# SensorLLM: Aligning Large Language Models with Motion Sensors for Human Activity Recognition

Anonymous ACL submission

## Abstract

We introduce SensorLLM, a two-stage framework that enables Large Language Models (LLMs) to perform human activity recognition (HAR) from sensor data. Despite their strong reasoning and generalization capabilities, LLMs remain underutilized for motion sensor data due to the lack of semantic context in time-series, computational constraints, and challenges in processing numerical inputs. SensorLLM addresses these limitations through a Sensor-Language Alignment stage, where we introduce special tokens for each sensor channel and automatically generate textual trend descriptions. This alignment enables LLMs to capture numerical variations, channel-specific features, and data of varying durations-without requiring human annotations. In the subsequent Task-Aware Tuning stage, we refine the model for HAR classification, achieving performance that matches or surpasses state-of-the-art methods. Our results demonstrate that SensorLLM evolves into an effective sensor learner, reasoner, and classifier through Sensor-Language Alignment, generalizing across diverse HAR datasets. We believe this work establishes a foundation for future research on time-series and text alignment, paving the way for foundation models in sensor data analysis. Our codes are available at https://anonymous.4open. science/r/sensorllm\_code-E0FC.

#### 1 Introduction

002

006

013

016

017

021

022

024

031

040

043

Human Activity Recognition (HAR) is a timeseries classification task that maps sensor signals, such as accelerometer and gyroscope data, to human activities. Traditional models like LSTM (Guan and Plötz, 2017; Hammerla et al., 2016) and DeepConvLSTM (Ordóñez and Roggen, 2016) learn high-level features but are task-specific and struggle to generalize across different sensor configurations and activity sets. In contrast, Large Language Models (LLMs) (Han et al., 2021) have



Figure 1: SensorLLM can analyze and summarize trends in captured sensor data, facilitating human activity recognition tasks.

shown remarkable success in integrating diverse data types (Wu et al., 2023b; Yin et al., 2023), including text and images.

044

045

046

047

051

053

058

059

060

061

062

063

064

065

066

067

Enabling LLMs to process sensor data (Jin et al., 2023) requires either (1) pretraining or fine-tuning on time-series data (Zhou et al., 2023a), which demands substantial computational resources and is hindered by limited and imbalanced labeled data, or (2) leveraging zero-shot and few-shot prompting by converting sensor data into text (Kim et al., 2024; Ji et al., 2024). The latter approach avoids retraining but introduces key challenges: (i) Numerical encoding issues-LLM tokenizers, designed for text, struggle with numerical values, treating consecutive numbers as independent tokens (Nate Gruver and Wilson, 2023) and failing to preserve temporal dependencies (Spathis and Kawsar, 2024). (ii) Sequence length constraints-sensor data often exceeds LLMs' maximum context length, leading to truncation, information loss, and increased computational costs. (iii) Multi-channel complexity-LLMs process univariate inputs, making it difficult to encode multi-sensor data in a way that retains inter-channel dependencies. (iv) Prompt engineering challenges—designing effective prompts that enable LLMs to interpret numerical sensor readings, detect trends and classify activities remains a challenge (Liu et al., 2023b).

069

070

087

090

091

093

097

100

101

102

103

104

105

106

107

109

110

111

112

113

114

115

116

117

To address these challenges, we propose SensorLLM, a framework that aligns sensor data with natural language, allowing LLMs to analyze sensor data through text-based interactions (see Figure 1) without modifying the LLM itself. A major challenge is the annotation bottleneck, further compounded by the complexity and heterogeneity of sensor data. Unlike image-text pairs, sensor data comprises multi-channel numerical signals with diverse characteristics, making direct interpretation and alignment difficult.

Existing methods (Jin et al., 2024a; Sun et al., 2024a) have explored condensed text prototypes for alignment, but these approaches often lack interpretability and require extensive tuning to select suitable prototypes. In contrast, we propose an automatic text generation approach that aligns with human intuition by deriving descriptive trend-based text directly from time-series data using statistical analyses and predefined templates. This method is precise, scalable, and interpretable, eliminating the need for manual annotations while preserving essential sensor characteristics. SensorLLM follows a two-stage framework:

Sensor-Language Alignment Stage. We automatically generate question-answer pairs to align sensor data with text while preserving temporal features using a pretrained encoder. The resulting embeddings are mapped into a space interpretable by the LLM, mitigating issues associated with textspecific tokenization. Additionally, we introduce *special tokens* for sensor channels, enabling LLMs to effectively capture multi-channel dependencies.

**Task-Aware Tuning Stage.** The aligned embeddings are utilized for HAR, leveraging the LLM's reasoning capabilities while keeping its parameters frozen. This design extends LLMs beyond their original training, addressing concerns raised by Tan et al. (2024) regarding their applicability to time-series data. To our knowledge, this is the first approach to integrate sensor data into LLMs for sensor data analysis and HAR tasks.

To our knowledge, this is the first approach to integrate sensor data into LLMs for sensor data analysis and HAR tasks. The key contributions of this work are as follows: • We propose a fully automated, humanintuitive approach for aligning time-series data with descriptive text, eliminating the need for manual annotations. Using text similarity metrics, human evaluations, and LLMbased assessments, we demonstrate that SensorLLM effectively captures temporal patterns and channel-specific features, enabling robust multimodal understanding. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

- SensorLLM achieves competitive results across five HAR datasets, matching or surpassing state-of-the-art models. Experiments further validate that *modality alignment* and *task-specific prompts* significantly enhance the LLM's ability to interpret and classify sensor data.
- We show that SensorLLM maintains strong performance in the Task-Aware Tuning Stage, even when applied to datasets distinct from those used during alignment, highlighting its robustness and generalizability in HAR tasks.

# 2 Related Work

In this section, we discuss recent developments in leveraging LLMs for time-series data, specifically focusing on two categories: (1) LLMs for time series as text and (2) Multimodal Large Language Models (MLLMs) for sensor data. A broader overview of other related works, including deep learning approaches to HAR and additional LLMbased forecasting methods, is provided in Appendix A.1.

LLMs for Time Series as Text. While LLMs excel in processing natural language, applying them directly to time-series data poses unique challenges (Spathis and Kawsar, 2024). Certain methods address this by treating time-series signals as raw text, using the same tokenization as natural language. Notable examples include PromptCast (Xue and Salim, 2023), which transforms numeric inputs into textual prompts for zero-shot forecasting, and LLMTime (Gruver et al., 2024), which encodes time-series as numerical strings for GPT-like models. However, due to the lack of specialized tokenizers for numeric sequences, LLMs may fail to capture crucial temporal dependencies and repetitive patterns (Spathis and Kawsar, 2024). To mitigate these issues, several works employ time-series encoders before mapping the resulting embeddings to language model spaces (Liu et al., 2024a; Zhou

256

257

258

259

260

261

262

263

216

217

218

et al., 2023c; Xia et al., 2024), thus aligning sensor embeddings with textual embeddings in a contrastive or supervised manner.

MLLMs for Sensor Data. Extending LLMs to non-textual domains has gained traction, particularly through MLLMs that accept inputs beyond text, such as images or speech. For sensor data, the challenge lies in representing continuous signals effectively. Yoon et al. (2024) propose to ground MLLMs with sensor data via visual prompting. Sensor signals are first visualized as images, guiding the MLLM to analyze the visualized sensor traces alongside task descriptions, which also lower token costs compared to raw-text baselines. Similarly, Moon et al. (2023) introduce IMU2CLIP, which aligns inertial measurement unit streams with text and video in a joint representation space. This approach enables wearable AI applications like motion-based media search and LM-based multimodal reasoning, showcasing how sensor data can be integrated into broader multimodal frameworks.

#### 3 Methods

170

171

172

173

174

175

177

178

179

181

182

183

184

186

188

190

191

192

193 194

195

196

199

207

208

210

211

In this work, we propose SensorLLM, a framework that aligns sensor data with descriptive text through automatically generated question-answer pairs. Our goal is to develop a multimodal model with reasoning capabilities for analyzing wearable sensor data. As shown in Figure 2, SensorLLM comprises three core components: (1) a pretrained LLM, (2) a pretrained time-series (TS) embedder, and (3) an alignment module MLP, all operating in a two-stage framework.

In the Sensor-Language Alignment Stage, a generative model aligns sensor readings with text based on user instructions (Liu et al., 2023a) and questions. In the Task-Aware Tuning Stage, a classifier is added on top of the LLM to perform HAR. Crucially, only the MLP in both stages and the classifier in the Task-Aware Tuning Stage are trainable, while the backbone LLM and TS embedder remain frozen. This design yields a highly efficient and lightweight training process, requiring only 5.67% (535.9M) of the parameters to be trainable in the alignment stage and 0.12% (10.5M) in the tuning stage.

#### 212 **3.1** Sensor-Text Data Generation

213Aligning time-series data with text for LLM-based214tasks is challenging due to the lack of rich semantic215labels beyond class annotations, making manual

annotation impractical (Deldari et al., 2024; Haresamudram et al., 2024). While prior works rely on predefined text prototypes (Sun et al., 2024b; Jin et al., 2024a), we aim for a more human-intuitive representation of sensor data.

We argue that time-series data inherently contains semantic patterns that can be expressed through descriptive text, from simple numerical trends to statistical insights. To achieve this, we automatically generate descriptive text by analyzing observed trends and fluctuations in the data. Using predefined templates, we construct diverse question-answer (QA) pairs that capture trend changes while ensuring accuracy and scalability. These templates (Appendix A.2) are randomly combined to enhance diversity. For example:

- (1) The time-series data represents readings taken from a <S> sensor between  $<t_s>$  and  $<t_e>$  seconds.
- (2) To sum up, the data exhibited a <T> trend for a cumulative period of <t<sub>t</sub>> seconds.

where T and S denote specific trends and sensor types, and t corresponds to numerical values.

#### 3.2 Sensor-Language Alignment

As shown in Figure 2 (a), the Sensor-Language Alignment stage employs a generative model to create multimodal sentences that combine singlechannel sensor readings with textual descriptions. The sensor data is represented as a matrix  $\mathbf{X} \in$  $\mathbb{R}^{C \times T}$ , where C is the number of sensor channels and T is the sequence length. Each channel's data, denoted as  $\mathbf{X}^c$  for channel c, is processed independently to preserve its unique characteristics. The data is then divided into non-overlapping segments,  $\mathbf{X}_{S}^{c}$ , where S is the total number of segments. Each segment  $x_s$  is assigned a random length l within a predefined range, allowing the model to learn from varying temporal patterns and trend variations. This segmentation strategy ensures that both long-term trends and short-term fluctuations are effectively captured in the generated multimodal sentences.

