3 ,4 4 1 ,4 7 6 , -1 3 0 , -7 5 5 , -1 2 8 7 , -1 0
1345
       7 6 , -7 8 8 , -8 2 2 , -7 0 1 , -3 5 9 , -5 0 1 , -4 2 1 , -6 6 5 , -6
1346
       1347
       -1 0 0 5 , -9 9 0 , -9 0 0 , -8 1 5 , -7 3 4 , -6 1 1 ,...
1348
       The following is the accelerometer data on the z-axis, it has mean 6.3687
       and std 3.9830:\n -2 9 0 , -1 0 4 ,2 1 9 ,8 9 0 ,9 8 7 ,1 1 2 9 ,1 2 2 5
1349
```

,1 2 9 0 ,1 0 3 4 ,9 8 6 ,9 3 3 ,8 3 3 ,6 7 2 ,7 5 5 ,9 2 3 ,1 0 1 4 ,1

```
1350
      0 2 9 ,1 0 2 4 ,1 1 7 7 ,1 1 5 5 ,9 7 8 ,6 3 9 ,6 6 2 ,9 9 4 ,9 7 6 ,8 1
1351
      4,886,1119,1230,1087,669,529,880,986,...
      Possible activity labels are:
1353
      biking, lying, running, sitting, standing, walking, walking_down,
1354
      walking_up.
1355
      - Please now write your rationale. Make sure that your last word is the
1356
      answer. You MUST end your response with "Answer:"
1357
1358
      GENERATED ANSWER (TEXT-ONLY BASELINE)
1359
      The following is the accelerometer data on the x-axis, it has mean 1.0000
1360
      and std 1.0000: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22
1361
      23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46
1362
      47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70
1363
      71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94
      95 96 97 98 99
1364
      100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117
1365
      118 119 [... continues to count]
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
```

A.4 FURTHER EXAMPLES FOR RATIONALES GENERATED BY OPENTSLM MODELS

A.4.1 M4 TIME-SERIES CAPTIONING



Caption

The time-series graph illustrates a progressive upward trend over a period of 75 time intervals. Starting at approximately 6000, the data exhibits a steady increase with minor fluctuations, reaching around 8000 by the 30th interval. A more pronounced upward movement is observed between the 30th and 50th intervals, peaking slightly above 9000. Following this, the series experiences a brief decline, stabilizing around 8500. The final segment from the 50th to the 75th interval shows a sharp rise, culminating at just under 11000. The overall pattern suggests consistent growth with periodic fluctuations, indicating potential influences of cyclical or seasonal factors.

(a) M4 Caption (Gemma3-1B)



The time-series graph illustrates a clear upward trend over the observed period, with data points ranging from 0 to 130 on the x-axis and values from 0 to 9000 on the y-axis. Initially, the series shows a gradual increase, with minor fluctuations around the 1000 mark. Around the midpoint, the data exhibits more pronounced volatility, with several peaks and troughs. Despite these fluctuations, the overall trajectory remains positive. As the series progresses, the upward trend becomes more pronounced, with values exceeding 8000 towards the end. The graph suggests a strong growth pattern, possibly indicating and underlying factor driving the increase in values over time.

(c) M4 Caption (Llama3.2-1B)



The time-series graph illustrates a progressive upward trend over a period, with the data points ranging from approximately 6,000 to 10,000. Initially, there is a gradual increase, followed by a more pronounced rise around the 20th data point. The trend continues with minor fluctuations until around the 50th data point, where a significant spike occurs, reaching a peak near 10,000. This is followed by a slight decline and stabilization. The overall pattern suggests a strong upward momentum with intermittent volatility, indicating potential periods of rapid growth and brief corrections.

(b) M4 Caption (Llama3.2-1B)



The time-series graph displays data points over a period, showcasing fluctuations between values of approximately 1300 and 1900. The series begins with moderate variability, experiencing a sharp increase around the 10th data point, reaching a peak near 1900. Following this peak, the data exhibits a downward trend with intermittent spikes and drops. Notably, there are significant drops around the 50th and 80th data points, where the values dip close to 1300. The latter part the series shows continued volatility with smaller peaks and troughs, indicating persistent fluctuations. The overall pattern suggests a high degree of variability, with no clear long-term trend.

(d) M4 Caption (Llama3.2-1B)

Figure 11: Examples of time series captioning on our generated M4 caption dataset.

A.4.2 HUMAN ACTIVITY RECOGNITION RATIONALES COT (a) Standing (Correct Prediction) (b) Sitting (Correct Prediction) (c) Standing (Correct Prediction) (d) Running (Correct Prediction) (e) Running (Correct Prediction) (f) Walking (Correct Prediction)

Figure 12: Example chain of thought rationales for human activity recognition, generated with OpenTSLMFlamingo/Llama3.2-1B.