We use Chronos (Ansari et al., 2024) as the TS embedder, generating segment embeddings  $\hat{x}_s \in \mathbb{R}^{(l+1) \times d_{ts}}$ , where  $d_{ts}$  is the feature dimension, and (l+1) accounts for the [EOS] token appended during Chronos tokenization. Prior to inputting sensor segments into Chronos (Appendix A.3), we apply



Figure 2: Our proposed SensorLLM framework: (a) Sensor-Language Alignment Stage, where a generative model aligns sensor readings with automatically generated text; (b) Task-Aware Tuning Stage, where a classification model leverages the aligned modalities to perform HAR.

instance normalization  $\tilde{x}_s = \frac{x_s - \text{mean}(x_s)}{\text{std}(x_s)}$  to standardize the data. Llama3-8B (Touvron et al., 2023) serves as our LLM backbone.

265

270

271

274

275

276

277

278

281

282

291

Alignment Module. To transform TS embeddings  $\hat{x}_s$  into text-aligned embeddings  $\hat{a}_s \in$  $\mathbb{R}^{(l+1) \times D}$  for downstream tasks, we introduce an alignment projection module. This module, implemented as a multi-layer perceptron (MLP), first maps sensor embeddings to an intermediate space of dimension  $d_m$  and then projects them to the target dimension D. Formally,

$$\hat{a}_s = \mathbf{W}_2 \cdot \sigma(\mathbf{W}_1 \hat{x}_s + \mathbf{b}_1) + \mathbf{b}_2, \qquad (1)$$

where  $\mathbf{W}_1 \in \mathbb{R}^{d_m imes d_{ts}}$  and  $\mathbf{W}_2 \in \mathbb{R}^{D imes d_m}$  are learnable weights,  $\mathbf{b}_1$  and  $\mathbf{b}_2$  are biases, and  $\sigma$  is the GELU activation function (Hendrycks and Gimpel, 2016). This projection ensures that the transformed embeddings  $\hat{a}_s$  are semantically aligned with the text embedding space, making them suitable for tasks such as text generation and classification.

Input Embedding. To integrate sensor data into the LLM, we introduce two special tokens per sensor channel (e.g., <x\_acc\_start> and 286 <x\_acc\_end> for the x-axis accelerometer), extending the LLM's embedding matrix from  $\mathbf{E} \in \mathbb{R}^{V \times D}$ to  $\mathbf{E} \in \mathbb{R}^{V' \times D}$ , where V' = V + 2c, with V as the vocabulary size and c as the number of channels. 290 These special token embeddings are concatenated with the aligned sensor embeddings. The final combined sensor representation  $\hat{o}_s \in \mathbb{R}^{(l+3) \times D}$  is then concatenated with instruction and question embed-294

dings to form the full input sequence  $\hat{z} \in \mathbb{R}^{k \times D}$ , where k is the total number of tokens.

295

297

298

299

300

301

302

303

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

324

Loss Function. SensorLLM processes an input sequence  $\mathbf{Z}_s = \{z_s^i\}_{i=1}^K$  consisting of sensor and text embeddings and generates an output sequence  $\mathbf{Z}_t = \{z_t^i\}_{i=1}^N$ , where  $z_s^i, z_t^i \in V'$ , and K and N represent the number of input and output tokens, respectively. The model is trained using a causal language modeling objective, predicting the next token based on previous ones. The optimization minimizes the negative log-likelihood:

$$\mathcal{L}_{gen} = -\sum_{i=0}^{N-1} \log P(z_t^i | Z_t^{< i}, z_s).$$
(2)

Loss is computed only on generated tokens, ensuring SensorLLM effectively integrates sensor and text embeddings to produce coherent, contextually appropriate responses.

#### 3.3 Task-Aware Tuning

As shown in Figure 2 (b), the Task-Aware Tuning stage refines the multimodal sensor-text embeddings for HAR. This stage integrates multi-channel sensor readings with activity labels, aligning temporal patterns with human activities. The input sensor data X is segmented into overlapping windows of size L with a 50% overlap (Li et al., 2018), forming segments  $\mathbf{X}_S \in \mathbb{R}^{S \times C \times L}$ , where S is the number of segments and C is the number of channels. The pretrained alignment module from the first stage maps sensor data to activity labels, preserving inter-channel dependencies while learning activity-related patterns.

**Input Embedding.** For each sensor channel c, we retrieve its aligned sensor embeddings  $\hat{o}_s^c$ . These embeddings are then concatenated across all channels, along with their corresponding statistical features (mean and variance), to form the final input embedding:

325

330

332

336

337

338

340

341

342

343

354

357

361

369

$$\hat{z} = \hat{o}_s^1 \oplus \hat{o}_s^2 \oplus \dots \oplus \hat{o}_s^C \oplus \hat{z}_{\text{stat}}, \qquad (3)$$

where  $\hat{z}_{\text{stat}}$  represents the statistical information, and *C* is the number of channels. This ensures the model integrates both temporal and statistical characteristics for HAR.

**Loss Function.** The input token sequence is processed by the LLM, yielding a latent representation  $\mathbf{H} \in \mathbb{R}^{K \times D}$ , where K is the number of tokens and D is the embedding dimension. Due to causal masking, we extract the final hidden state,  $\mathbf{h} = \mathbf{H}_K$ , which encodes all preceding token information. This pooled vector is passed through a fully connected layer to produce a prediction vector of size M, where M is the number of activity classes. The final class probabilities  $\hat{y}_i$  are obtained via the softmax function, and the model is optimized using cross-entropy loss:

$$\mathcal{L}_{cls} = -\sum_{i=0}^{M-1} y_i \log \hat{y}_i, \qquad (4)$$

where  $y_i$  is the ground truth label.

## 4 Experiments

In this section, we evaluate SensorLLM in enabling LLMs to interpret, reason about, and classify sensor data for HAR tasks. All experiments are conducted on NVIDIA A100-80G GPUs. To assess the LLM's ability to learn and generalize from raw sensor inputs, we ensure that the same training and testing subjects are used in both the Sensor-Language Alignment and Task-Aware Tuning stages. This guarantees that test data in the second stage remains unseen during alignment, ensuring a fair evaluation of generalization. We select Chronos as the TS embedder because it has not been pre-trained on motion sensor data, making it an ideal candidate for evaluating our approach's robustness in learning directly from unprocessed sensor signals.

#### 4.1 Datasets

To evaluate the effectiveness and generalizability of SensorLLM, we conduct experiments on five publicly available HAR datasets: USC-HAD (Zhang and Sawchuk, 2012), UCI-HAR (Anguita et al., 2013), PAMAP2 (Reiss and Stricker, 2012), MHealth (Baños et al., 2014), and CAPTURE-24 (Chan et al., 2024). These datasets differ in sensor placement, sampling rates, channel configurations, and activity types, covering both controlled laboratory conditions and free-living environments. Additionally, they vary in scale, with USC-HAD, UCI-HAR, PAMAP2, and MHealth collected from a limited number of subjects, whereas CAPTURE-24 features a large-scale dataset with sensor recordings from 151 participants in real-world settings. Full dataset details, including subject count, sensor configurations, data splits, activity classes, preprocessing steps, and windowing strategies, are provided in Appendix A.6.

370

371

372

373

374

375

376

377

378

379

381

382

383

384

385

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

#### 4.2 Sensor Data Understanding

**Setup.** We train SensorLLM on all five datasets using the same following parameters in the Sensor-Language Alignment Stage: a learning rate of 2e-3, 8 epochs, batch size of 4, gradient accumulation steps of 8, and a maximum sequence length of 8192 for CAPTURE-24 and 4096 for the other datasets.

**Evaluation Metrics.** We assess SensorLLM's ability to generate trend descriptions from sensor data, comparing it with the advanced GPT-40<sup>1</sup> to evaluate Sensor-Language Alignment. GPT-40 generates responses based on a predefined prompt (Appendix A.4), following our text template. We employ three evaluation methods:

- NLP Metrics. We measure surface-level similarity and n-gram overlap using BLEU-1 (Papineni et al., 2002), ROUGE-1, ROUGE-L (Lin, 2004), and METEOR (Banerjee and Lavie, 2005). We also use SBERT (Reimers and Gurevych, 2019) and SimCSE (Gao et al., 2021) to evaluate semantic similarity between model outputs and ground truth.
- **GPT-4o Evaluation.** GPT-4o rates the generated trend descriptions on a scale of 1 to 5 (with 5 being the highest) by comparing each output to ground truth and providing explanatory feedback. As an advanced LLM, its evaluation ensures a semantic assessment of trend comprehension.
- Human Evaluation. Five time-series experts (PhD students, postdocs, and academics)

<sup>&</sup>lt;sup>1</sup>gpt-4o-2024-08-06(OpenAI, 2024)

	USC-H	IAD	UCI-H	IAR	PAMA	AP2	MHe	alth	CAPTU	RE-24
Metric	GPT-40	Ours								
BLEU-1	41.43	57.68	37.97	56.78	46.35	60.20	49.97	61.38	46.58	57.10
ROUGE-1	54.92	68.32	51.24	67.63	58.08	69.92	61.11	71.20	58.21	68.11
ROUGE-L	49.00	64.17	44.88	63.05	50.30	66.25	51.99	67.83	48.88	60.90
METEOR	30.51	45.95	26.93	45.81	37.17	52.21	38.50	51.73	31.16	40.51
SBERT	77.22	86.09	76.05	85.01	82.71	87.31	83.15	86.66	83.11	84.83
SimCSE	86.96	93.09	90.23	92.51	89.64	93.82	92.10	93.38	90.10	92.20
GPT-40	1.67	3.11	1.61	3.20	1.90	3.77	1.69	3.69	1.70	2.32
Human	2.10	4.16	1.94	4.04	2.38	4.70	1.74	4.56	2.30	3.10

Table 1: Evaluation of Sensor Data Trend Analysis Tasks for SensorLLM and GPT-40. The assessment includes human and GPT-40 ratings (from 1 to 5, with 5 being the highest), as well as BLEU-1, ROUGE-1, ROUGE-L, METEOR, SBERT, and SimCSE (in %). The column *GPT-40* refers to the trend analysis generated by GPT-40 itself, while the row *GPT-40* refers to GPT-40's evaluation of the generated outputs.

score accuracy and quality using the same criteria as GPT-40, providing a human-centered perspective on the model's outputs.

Appendix A.5 details all metrics and scoring criteria. We randomly sample 200 instances per dataset for both SensorLLM and GPT-40, then average the results for comparison. Because reading and comparing lengthy sequences is difficult for human annotators, we conduct human evaluation on 20 shorter sequences per dataset (each containing at most 50 time steps).

**Results.** Table 1 compares SensorLLM and GPT-40 on the Sensor Data Trend Analysis tasks, showing that our model consistently outperforms GPT-40 across all metrics. BLEU-1, ROUGE-1, ROUGE-L, and METEOR primarily focus on surface-level lexical or n-gram overlaps and SBERT and SimCSE can capture factual correctness or deeper semantic similarities. Across all metrics, SensorLLM generates trend descriptions more closely aligned with the ground truth. GPT-40 evaluations further highlight SensorLLM's superior ability to capture trend details and coherence, whereas GPT-40 struggles with complex numerical data and trend observations (Yehudai et al., 2024). Human evaluation also favors SensorLLM, particularly for shorter sequences. CAPTURE-24 results are weaker compared to other datasets, likely due to its longer sequences being trained with the same parameters. Overall, these findings validate the effectiveness of our Sensor-Language Alignment method in enhancing LLMs' ability to interpret complex numerical sequence. Appendix A.9 provides qualitative examples of outputs from both

models.

#### 4.3 Human Activity Recognition

**Setup.** In this section, we evaluate the performance of SensorLLM on HAR tasks. Each experiment runs for five trials, using 8 training epochs, a batch size of 4, gradient accumulation steps of 8, and a maximum sequence length of 4096. We report the F1 macro score A.8 to account for class imbalance across different activity categories.

**Baselines.** We benchmark SensorLLM against 11 baselines across two categories: (i) *TS models*—Transformer (Vaswani et al., 2017), Informer (Zhou et al., 2021), NS-Transformer (Liu et al., 2022), PatchTST (Nie et al., 2023), TimesNet (Wu et al., 2023a), and iTransformer (Liu et al., 2024c); (ii) *HAR models*—DeepConvLSTM (Ordóñez and Roggen, 2016), DeepConvLSTMAttn (Murahari and Plötz, 2018), and Attend (Abedin et al., 2021). We also include Chronos+MLP and GPT4TS (Zhou et al., 2023a) for a more comprehensive comparison. Full baseline details are in Appendix A.7.

**Results.** Table 2 reports the macro F1 scores (%) averaged over five random runs. SensorLLM achieves the highest performance on four datasets (USC-HAD, PAMAP2, MHealth, CAPTURE-24) and ranks second on UCI-HAR, demonstrating its effectiveness in handling diverse sensor data.

Notably, SensorLLM achieves a notable improvement on the CAPTURE-24 dataset, surpassing all baselines by a significant margin with a mean F1-macro score of 48.6%, which is 5.0% higher than Attend (43.6%). On USC-HAD, Sen-

Method	USC-HAD	UCI-HAR	PAMAP2	MHealth	CAPTURE-24
PatchTST	$ 45.2_{\pm 1.48} $	$86.8_{\pm 0.84}$	$82.0_{\pm 0.71}$	$80.0_{\pm 1.58}$	$35.6_{\pm 0.89}$
Ns-Transformer	$52.6_{\pm 2.30}$	$88.0_{\pm 0.71}$	$78.8_{\pm0.84}$	$77.2_{\pm 1.48}$	$34.8_{\pm 1.10}$
Informer	$51.2_{\pm 1.30}$	$86.6_{\pm 1.14}$	$78.0_{\pm 1.58}$	$74.0_{\pm0.71}$	$35.6_{\pm 0.55}$
Transformer	$49.6_{\pm 1.67}$	$85.4_{\pm 0.89}$	$77.0_{\pm 0.71}$	$75.2_{\pm1.30}$	$32.8_{\pm 0.84}$
iTransformer	$48.4_{\pm 1.82}$	$81.8_{\pm 0.84}$	$76.6_{\pm 0.55}$	$80.4_{\pm 1.14}$	$19.8_{\pm 0.84}$
TimesNet	$52.2_{\pm 2.39}$	$87.4_{\pm 1.14}$	$76.2_{\pm 1.92}$	$78.4_{\pm 1.52}$	$34.8_{\pm 0.84}$
GPT4TS	$54.2_{\pm 2.05}$	$88.2_{\pm 0.84}$	$80.4_{\pm 0.89}$	$76.4_{\pm 1.14}$	$32.8_{\pm 1.10}$
Chronos+MLP	$44.2_{\pm 1.30}$	$82.2_{\pm 0.84}$	$79.8_{\pm 0.45}$	$83.0_{\pm0.71}$	$38.0_{\pm 0.71}$
DeepConvLSTM	$48.8_{\pm 2.39}$	$89.2_{\pm 0.84}$	$78.4_{\pm 1.52}$	$75.0_{\pm 1.87}$	$40.4_{\pm 0.89}$
DeepConvLSTMAtt	$54.0_{\pm 2.12}$	89.6 $_{\pm 1.14}$	$79.2_{\pm 1.30}$	$77.4_{\pm 2.19}$	$41.4_{\pm 0.55}$
Attend	<u>60.2</u> $\pm 2.17$	$\textbf{93.2}_{\pm 0.84}$	$\underline{84.6}_{\pm 1.14}$	$\underline{83.4}_{\pm 1.14}$	$\underline{43.6}_{\pm 0.55}$
SensorLLM	<b>61.2</b> ±3.56	$\underline{91.2}_{\pm 1.48}$	$86.2_{\pm 1.48}$	$\textbf{89.4}_{\pm 3.85}$	<b>48.6</b> $_{\pm 1.14}$

Table 2: F1-macro results (%) for the Task-Aware Tuning Stage, presented as the mean and standard deviation over 5 random repetitions. The top results for each dataset are highlighted as follows: **Bold** for the best and <u>underline</u> for the second-best.

	Task-	only	SensorLLM		
Dataset	w/o prompts	w/ prompts	w/o prompts	w/ prompts	
USC-HAD	43.4 <sub>±2.88</sub>	$45.0_{\pm 1.58}$	$49.6_{\pm 1.67}$	$61.2_{\pm 3.56}$	
UCI-HAR	$80.0_{\pm 2.12}$	$82.0_{\pm 1.58}$	$89.2_{\pm 1.10}$	$91.2_{\pm 1.48}$	
PAMAP2	$74.2_{\pm 2.28}$	$75.4_{\pm 3.05}$	$83.0_{\pm 0.71}$	$86.2_{\pm 1.48}$	
MHealth	$76.6_{\pm 1.34}$	$77.4_{\pm 3.13}$	$86.6_{\pm 1.14}$	$89.4_{\pm 3.85}$	
CAPTURE-24	$44.8_{\pm 0.84}$	$46.0_{\pm0.71}$	$47.2{\scriptstyle\pm0.84}$	$\textbf{48.6}_{\pm 1.14}$	

Table 3: The results for SensorLLM trained with/without text prompts. *Task-only* refers to conduct-ing HAR directly bypassing sensor-language alignment.

sorLLM achieves the highest score of 61.2%, outperforming Attend, the second-best baseline, by 1.0%. Similarly, on PAMAP2, SensorLLM achieves a score of 86.2%, exceeding Attend (84.6%) by 1.6%. On MHealth, SensorLLM sets a new state-of-the-art with a score of 89.4%, surpassing Attend (83.4%) by 6.0%. These results highlight SensorLLM's ability to consistently outperform existing methods across diverse datasets.

For UCI-HAR, SensorLLM achieves the secondbest score (91.2%), slightly trailing Attend (93.2%). In contrast, Chronos+MLP shows only a slight improvement over iTransformer, the lowestperforming baseline (82.2% vs. 81.8%), indicating that Chronos embeddings alone have limited utility for HAR on this dataset. However, our framework significantly enhances their effectiveness, highlighting the robustness of our alignment approach.

#### 5 Ablation Studies

484

485

486

487

488

490

491

492

493

494

495

496

497

498

499

501

502

504

505

**Impact of Alignment.** To evaluate the impact of alignment, we included the Chronos+MLP base-line in Section 4.3 to show that SensorLLM's per-

formance is not solely driven by Chronos embeddings. Additionally, we compared SensorLLM with a Task-only model, which skips the Sensor-Language Alignment Stage and directly applies Chronos embeddings and the LLM for HAR. As shown in Table 3, SensorLLM consistently outperforms the Task-only model across all five datasets, regardless of additional textual input. Notably, the Task-only model often performs on par with or worse than traditional TS baselines, underscoring the importance of our alignment method. These results confirm that Chronos embeddings alone are insufficient for optimal HAR performance and that alignment is essential for enabling LLMs to effectively understand sensor data. 506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

Impact of Prompts. To assess the role of additional textual information (e.g., statistical features for each sensor channel) in the Task-Aware Tuning Stage, we compared SensorLLM's performance with and without prompts. As shown in Table 3, incorporating prompts consistently improves F1-macro scores across all datasets, with a more pronounced effect in the full SensorLLM architecture. This demonstrates that the model effectively integrates sensor and textual data, enhancing its ability to capture complex temporal patterns. The results highlight the benefits of multimodal inputs, which enrich sensor data representations and improve HAR accuracy. More broadly, the ability to process both sensor signals and textual prompts not only enhances classification performance but also extends LLMs' potential for tackling complex



Figure 3: Effect of the number of alignment module layers.



Figure 4: Effect of Model Size.

sensor-driven tasks in future applications.

539

540

541

542

543

544

546

547

Alignment Module Layers. We investigate the effect of hidden layer depth in the alignment module MLP on UCI-HAR, PAMAP2, and MHealth. As shown in Figure 3, increasing the number of hidden layers from one  $(1024 \rightarrow 2048 \rightarrow 4096)$  to two  $(1024 \rightarrow 2048 \rightarrow 3072 \rightarrow 4096)$  led to mixed results. The average F1-macro scores improved on UCI-HAR (91.2%  $\rightarrow$  92.0%) and MHealth (89.4%  $\rightarrow$  90.2%), but slightly dropped on PAMAP2 (86.2%  $\rightarrow$  85.8%). These results suggest that the optimal number of hidden layers varies by dataset.