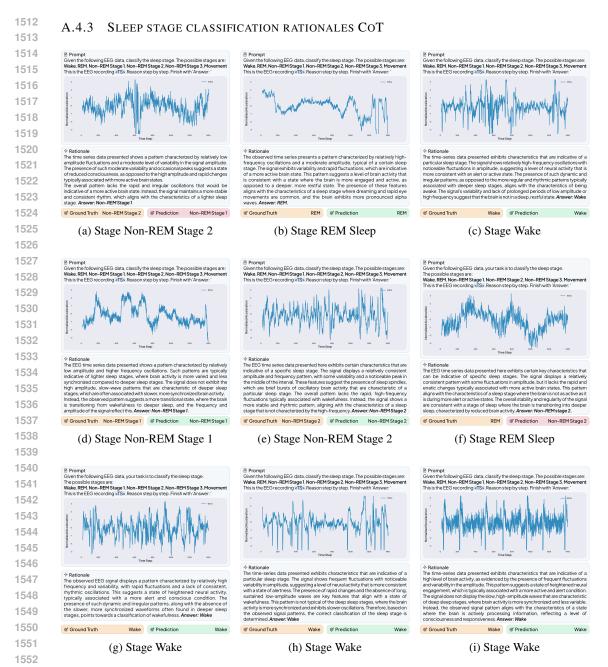


Figure 13: Example rationales for sleep stage classification, generated with OpenTSLM-Flamingo/Llama3.2-1B.

A.5 ECG EVALUATION RUBRIC

These are the questions asked to clinicians during evaluation of ECG-QA rationales generated by OpenTSLMFlamingo/Llama3.2-3B. See Appendix 4.5 for details.

Assessment Criteria	Description	Options
1. ECG Pattern	Did the model correctly identify	Yes; Some but not all; None
Recognition Accuracy	the relevant ECG features needed	identified
	to answer the question?	
2. Clinical Reasoning	Did the model appropriately con-	Yes; Some incorrect logic;
Quality	nect the identified ECG features	Completely incorrect logic
	to the final answer?	
3. Clinical Context	Did the model appropriately in-	Yes; Used some key
Integration	corporate patient clinical back-	background; No did not use
	ground (age, recording condi-	any relevant background
	tions, artifacts) in its interpreta-	
	tion?	

Table 8: Assessment Criteria for ECG Interpretation Reasoning

A.5.1 ECG REVIEW FORM

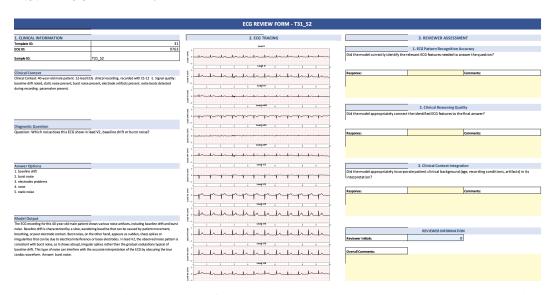


Figure 14: ECG Review Form. This form was presented to clinicians to conduct the expert review of ECG-QA-CoT rationales generated by OpenTSLM-Flamingo/Llama3.2-3B (best model during evaluation, see Table 2).

A.5.2 REVIEWER DISAGREEMENT PATTERNS

?? shows disagreement of reviewers on generated ECG-rationales (see Appendix 4.5).

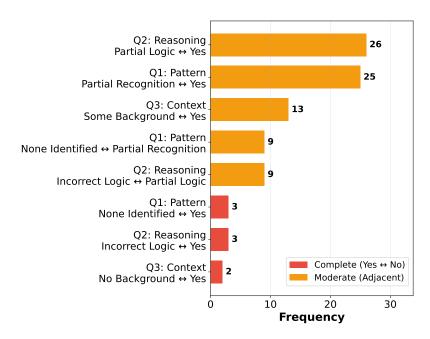


Figure 15: Disagreement Patterns

A.6 EVALUATION OF MEMORY CONSUMPTION

We complement the main results with detailed tables and plots. Figure 16 illustrates scaling trends, while the following subsections report detailed VRAM usage for both CoT datasets and synthetic simulation data.

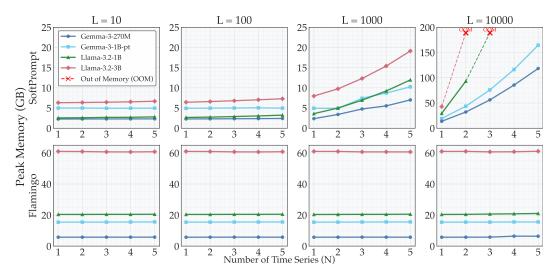


Figure 16: Simulation of memory scaling with total sequence length $(N \times L)$.

A.6.1 MEMORY USAGE ON COT DATASETS

Table 9 reports VRAM for TSQA, HAR-CoT, Sleep-CoT, ECG-QA-CoT datasets. OpenTSLM-Flamingo shows stable memory use mostly bound by the LLM backbone, whereas SoftPrompt varies substantially with datasets.