Effect of Model Size on Performance. To assess 552 the impact of model size, we tested SensorLLM-3b, a resource-efficient variant using Chronos-base 553 and Llama3.2-3b. Experiments were conducted on USC-HAD, UCI-HAR, and MHealth. As shown in Figure 4, SensorLLM-3b achieves slightly lower performance than SensorLLM-8b across all datasets, illustrating the trade-off between model 558 size and accuracy. However, SensorLLM-3b re-559 mains competitive, outperforming Attend on USC-HAD and MHealth, while trailing it only on UCI-HAR. These results indicate that SensorLLM-3b offers a resource-efficient alternative, maintaining 563 strong performance relative to other baselines. 564

565 **Cross-Dataset Generalization.** To assess the 566 robustness of SensorLLM, we conducted cross-

Stage 1	Stage 2	Results
USC-HAD	UCI-HAR	$91.0_{\pm 1.41}$
UCI-HAR	USC-HAD	$61.6_{\pm 2.07}$

Table 4: Cross-dataset experiments.

567

568

570

571

572

573

574

575

576

577

578

579

581

582

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

600

601

602

603

604

605

606

607

608

dataset experiments, training the Sensor-Language Alignment Stage on USC-HAD and the Task-Aware Tuning Stage on UCI-HAR, and vice versa. These datasets share the same sensor channels but differ in sampling rates. As shown in Table 4, SensorLLM achieves performance comparable to training both stages on the same dataset. This indicates that once modality alignment is established, retraining on new datasets is unnecessary. SensorLLM effectively generalizes to downstream tasks across diverse datasets, demonstrating that the alignment stage enables the LLM to truly understand sensor data, rather than memorizing dataset-specific features. These results highlight SensorLLM's potential for broad cross-dataset generalization, laying the groundwork for future TS-LLM models.

#### 6 Conclusions

We introduced SensorLLM, a multimodal framework that aligns sensor data with automatically generated text at a human-perception level, moving beyond machine-level alignment. SensorLLM effectively captures complex sensor patterns, achieving superior performance in HAR tasks. Experiments across diverse datasets demonstrate its robustness in handling variable-length sequences, multi-channel inputs, and textual metadata. Crossdataset results further highlight its strong generalizability without requiring dataset-specific alignment. This work establishes a foundation for Sensor-Text MLLMs, with potential applications for sensor data analysis. We release our code and data generation pipeline to facilitate future research on integrating time-series and text, particularly in low-resource domains.

## 7 Limitations

While SensorLLM demonstrates strong performance in aligning sensor data with LLMs, certain limitations remain, offering directions for future exploration.

**Classifier-Based Design.** To ensure fair comparisons with existing HAR models (classification models), we adopt a classifier for downstream

- tasks rather than fully exploiting the LLM's gener-609 ative abilities. Although our results verify that the 610 Sensor-Language Alignment Stage can generalize 611 across datasets, relying on a fixed-class classifier 612 may constrain the model's adaptability to new activity categories. Future work could explore using 614 the LLM in a generative or prompt-based capac-615 ity, enabling broader application scenarios such as 616 activity discovery or open-set recognition.
- Scope of Sensor-Text Alignment. Our align-618 ment focuses on mapping sensor data to trend-619 descriptive text, demonstrating clear benefits for LLM-based HAR. However, human-intuitive de-621 scriptions of sensor data extend beyond trend 622 changes-incorporating frequency-domain features, periodicity, and higher-order patterns may further enhance an LLM's ability to interpret timeseries data. Future research could investigate whether aligning text with alternative sensor char-627 acteristics improves time-series reasoning. This 628 could expand the potential of multimodal NLP applications in sensor-driven tasks beyond activity recognition.

#### References

632

633

635

636

637 638

639

641

643

648

651

653

656

657

- Alireza Abedin, Mahsa Ehsanpour, Qinfeng Shi, Hamid Rezatofighi, and Damith C. Ranasinghe. 2021. Attend and discriminate: Beyond the state-of-the-art for human activity recognition using wearable sensors. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 5(1).
- D. Anguita, Alessandro Ghio, L. Oneto, Xavier Parra, and Jorge Luis Reyes-Ortiz. 2013. A public domain dataset for human activity recognition using smartphones. In *The European Symposium on Artificial Neural Networks*.
- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Syndar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, Jasper Zschiegner, Danielle C. Maddix, Hao Wang, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. 2024. Chronos: Learning the language of time series. arXiv preprint arXiv:2403.07815.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Oresti Baños, Rafael García, Juan Antonio Holgado Terriza, Miguel Damas, Héctor Pomares, Ignacio Rojas, Alejandro Saez, and Claudia Villalonga. 2014. mhealthdroid: A novel framework for agile development of mobile health applications. In *International Workshop on Ambient Assisted Living and Home Care*. 661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

- Antonio Bevilacqua, Kyle MacDonald, Aamina Rangarej, Venessa Widjaya, Brian Caulfield, and Tahar Kechadi. 2019. *Human Activity Recognition with Convolutional Neural Networks*, page 541–552. Springer International Publishing.
- Defu Cao, Furong Jia, Sercan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu. 2024. Tempo: Prompt-based generative pre-trained transformer for time series forecasting. *Preprint*, arXiv:2310.04948.
- Shing Chan, Hang Yuan, Catherine Tong, Aidan Acquah, Abram Schonfeldt, Jonathan Gershuny, and Aiden Doherty. 2024. Capture-24: A large dataset of wrist-worn activity tracker data collected in the wild for human activity recognition. *Preprint*, arXiv:2402.19229.
- Ching Chang, Wei-Yao Wang, Wen-Chih Peng, and Tien-Fu Chen. 2024. Llm4ts: Aligning pretrained llms as data-efficient time-series forecasters. *Preprint*, arXiv:2308.08469.
- Shohreh Deldari, Dimitris Spathis, Mohammad Malekzadeh, Fahim Kawsar, Flora D. Salim, and Akhil Mathur. 2024. Crossl: Cross-modal self-supervised learning for time-series through latent masking. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, WSDM '24, page 152–160.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Empirical Methods in Natural Language Processing (EMNLP)*.
- Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. 2024. Large language models are zero-shot time series forecasters. *Advances in Neural Information Processing Systems*, 36.
- Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew Gordon Wilson. 2023. Large language models are zero-shot time series forecasters. *Preprint*, arXiv:2310.07820.
- Yu Guan and Thomas Plötz. 2017. Ensembles of deep lstm learners for activity recognition using wearables. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 1(2).
- Sojeong Ha and Seungjin Choi. 2016. Convolutional neural networks for human activity recognition using multiple accelerometer and gyroscope sensors. In 2016 International Joint Conference on Neural Networks (IJCNN), pages 381–388.

818

819

820

821

822

823

771

Nils Y. Hammerla, Shane Halloran, and Thomas Plötz. 2016. Deep, convolutional, and recurrent models for human activity recognition using wearables. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, IJCAI'16, page 1533–1540. AAAI Press.

715

716

717

719

721

724

725

726

727

728

729

730

731

734

735

738

739

740

741

742

743

744

745

746

747

748

749

754

755

756

758

759

761

764

770

- Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao, Ao Zhang, Liang Zhang, Wentao Han, Minlie Huang, Qin Jin, Yanyan Lan, Yang Liu, Zhiyuan Liu, Zhiwu Lu, Xipeng Qiu, Ruihua Song, and 5 others. 2021. Pretrained models: Past, present and future. *Preprint*, arXiv:2106.07139.
- Harish Haresamudram, David V. Anderson, and Thomas Plötz. 2019. On the role of features in human activity recognition. In *Proceedings of the 2019 ACM International Symposium on Wearable Computers*, ISWC '19, page 78–88, New York, NY, USA. Association for Computing Machinery.
- Harish Haresamudram, Apoorva Beedu, Mashfiqui Rabbi, Sankalita Saha, Irfan Essa, and Thomas Ploetz. 2024. Limitations in employing natural language supervision for sensor-based human activity recognition–and ways to overcome them. *arXiv preprint arXiv:2408.12023*.
- Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415.*
- Sijie Ji, Xinzhe Zheng, and Chenshu Wu. 2024. Hargpt: Are llms zero-shot human activity recognizers? *Preprint*, arXiv:2403.02727.
- Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y Zhang, Xiaoming Shi, Pin-Yu Chen, Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. 2024a. Time-LLM: Time series forecasting by reprogramming large language models. In *International Conference on Learning Representations (ICLR)*.
- Ming Jin, Qingsong Wen, Yuxuan Liang, Chaoli Zhang, Siqiao Xue, Xue Wang, James Zhang, Yi Wang, Haifeng Chen, Xiaoli Li, Shirui Pan, Vincent S. Tseng, Yu Zheng, Lei Chen, and Hui Xiong. 2023. Large models for time series and spatiotemporal data: A survey and outlook. *Preprint*, arXiv:2310.10196.
- Ming Jin, Yifan Zhang, Wei Chen, Kexin Zhang, Yuxuan Liang, Bin Yang, Jindong Wang, Shirui Pan, and Qingsong Wen. 2024b. Position paper: What can large language models tell us about time series analysis. *Preprint*, arXiv:2402.02713.
- Panagiotis Kasnesis, Charalampos Z. Patrikakis, and Iakovos S. Venieris. 2019. Perceptionnet: A deep convolutional neural network for late sensor fusion. In *Intelligent Systems and Applications*, pages 101– 119, Cham. Springer International Publishing.
- Yubin Kim, Xuhai Xu, Daniel McDuff, Cynthia Breazeal, and Hae Won Park. 2024. Health-llm:

Large language models for health prediction via wearable sensor data. In *Proceedings of the fifth Conference on Health, Inference, and Learning*, volume 248 of *Proceedings of Machine Learning Research*, pages 522–539. PMLR.

- Jennifer R. Kwapisz, Gary M. Weiss, and Samuel A. Moore. 2011. Activity recognition using cell phone accelerometers. *SIGKDD Explor. Newsl.*, 12(2):74–82.
- Hong Li, Gregory D. Abowd, and Thomas Plötz. 2018. On specialized window lengths and detector based human activity recognition. In *Proceedings of the* 2018 ACM International Symposium on Wearable Computers, ISWC '18, page 68–71, New York, NY, USA. Association for Computing Machinery.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Che Liu, Zhongwei Wan, Sibo Cheng, Mi Zhang, and Rossella Arcucci. 2024a. Etp: Learning transferable ecg representations via ecg-text pre-training. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8230–8234. IEEE.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning.
- Xin Liu, Daniel McDuff, Geza Kovacs, Isaac Galatzer-Levy, Jacob Sunshine, Jiening Zhan, Ming-Zher Poh, Shun Liao, Paolo Di Achille, and Shwetak Patel. 2023b. Large language models are few-shot health learners. *Preprint*, arXiv:2305.15525.
- Xu Liu, Junfeng Hu, Yuan Li, Shizhe Diao, Yuxuan Liang, Bryan Hooi, and Roger Zimmermann. 2024b. Unitime: A language-empowered unified model for cross-domain time series forecasting. *Preprint*, arXiv:2310.09751.
- Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. 2024c. itransformer: Inverted transformers are effective for time series forecasting. *International Conference on Learning Representations*.
- Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. 2022. Non-stationary transformers: Exploring the stationarity in time series forecasting.
- Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Aparajita Saraf, Amy Bearman, and Babak Damavandi. 2023. IMU2CLIP: Language-grounded motion sensor translation with multimodal contrastive learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13246– 13253, Singapore. Association for Computational Linguistics.

915

916

917

918

919

920

921

922

923

924

925

926

927

928

876

877

Vishvak S. Murahari and Thomas Plötz. 2018. On attention models for human activity recognition. In *Proceedings of the 2018 ACM International Symposium on Wearable Computers*, ISWC '18, page 100–103, New York, NY, USA. Association for Computing Machinery.

824

825

828

831

832

833

834

835

837

838

841

844

849

850

851

852

853

855

870

871

872

874

875

- Shikai Qiu Nate Gruver, Marc Finzi and Andrew Gordon Wilson. 2023. Large Language Models Are Zero Shot Time Series Forecasters. In Advances in Neural Information Processing Systems.
- Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. 2023. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations*.
- OpenAI. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- Francisco Javier Ordóñez and Daniel Roggen. 2016. Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition. *Sensors*, 16(1).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting on Association for Computational Linguistics, ACL '02, page 311–318, USA. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Conference on Empirical Methods in Natural Language Processing.
- Attila Reiss and Didier Stricker. 2012. Introducing a new benchmarked dataset for activity monitoring. In 2012 16th International Symposium on Wearable Computers, pages 108–109.
- Dimitris Spathis and Fahim Kawsar. 2024. The first step is the hardest: Pitfalls of representing and tokenizing temporal data for large language models. *Journal of the American Medical Informatics Association*, 31(9):2151–2158.
- Chenxi Sun, Hongyan Li, Yaliang Li, and Shenda Hong. 2024a. Test: Text prototype aligned embedding to activate llm's ability for time series.

- Chenxi Sun, Hongyan Li, Yaliang Li, and Shenda Hong. 2024b. Test: Text prototype aligned embedding to activate llm's ability for time series. *Preprint*, arXiv:2308.08241.
- Mingtian Tan, Mike A. Merrill, Vinayak Gupta, Tim Althoff, and Thomas Hartvigsen. 2024. Are language models actually useful for time series forecasting? *Preprint*, arXiv:2406.16964.
- 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. 2023. Llama: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. 2023a. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *International Conference on Learning Representations*.
- Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023b. Next-gpt: Any-to-any multimodal llm.
- Kang Xia, Wenzhong Li, Shiwei Gan, and Sanglu Lu. 2024. Ts2act: Few-shot human activity sensing with cross-modal co-learning. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 7(4):1–22.
- Qingxin Xia, Takuya Maekawa, and Takahiro Hara. 2023. Unsupervised human activity recognition through two-stage prompting with chatgpt. *Preprint*, arXiv:2306.02140.
- Cheng Xu, Duo Chai, Jie He, Xiaotong Zhang, and Shihong Duan. 2019. Innohar: A deep neural network for complex human activity recognition. *IEEE Access*, 7:9893–9902.
- Hao Xue and Flora D Salim. 2023. Promptcast: A new prompt-based learning paradigm for time series forecasting. *IEEE Transactions on Knowledge and Data Engineering*.
- Gilad Yehudai, Haim Kaplan, Asma Ghandeharioun, Mor Geva, and Amir Globerson. 2024. When can transformers count to n? *Preprint*, arXiv:2407.15160.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*.

929Hyungjun Yoon, Biniyam Aschalew Tolera, Taesik930Gong, Kimin Lee, and Sung-Ju Lee. 2024. By my931eyes: Grounding multimodal large language models932with sensor data via visual prompting. In Proceed-933ings of the 2024 Conference on Empirical Methods934in Natural Language Processing, pages 2219–2241,935Miami, Florida, USA. Association for Computational936Linguistics.

938

939 940

941

942

943 944

945

946

947

951

952

953

954 955

956

957

959

960

961

962

963 964

965

966

967

- Yuta Yuki, Junto Nozaki, Kei Hiroi, Katsuhiko Kaji, and Nobuo Kawaguchi. 2018. Activity recognition using dual-convlstm extracting local and global features for shl recognition challenge. In Proceedings of the 2018 ACM International Joint Conference and 2018 International Symposium on Pervasive and Ubiquitous Computing and Wearable Computers, UbiComp '18, page 1643–1651, New York, NY, USA. Association for Computing Machinery.
  - Mi Zhang and Alexander A. Sawchuk. 2012. Usc-had: a daily activity dataset for ubiquitous activity recognition using wearable sensors. In *Proceedings of the 2012 ACM Conference on Ubiquitous Computing*, UbiComp '12, page 1036–1043, New York, NY, USA. Association for Computing Machinery.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 2021. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Virtual Conference*, volume 35, pages 11106–11115. AAAI Press.
  - Tian Zhou, Peisong Niu, Xue Wang, Liang Sun, and Rong Jin. 2023a. One Fits All: Power general time series analysis by pretrained lm. In *NeurIPS*.
  - Tian Zhou, PeiSong Niu, Xue Wang, Liang Sun, and Rong Jin. 2023b. One fits all:power general time series analysis by pretrained lm. *Preprint*, arXiv:2302.11939.
  - Yunjiao Zhou, Jianfei Yang, Han Zou, and Lihua Xie. 2023c. Tent: Connect language models with iot sensors for zero-shot activity recognition. *arXiv preprint arXiv:2311.08245*.

#### 970 A

971

973

974

975

977

979

981

987

988

991

992

994

997

998

1000

1002

1003

# A.1 More related work

Appendix

Deep learning in human activity recognition. Over the last decade, HAR has transitioned from hand-crafted feature extraction to deep learning models capable of automatic feature learning. Early work by Kwapisz et al. (2011) utilized machine learning techniques, such as decision trees and MLPs, to classify activities using features extracted from wearable sensor data. Later, Haresamudram et al. (2019) demonstrated that optimized feature extraction within the Activity Recognition Chain (ARC) could rival or outperform endto-end deep learning models. Deep learning models, particularly CNNs and LSTMs, have since become dominant in HAR. Bevilacqua et al. (2019) developed a CNN-based model for HAR, while Ha and Choi (2016) introduced CNN-pf and CNNpff architectures that apply partial and full weight sharing for better feature extraction. Other notable works include Perception-Net Kasnesis et al. (2019), which leverages 2D convolutions for multimodal sensor data, and InnoHAR (Xu et al., 2019), which combines Inception CNN and GRUs for multiscale temporal feature learning. A dual-stream network utilizing convolutional layers and LSTM units, known as ConvLSTM, was employed by Yuki et al. (2018) to analyze complex temporal hierarchies with streams handling different time lengths. The combination of attention mechanisms with recurrent networks to enhance the computation of weights for hidden state outputs has also been demonstrated by DeepConvLSTM (Kasnesis et al., 2019) in capturing spatial-temporal features.

Large Language Models for Time-Series Fore-1004 casting. LLMs have achieved remarkable suc-1005 cess in text-related tasks, and their utility has ex-1006 panded into time-series forecasting. Xue and Salim 1007 (2023) presents PromptCast, which redefines time-1008 series forecasting as a natural language generation 1009 task by transforming numerical inputs into textual prompts, enabling pre-trained language models to 1011 handle forecasting tasks with superior generaliza-1012 tion in zero-shot settings. Gruver et al. (2023) ex-1013 plores encoding time-series as numerical strings, 1014 1015 allowing LLMs like GPT-3 and LLaMA-2 to perform zero-shot forecasting, matching or surpass-1016 ing the performance of specialized models, while 1017 highlighting challenges in uncertainty calibration due to model modifications like RLHF. Zhou et al. 1019

(2023b) demonstrates that pre-trained language and 1020 image models, such as a Frozen Pretrained Trans-1021 former (FPT), can be adapted for diverse time-1022 series tasks like classification, forecasting, and 1023 anomaly detection, leveraging self-attention mecha-1024 nisms to bridge the gap between different data types 1025 and achieving state-of-the-art performance across 1026 various tasks. Jin et al. (2024b) highlights the trans-1027 formative potential of LLMs for time-series analy-1028 sis by integrating language models with traditional 1029 analytical methods. Jin et al. (2024a) introduces a 1030 reprogramming framework that aligns time-series 1031 data with natural language processing capabilities, 1032 enabling LLMs to perform time-series forecast-1033 ing without altering the core model structure. Cao 1034 et al. (2024) presents TEMPO, a generative trans-1035 former framework based on prompt tuning, which 1036 adapts pre-trained models for time-series forecast-1037 ing by decomposing trends, seasonality, and resid-1038 ual information. Sun et al. (2024b) proposes TEST, 1039 an innovative embedding technique that integrates 1040 time-series data with LLMs through instance-wise, 1041 feature-wise, and text-prototype-aligned contrast, 1042 vielding improved or comparable results across var-1043 ious applications. Chang et al. (2024) develops 1044 a framework that enhances pre-trained LLMs for 1045 multivariate time-series forecasting through a two-1046 stage fine-tuning process and a novel multi-scale 1047 temporal aggregation method, outperforming tradi-1048 tional models in both full-shot and few-shot scenar-1049 ios. Finally, Liu et al. (2024b) introduces UniTime, 1050 a unified model that leverages language instruc-1051 tions and a Language-TS Transformer to handle 1052 multivariate time series across different domains, 1053 demonstrating enhanced forecasting performance 1054 and zero-shot transferability. 1055

LLMs for Human Activity Recognition. While 1056 LLMs like ChatGPT have demonstrated remark-1057 able performance in various NLP tasks, their effec-1058 tiveness in HAR remains limited due to challenges 1059 in interpreting sensor data. These models often 1060 struggle to distinguish between activities that share similar objects, requiring more advanced prompt 1062 engineering to highlight activity-specific details. 1063 (Xia et al., 2023) proposed an unsupervised ap-1064 proach to HAR using ChatGPT, leveraging two-1065 stage prompts to infer activities from object se-1066 quences without manual descriptions. The method 1067 demonstrates superior performance on three bench-1068 mark datasets, marking a significant advancement 1069 in applying language models to activity recognition 1070 1071tasks. Similarly, Ji et al. (2024) explored LLMs1072for zero-shot HAR using raw IMU data, showing1073that GPT-4 can outperform both traditional and1074deep learning models in simple HAR tasks without1075domain-specific adaptations, highlighting LLMs'1076potential in sensor-based systems.