1689 1690

1692

1693 1694

1695

1696

1726 1727

Table 9: VRAM Usage (GB) for Regular Datasets

Method	Model	TSQA	HAR-CoT	SleepEDF- CoT	ECG-QA- CoT
pt M	Llama-3.2-1B	4.4	9.6	15.9	64.9
SC	Llama-3.2-3B	8.1	14.3	20.3	87.1
nT. Pr	Gemma-3-270M	2.4	8.6	20.1	24.1
OpenTSLM SoftPrompt	Gemma-3-1B-pt	5.1	6.1	14.7	32.7
M o	Llama-3.2-1B	20.5	22.0	21.6	20.9
SL	Llama-3.2-3B	61.1	63.5	63.4	71.6
n. m.	Gemma-3-270M	5.7	6.4	6.3	7.3
OpenTSLM Flamingo	Gemma-3-1B-pt	15.6	16.3	15.7	18.4

A.6.2 MEMORY USAGE FOR SIMULATION DATA

Table 10 shows results for simulated datasets, using permutations of N=[1,2,3,4,5] and L=[10,100,1000,10000]. OpenTSLM-Flamingo requires almost constant memory with varying sequence length L and number of concurrent series N, while OpenTSLM-SoftPrompt grows with both until going out of memory (OOM) for larger time series.

Simulation dataset generation. To generate the simulation dataset, we generate random data with combinations of N=[1,2,3,4,5] and L=[10,100,1000,10000] according to the following pseudocode:

```
1697
      num\_series = n
1698
      series_length = 1
1699
      simulation dataset = []
1700
      for element id in 1..200:
1701
          time_series_texts = []
1702
          time_series_simulations = []
          for i in 1..num_series:
1703
               series_i = random_normal(series_length)
1704
               series_mean = mean(series_i)
1705
               series_std = std(series_i)
               normalized_i = normalize(series_i)
1707
               time_series_simualtions.append(
                   normalized i
1709
               )
1710
               time_series_texts.append(
1711
                   "This is a time series with mean {series mean} "
1712
                   "and std {series_std}."
1713
               )
1714
          simulation_dataset.append([
1715
               {
                   "Series": time_series_simualtions,
1716
                   "Texts": time_series_texts,
1717
                   "PrePrompt": "You are given different time series. "
1718
                                 "All have the same length"
1719
                                 "of {length} data points.",
1720
                   "PostPrompt": "Predict the pattern "
1721
                   "of the time series. Answer:",
1722
                   "Answer": "This is a random pattern."
1723
1724
          1)
1725
```

Table 10: VRAM Usage (GB) for Simulation Datasets

		OpenTSLM-SoftPrompt				OpenTSLM-Flamingo			
		LI	_L aMA	Gemma		LLaMA		Gemma	
L	N	1B	3B	270M	1B	1B	3B	270M	1B
10	1	2.6	6.3	2.3	5.0	20.4	61.0	5.7	15.4
10	2	2.6	6.4	2.3	5.0	20.4	60.9	5.7	15.5
10	3	2.7	6.4	2.3	4.9	20.4	60.7	5.8	15.5
10	4	2.7	6.5	2.3	5.0	20.5	60.7	5.8	15.5
10	5	2.8	6.7	2.3	5.0	20.5	60.8	5.8	15.6
100	1	2.7	6.4	2.3	4.9	20.4	61.0	5.7	15.4
100	2	2.8	6.6	2.3	5.0	20.4	60.9	5.7	15.5
100	3	2.9	6.8	2.3	5.0	20.5	60.7	5.8	15.5
100	4	3.0	7.0	2.4	5.0	20.5	60.7	5.8	15.5
100	5	3.2	7.3	2.4	5.0	20.5	60.8	5.7	15.5
1000	1	3.6	8.0	2.4	5.0	20.4	61.0	5.7	15.4
1000	2	5.0	9.8	3.4	4.9	20.4	61.0	5.7	15.4
1000	3	6.9	12.3	4.8	7.4	20.4	60.7	5.8	15.5
1000	4	9.2	15.4	5.5	8.7	20.5	60.7	5.8	15.6
1000	5	12.0	19.1	7.0	10.2	20.6	60.7	5.7	15.6
10000	1	29.5	42.7	13.7	19.2	20.4	61.0	5.7	15.4
10000	2	93.3	191.4	32.1	43.6	20.4	61.0	5.7	15.4
10000	3	OOM^{*1}	OOM	56.1	76.0	20.6	60.7	5.8	15.5
10000	4	OOM	OOM	85.6	116.4	20.8	60.8	6.4	15.5
10000	5	OOM	OOM	118.4	164.5	21.0	61.1	6.4	15.5

Table 11: *1 OOM: Out of memory; OpenTSLM-SoftPrompt requires more tokens for longer time series, and separate tokens for separate time series. Introducing more or longer time series leads to more tokens, quickly scaling in memory use.