#### A.2 Data Generation

1079

1080

1081

1082

1083

1084

1086

1087

1088 1089

1090

1091

1092

1093

1094

1095

1096 1097

1098

1099

1100

1101

We generate text data from sensor readings using predefined sentence templates (5, Tables 6,, 7). These templates are randomly selected to create diverse question-answer (QA) pairs. To enhance variability, we employ GPT-40 to generate synonymous variations. Each sentence contains placeholders for numerical values (e.g., timestamps, sensor readings) or textual information, which are dynamically replaced to produce coherent QA pairs aligned with the sensor data.

#### **Trend Description Templates**

- {start\_time}s to {end\_time}s: {trend}
- {start\_time} seconds to {end\_time} seconds: {trend}
- {start\_time} to {end\_time} seconds: {trend}
- {start\_time}-{end\_time} seconds: {trend}
- {start\_time}-{end\_time}s: {trend}
- {start\_time}s-{end\_time}s: {trend}

Table 5: Examples of answer templates used for trend descriptions.

The system prompt instructs the model on how to respond to generated questions, incorporating dataset-specific attributes such as sensor frequency and sampling rate. These tailored prompts ensure responses align with the unique characteristics of each dataset. Below is the system prompt template used for all datasets:

• A dialogue between a researcher and an AI assistant. The AI analyzes a sensor timeseries dataset (*N* points, sampled at {sample\_rate}Hz) to answer specific questions, demonstrating its analytical capabilities and the potential for human-AI collaboration in interpreting sensor data.

#### A.3 Chronos

Chronos (Ansari et al., 2024) is a pretrained proba-1103 bilistic time-series framework that tokenizes real-1104 valued time-series data into discrete representations 1105 for language model training. It utilizes scaling and 1106 quantization to transform time-series data into a 1107 fixed vocabulary, enabling T5-based (Raffel et al., 1108 2020) models to learn from tokenized sequences us-1109 ing cross-entropy loss. Pretrained on diverse public 1110 and synthetic datasets, Chronos surpasses exist-1111 ing models on familiar datasets and demonstrates 1112 strong zero-shot performance on unseen tasks, mak-1113 ing it a versatile tool for time-series forecasting 1114 across domains. 1115

1102

1116

1117

1118

1119

1120

1121

**Time-Series Tokenization and Quantization.** Chronos converts time-series data into discrete tokens through a two-step process: normalization and quantization. Mean scaling is first applied to ensure consistency across different time series:

$$\tilde{x} = \frac{x}{\operatorname{mean}(|x|)} \tag{5}$$

Next, the normalized values are quantized using1122B bin centers  $c_1, \ldots, c_B$  and corresponding bin1123edges  $b_1, \ldots, b_{B-1}$ , mapping real values to discrete1124tokens via:1125

$$q(x) = \begin{cases} 1 & \text{if } -\infty \le x < b_1, \\ 2 & \text{if } b_1 \le x < b_2, \\ \vdots & & \\ B & \text{if } b_{B-1} \le x < \infty. \end{cases}$$
(6) 112

Special tokens such as PAD and EOS are added1127to handle sequence padding and denote the end of1128sequences, allowing Chronos to process variable-1129length inputs efficiently within language models.1130

**Objective Function.** Chronos models the tok-<br/>enized time series using a categorical distribution1131over the vocabulary  $V_{ts}$ , minimizing the cross-<br/>entropy loss:1132

$$\ell(\theta) = -\sum_{h=1}^{H+1} \sum_{i=1}^{|V_{ts}|} \mathbf{1}(z_{C+h+1} = i)$$

$$\cdot \log p_{\theta}(z_{C+h+1} = i \mid z_{1:C+h})$$
(7) 1135

where C is the historical context length, H is 1136 the forecast horizon, and  $p_{\theta}$  is the predicted token 1137 distribution. 1138

## **Trend Description Templates**

- Kindly provide a detailed analysis of the trend changes observed in the {data}.
- Please offer a comprehensive description of how the trends in the {data} have evolved.
- I would appreciate a thorough explanation of the trend fluctuations that occurred within the  $\{data\}$ .
- Could you examine the {data} in depth and explain the trend shifts observed step by step?
- Detail the {data}'s trend transitions.
- Could you assess the {data} and describe the trend transformations step by step?
- Could you analyze the trends observed in the {data} over the specified period step by step?
- Can you dissect the {data} and explain the trend changes in a detailed manner?
- What trend changes can be seen in the {data}?

# **Summary Templates**

- Could you provide a summary of the main features of the input {data} and the distribution of the trends?
- Please give an overview of the essential attributes of the input {data} and the spread of the trends.
- Describe the salient features and trend distribution within the {data}.
- Give a summary of the {data}'s main elements and trend apportionment.
- Summarize the {data}'s core features and trend dissemination.
- Outline the principal aspects and trend allocation of the {data}.
- Summarize the key features and trend distribution of the {data}.
- I need a summary of {data}'s main elements and their trend distributions.

Table 6: Examples of question templates used for trend description and summary generation.

# **Summary 1: Trend Count**

- Number of {trend} trends: {num}
- Count of {trend} trends: {num}
- Number of {trend} segments: {num}
- Count of {trend} segments: {num}

# Summary 2: Sensor Data Context

- The given {data\_name} represents {sensor\_name} sensor readings from {start\_time}s to {end\_time}s.
- The {data\_name} contains {sensor\_name} sensor readings recorded between {start\_time} and {end\_time} seconds.
- The {sensor\_name} sensor readings collected from {start\_time} to {end\_time} seconds are presented in this {data\_name}.

#### **Summary 3: Trend Change Statistics**

- The data exhibits {trend\_num} distinct trends, with {change\_num} trend changes observed.
- Across {trend\_num} trends, the data shows {change\_num} occurrences of trend shifts.
- {trend\_num} trends are present, with {change\_num} instances of trend changes.

#### **Summary 4: Cumulative Trend Analysis**

- To sum up, the data exhibited a {trend\_type} trend for a total duration of {total\_time} seconds.
- Overall, the data showed a {trend\_type} trend spanning {total\_time} seconds.
- In conclusion, the trend was {trend\_type} over {total\_time} seconds.

#### **Summary 5: Overall Trend Summary**

- The overall trend is {overall\_trend}.
- The primary trend detected is {overall\_trend}.
- Looking at the broader pattern, the trend is {overall\_trend}.

Table 7: Examples of answer templates used for summaries.

1139This approach offers two key advantages: (i)1140Seamless integration with language models, requir-1141ing no architectural modifications, and (ii) Flexible1142distribution learning, enabling robust generaliza-1143tion across diverse time-series datasets.

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

# A.4 GPT-40 Prompt for Sensor Data Trend Analysis

Table 8 presents the system prompt used to generate trend-descriptive texts from sensor data, providing a structured framework for GPT-40 to analyze and respond to specific questions. This standardized prompt ensures consistency in GPT-40's interpretation of time-series data, allowing direct comparison with descriptions produced by SensorLLM.

> **Prompt** A dialogue between a curious researcher and an AI assistant. The AI analyzes a sensor time-series dataset (N points, {sr}Hz sampling rate) to answer specific questions.

> > Please output your answer in the format like this example: {example from ground-truth}

> > Now, analyze the following: Input: {sensor\_data} How trends in the given sensor data evolve? Output:

Table 8: Prompt for GPT-40 to generate descriptive texts based on the given numerical sensor data.

We evaluate GPT-4o's ability to interpret numerical sensor data by assessing its responses against human evaluations and NLP metrics. This comparison benchmarks GPT-4o's performance against SensorLLM, highlighting differences in how both models process time-series data trends. The results demonstrate the effectiveness of SensorLLM's Sensor-Language Alignment Stage.

# A.5 Evaluation Metrics for Sensor-Language Alignment Stage

1163In this section, we describe the various evaluation1164metrics used to assess the performance of Sensor-1165LLM in generating trend descriptions from sen-1166sor data. Each metric offers a distinct perspective1167on model performance, ranging from surface-level1168textual similarity to more complex semantic alignment.

BLEU-1 (Papineni et al., 2002). BLEU (Bilin-1170 gual Evaluation Understudy) is a precision-based 1171 metric commonly used to evaluate machine-1172 generated text by comparing it to reference texts. 1173 BLEU-1 focuses on unigram (single-word) overlap, 1174 assessing the lexical similarity between the gener-1175 ated and reference text. While useful for measur-1176 ing word-level matches, BLEU-1 does not capture 1177 deeper semantic meaning, making it most effective 1178 for surface-level alignment. 1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

**ROUGE-1 and ROUGE-L (Lin, 2004).** ROUGE (Recall-Oriented Understudy for Gisting Evaluation) evaluates the recall-oriented overlap between generated text and reference text. ROUGE-1 focuses on unigram recall, similar to BLEU-1 but emphasizing how much of the reference text is captured. ROUGE-L measures the longest common subsequence, assessing both precision and recall in terms of structure and content overlap, though it does not evaluate semantic accuracy.

**METEOR (Banerjee and Lavie, 2005).** ME-TEOR (Metric for Evaluation of Translation with Explicit Ordering)combines precision and recall, with additional alignment techniques such as stemming and synonym matching. Unlike BLEU and ROUGE, METEOR accounts for some degree of semantic similarity. However, its emphasis is still on word-level alignment rather than factual accuracy or meaning.

**SBERT** (Reimers and Gurevych, 2019). SBERT (Sentence-BERT) <sup>2</sup> is a metric that generates sentence embeddings using the BERT architecture. It computes cosine similarity between embeddings of the generated and reference texts, providing a deeper assessment of semantic similarity beyond lexical matches.

**SimCSE (Gao et al., 2021).** SimCSE (Simple Contrastive Sentence Embedding) <sup>3</sup> introduces a contrastive learning approach to fine-tune language models for sentence embeddings. By applying different dropout masks to the same sentence, it generates positive examples, encouraging similar embeddings for semantically identical sentences while distinguishing different ones.

 $<sup>^{2}</sup> https://huggingface.co/sentence-transformers/all-mpnet-base-v2$ 

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/princeton-nlp/sup-simcse-roberta-large

**GPT-40 Evaluation.** In addition to the NLP met-1215 rics, we also employed GPT-40 as a human-like 1216 evaluator. Given its strong reasoning and compre-1217 hension abilities, GPT-40 was tasked with scoring 1218 the generated text based on its alignment with the ground truth. GPT-40 evaluated the correctness, 1220 completeness, and coherence of the trend descrip-1221 tions and assigned a score from 1 to 5, accompanied 1222 by an explanation (see Table 9). This type of evalu-1223 ation provides insights into how well the generated 1224 outputs capture the nuances of sensor data trends 1225 in a manner similar to human understanding. 1226

Human Evaluation. Finally, five human experts assessed the correctness and quality of the generated trend descriptions. Following the same criteria as GPT-40, they rated the outputs on a scale from 1 to 5, focusing on the factual accuracy and coherence of the descriptions. This manual evaluation serves as an important benchmark for the model's performance from a human perspective, ensuring that the generated outputs are not only technically correct but also practically useful for human interpretation.

#### A.6 Datasets

1227

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1261

1263

We used five datasets in our study:

USC Human Activity Dataset (USC-HAD). USC-HAD (Zhang and Sawchuk, 2012) consists of six sensor readings from body-worn 3-axis accelerometers and gyroscopes, collected from 14 subjects. The data is sampled at 100 Hz across six channels and includes 12 activity class labels. For evaluation, we use data from subjects 13 and 14 as the test set, while the remaining subjects' data are used for training. A window size  $w \in [5, 200]$  is used in alignment stage, and w = 200 with stride of 100 are used in HAR.

UCI Human Activity Recognition Dataset (UCI-HAR). UCI-HAR (Anguita et al., 2013) includes data collected from 30 volunteers performing six activities while wearing a smartphone on their waist. The embedded accelerometer and gyroscope sensors sampled data at 50 Hz across six channels. The dataset was partitioned into 70% for training and 30% for testing. A window size  $w \in [5, 200]$  is used in alignment stage, and w = 128 with stride of 64 is used in HAR.

Physical Activity Monitoring Dataset (PAMAP2). PAMAP2 (Reiss and Stricker, 2012) includes data from nine subjects wearing

IMUs on their chest, hands, and ankles. IMUs capture the acceleration, gyroscope, and magnetometer data across 27 channels and include 12 activity class labels. For our experiments, data from subjects 105 and 106 are used as the test set, with the remaining subjects' data used for training. The sample rate is downsampled from 100 Hz to 50 Hz. A window size  $w \in [5, 100]$  is used in alignment stage, and w = 100 with stride of 50 in HAR.

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1283

1284

1285

1286

1287

1289

1290

1291

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1305

1306

1310

Mobile Health Dataset (MHealth). MHealth (Baños et al., 2014) contains body motion and vital sign recordings from ten volunteers. Sensors were placed on the chest, right wrist, and left ankle of each subject. For our experiments, we used acceleration data from the chest, left ankle, and right lower arm, along with gyroscope data from the left ankle and right lower arm, resulting in a total of 15 channels. The data is sampled at 50 Hz and includes 12 activity class labels. Data from subjects 1, 3, and 6 is used as the test set, while the remaining subjects' data are used for training. We use a window size  $w \in [5, 100]$  in alignment stage and w = 100 with stride of 50 in HAR.

CAPTURE-24. CAPTURE-24 (Chan et al., 2024) is a large-scale dataset featuring 3-channel wrist-worn accelerometer data collected in freeliving settings for over 24 hours per participant. It includes annotated data from 151 participants, making it significantly larger than existing datasets. We used the first 100 participants as the training set and the remaining 51 as the test set. For each subject, sequences were windowed, and 5% of the data was randomly selected for training and testing. The sample rate was downsampled from 100 Hz to 50 Hz and it includes 10 activity class labels. During the alignment stage, we used a variable window size  $w \in [10, 500]$ , while in the HAR, we fixed w = 500 with a stride of 250.

Each dataset includes multiple activity classes, and the proportion of each class in the dataset is shown in Table 10.

#### A.7 Baselines for Task-Aware Tuning Stage

In Task-Aware Tuning Stage, we compare Sensor-1307 LLM against several state-of-the-art baseline mod-1308 els for time-series classification and human activity 1309 recognition (HAR). These models were selected for their strong performance in relevant tasks, pro-1311 viding a thorough benchmark for evaluating Sen-1312 sorLLM's effectiveness. 1313

Prompt	Please evaluate the model-generated trend descriptions against the ground truth. Rate each pair based on the degree of accuracy, using a scale from 1 to 5, where 1 represents the lowest correctness and 5 represents the highest. Deduct 1 point for minor errors in the trend description, and 2-3 points for moderate errors.
	Provide your score (1-5) and a brief explanation in the format: "score#reason" (e.g., 4#The description of trend changes slightly differs from the ground truth).
	Now, please proceed to score the following: Model: {model_output} Human: {ground_truth} Output:
Output example 1:	2#Significant discrepancies in segment durations and trend counts com- pared to ground-truth.
Output example 2:	5#The model's description matches the human-generated text accurately.

Table 9: Prompt and output examples for GPT-40 in evaluating model-generated texts and ground-truth.

Dataset	# Classes	Classes	Proportions (%)
USC-HAD	12	Sleeping, Sitting, Elevator down, Elevator up, Standing, Jumping, Walking downstairs, Walking right, Walking forward, Running forward, Walking upstairs, Walking left	12.97, 9.06, 6.04, 5.94, 8.6, 3.62, 7.61, 9.81, 13.15, 5.72, 8.22, 9.25
UCI-HAR	6	Standing, Sitting, Laying, Walking, Walking downstairs, Walking upstairs	18.69, 17.49, 19.14, 16.68, 13.41, 14.59
PAMAP2	12	Lying, Sitting, Standing, Ironing, Vacuum cleaning, Ascending stairs, Descending stairs, Walking, Nordic walking, Cycling, Running, Rope jumping	10.25, 9.52, 10.11, 11.82, 9.14, 6.3, 5.67, 12.77, 9.52, 8.42, 3.57, 2.91
MHealth	12	Climbing stairs, Standing still, Sitting and relaxing, Lying down, Walking, Waist bends forward, Frontal elevation of arms, Knees bending (crouching), Jogging, Running, Jump front & back, Cycling	8.91, 8.95, 8.95, 8.95, 8.95, 8.26, 8.7, 8.53, 8.95, 8.95, 2.96, 8.95
CAPTURE-24	10	Sleep, Household-chores, Walking, Vehicle, Standing, Mixed-activity, Sitting, Bicycling, Sports, Manual-work	37.45, 6.5, 6.16, 3.83, 3.25, 3.49, 37.07, 1.03, 0.43, 0.79

Table 10: Dataset classes and Proportions

**Transformer (Vaswani et al., 2017).** The Trans-1314 former model is a widely-used architecture in var-1315 ious tasks, including time-series forecasting and 1316 classification. It uses self-attention mechanisms to 1317 capture long-range dependencies in sequential data, 1318 making it highly effective for modeling complex 1319 temporal relationships. 1320

Informer (Zhou et al., 2021). Informer is a 1321 transformer-based model designed for long se-1322 quence time-series data. It addresses key limita-1323 tions of standard Transformers, such as high time 1324 complexity and memory usage, through three inno-1325 vations: ProbSparse self-attention, which reduces 1326 time complexity; self-attention distilling, which 1327 enhances efficiency by focusing on dominant pat-1328 terns; and a generative decoder that predicts entire 1329 sequences in a single forward pass. 1330

NS-Transformer (Liu et al., 2022). Non-1331 1332 stationary Transformers (NS-Transformer) tackles the issue of over-stationarization in time-series by 1333 balancing series predictability and model capability. 1334 It introduces Series Stationarization to normalize 1335 inputs and De-stationary Attention to restore in-1336 trinsic non-stationary information into temporal 1337 dependencies. 1338

PatchTST (Nie et al., 2023). PatchTST is a Transformer-based model for multivariate time se-1340 ries tasks, using subseries-level patches as input to-1341 kens and a channel-independent approach to reduce 1342 computation and improve efficiency. This design 1343 retains local semantics and allows for longer his-1344 torical context, significantly improving long-term 1345 forecasting accuracy. 1346

TimesNet (Wu et al., 2023a). TimesNet is a ver-1347 1348 satile backbone for time series analysis that transforms 1D time series into 2D tensors to better cap-1349 ture intraperiod and interperiod variations. This 2D 1350 transformation allows for more efficient modeling 1351 using 2D kernels. It also introduces TimesBlock to 1352 adaptively discovers multi-periodicity and extracts 1353 temporal features from transformed 2D tensors us-1354 ing a parameter-efficient inception block. 1355

iTransformer (Liu et al., 2024c). iTransformer 1356 1357 reimagines the Transformer architecture by applying attention and feed-forward networks to inverted 1358 dimensions. Time points of individual series are 1359 embedded as variate tokens, allowing the attention mechanism to capture multivariate correlations, 1361

while the feed-forward network learns nonlinear representations for each token.

1362

1363

1364

1365

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

1402

1407

1408

DeepConvLSTM (Ordóñez and Roggen, 2016). DeepConvLSTM integrates four consecutive convolutional layers followed by two LSTM layers to effectively capture both spatial and temporal dynamics in sensor data. The final output vector is passed through a fully connected layer, and the softmax function is applied to produce activity class probabilities as the model's final output.

DeepConvLSTMAttn (Murahari and Plötz, **2018).** DeepConvLSTMAttn enhances the original DeepConvLSTM by integrating an attention mechanism to improve temporal modeling in HAR tasks. Instead of using the last LSTM hidden state for classification, the attention mechanism is applied to the first 7 hidden states, representing historical temporal context. These states are transformed through linear layers to generate attention scores, which are passed through softmax to produce weights. The weighted sum of the hidden states is combined with the last hidden state to form the final embedding for classification.

Attend (Abedin et al., 2021). The Attend model use the latent relationships between multi-channel sensor modalities and specific activities, apply dataagnostic augmentation to regularize sensor data streams, and incorporate a classification loss criterion to minimize intra-class representation differences while maximizing inter-class separability. These innovations result in more discriminative activity representations, significantly improving HAR performance.

Chronos+MLP. Chronos (Ansari et al.. 2024)+MLP is a baseline designed to evaluate whether the performance gains in SensorLLM are 1397 solely attributable to Chronos and the MLP. In 1398 SensorLLM, Chronos is used to generate sensor 1399 embeddings, which are then mapped by the MLP 1400 for input into the LLM to perform HAR. Since 1401 Chronos does not natively support classification tasks and only processes single-channel data, we 1403 adapt it for HAR by inputting each channel's 1404 data separately into Chronos. The resulting 1405 sensor embeddings for all channels are then 1406 concatenated and fed into an MLP, which acts as a classifier. This setup allows us to benchmark against a simpler framework and validate the 1409 unique contributions of SensorLLM's design. 1410

GPT4TS (Zhou et al., 2023a). GPT4TS is a uni-1411 fied framework that leverages a frozen pre-trained 1412 language model (e.g., GPT-2 (Radford et al., 2019)) 1413 to achieve state-of-the-art or comparable perfor-1414 mance across various time-series analysis tasks, 1415 including classification, forecasting (short/long-1416 term), imputation, anomaly detection, and few-1417 shot/zero-sample forecasting. The authors also 1418 found that self-attention functions similarly to PCA, 1419 providing a theoretical explanation for the versatil-1420 ity of transformers. 1421

## A.8 Evaluation Metrics for Task-Aware Tuning Stage

1422

1423

1435

1436

1443

1444

In our evaluation, we use the F1-macro score to 1424 assess the model's performance across datasets. 1425 F1-macro is particularly suitable for datasets with 1426 imbalanced label distributions, which is common 1427 in Human Activity Recognition (HAR) tasks where 1428 certain activities are overrepresented while others 1429 have fewer samples. Unlike the micro F1 score, 1430 which emphasizes the performance on frequent 1431 classes, F1-macro treats each class equally by cal-1432 culating the F1 score independently for each class 1433 and then averaging them. 1434

The formula for the F1-macro score is:

$$F1\text{-macro} = \frac{1}{C} \sum_{i=1}^{C} F1_i$$
(8)

1437where C is the total number of classes, and  $Fl_i$  is1438the F1 score for class i. The F1 score for each class1439is calculated as:

1440 
$$F1_i = \frac{2 \times \operatorname{Precision}_i \times \operatorname{Recall}_i}{\operatorname{Precision}_i + \operatorname{Recall}_i}$$
(9)

1441The precision and recall for each class are defined1442as:

$$\operatorname{Precision}_{i} = \frac{\operatorname{TP}_{i}}{\operatorname{TP}_{i} + \operatorname{FP}_{i}}$$
(10)

$$\operatorname{Recall}_{i} = \frac{\operatorname{TP}_{i}}{\operatorname{TP}_{i} + \operatorname{FN}_{i}}$$
(11)

1445where  $TP_i$ ,  $FP_i$ , and  $FN_i$  represent the number of1446true positives, false positives, and false negatives1447for class i, respectively. This metric ensures that1448performance is evaluated fairly across all classes,1449regardless of the frequency of each label, making1450it a robust measure for imbalanced datasets.

#### A.9 Sensor-Language Alignment Stage Output Examples

Tables 11 and 12 present two examples of the trend 1453 analysis results generated by SensorLLM and GPT-1454 40 based on the input sensor data. From the results, 1455 it is evident that SensorLLM outperforms GPT-40 1456 across both shorter and medium-length sequences. 1457 This demonstrates that our approach enables LLMs 1458 to better understand numerical variations, as well as 1459 accurately compute the time duration represented 1460 by the input sequences based on their length and 1461 the given sample rate. In contrast, current large 1462 language models struggle with directly interpreting 1463 numerical data, as their tokenization methods are 1464 not well-suited for tasks such as comparing numer-1465 ical values or counting (Yehudai et al., 2024). 1466

1451

1452

Sensor readings:	[-9.8237, -9.4551, -10.007, -11.273, -11.258, -11.677, -11.774, -11.638, -11.195, -11.087, -10.833, -11.044, -11.393, -11.943, -12.168, -15.455, -12.967, -12.326, -12.515, -13.195, -12.634, -11.873, -12.002, -11.583, -10.859, -10.349, -9.831, -9.1622, -8.2721, -6.9299, -6.255, -5.5998]
Ground-truth:	0.0 seconds to 0.02 seconds: growing 0.02 seconds to 0.06 seconds: declining 0.06 seconds to 0.08 seconds: growing 0.08 seconds to 0.12 seconds: declining 0.12 seconds to 0.2 seconds: growing 0.2 seconds to 0.3 seconds: declining 0.3 seconds to 0.34 seconds: growing 0.34 seconds to 0.38 seconds: declining 0.38 seconds to 0.42 seconds: growing 0.42 seconds to 0.44 seconds: growing 0.44 seconds to 0.62 seconds: growing
	Total growing trends: 6 Total declining trends: 5
	From 0.0s to 0.62s, normalized left-ankle y-axis accelerometer data is showcased in the sensor data. Examining the data, we notice 2 clear trend characteristics, with the trend fluctuating a total of eleven times. The analysis reveals that the data's declining inclination persisted for a total of 0.24 seconds, and a growing trend within a span of 0.38 seconds. The general trend observed is growing.
SensorLLM:	0.0 to 0.02 seconds: ascending 0.02 to 0.08 seconds: descending 0.08 to 0.1 seconds: ascending 0.1 to 0.12 seconds: descending 0.12 to 0.2 seconds: ascending 0.2 to 0.28 seconds: descending 0.32 to 0.36 seconds: descending 0.36 to 0.4 seconds: ascending 0.4 to 0.42 seconds: ascending 0.42 to 0.62 seconds: ascending
	Count of ascending segments: 6 Count of descending segments: 5
	The time series data encapsulates normalized left-ankle y-axis accelerome- ter sensor readings from 0.0 seconds to 0.62 seconds. Two separate trends and nine trend shifts are observed in the data. The analysis reveals that the data's descending inclination persisted for a total of 0.22 seconds, and an ascending trend for a sum of 0.40 seconds. The trend overall is ascending.
GPT-40 Evaluation:	3#Moderate errors in trend durations and shifts compared to ground truth.
GPT-4o's output:	0.0s to 0.02s: rising 0.02s to 0.06s: falling 0.06s to 0.08s: rising

	0.08s to 0.1s: falling
	0.1s to 0.12s: rising
	0.12s to 0.14s: falling
	0.14s to 0.3s: rising
	Total rising segments: 4
	Total falling segments: 3
	The normalized left-ankle y-axis accelerometer sensor readings recorded within the 0.0 to 0.3 second timeframe are presented in this sensor data. The input data displays three individual trends, with a comprehensive change count reaching 7. The examination reveals that the data's falling inclination endured for an aggregate of 0.08 seconds, succeeded by a rising trend for a cumulative period of 0.22 seconds, and a steady pattern for a total of 0.00 seconds. The dominant trend is rising.
GPT-40 Evaluation:	2#Significant discrepancies in trend durations and counts between model and ground-truth.

Table 11: Example 1 of SensorLLM and GPT-40 Generated Descriptions with Ground-Truth and GPT-40 Evaluation

Sensor readings:	[0.53137, 0.53137, 0.53137, 0.51176, 0.51176, 0.51176, 0.45098, 0.45098, 0.45098, 0.45098, 0.45098, 0.45082, 0.45882, 0.45882]
Ground-truth:	0.0 seconds to 0.04 seconds: stable 0.04 seconds to 0.06 seconds: decreasing 0.06 seconds to 0.1 seconds: stable 0.1 seconds to 0.12 seconds: decreasing 0.12 seconds to 0.18 seconds: stable 0.18 seconds to 0.2 seconds: increasing 0.2 seconds to 0.24 seconds: stable
	Number of stable trends: 4 Number of decreasing trends: 2 Number of increasing trends: 1 The sensor data represents readings taken from a normalized right-lower-
	arm x-axis gyroscope sensor between 0.0 and 0.24 seconds. Analysis reveals three separate trends within the data, undergoing a cumulative total of seven shifts in direction. Encapsulating the outcomes, the data's decreasing trend stretched across a total time of 0.04 seconds, came after an increasing pattern observed over 0.02 seconds, and a stable trend for 0.18 seconds in total. The dominant trend is decreasing.
SensorLLM:	0.0s to 0.04s: consistent 0.04s to 0.06s: downward 0.06s to 0.1s: consistent 0.1s to 0.12s: downward 0.12s to 0.18s: consistent 0.18s to 0.2s: upward 0.2s to 0.24s: consistent

	Number of consistent segments: 4 Number of downward segments: 2 Number of upward segments: 1	
	Number of upward segments: 1	
	The sensor data illustrates normalized right-lower-arm x-axis gyroscope sensor readings between 0.0 and 0.24 seconds. The input data displays three individual trends, with a comprehensive change count reaching 7. To encapsulate, the data's downward trend spanned a combined duration of 0.04 seconds, and then an upward pattern for a sum of 0.02 seconds, and a consistent trend for an accumulated time of 0.18 seconds. The overarching trend is characterized as downward.	
GPT-40 Evaluation:	5#The model's description matches the ground truth accurately.	
GPT-4o's output:	0.0s to 0.06s: steady	
	0.06s to 0.12s: falling	
	0.12s to 0.14s: steady	
	0.14s to 0.16s: rising	
	0.16s to 0.18s: steady	
	Total steady segments: 3	
	Total rising segments: 1	
	Total falling segments: 1	
	The normalized right-lower-arm x-axis gyroscope sensor readings recorded within the 0.0 to 0.18 second timeframe are presented in this sensor data. The input data displays three individual trends, with a comprehensive change count reaching 5. The examination reveals that the data's falling inclination endured for an aggregate of 0.06 seconds, succeeded by a rising trend for a cumulative period of 0.02 seconds, and a steady pattern for a total of 0.10 seconds. The dominant trend is steady.	
GPT-40 Evaluation:	2#Significant discrepancies in trend durations and counts compared to ground-truth.	

Table 12: Example 2 of SensorLLM and GPT-40 Generated Descriptions with Ground-Truth and GPT-40 Evaluation