

# IoT-LLM: ENHANCING REAL-WORLD IOT TASK REASONING WITH LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Large Language Models (LLMs) have demonstrated remarkable capabilities across textual and visual domains but often generate outputs that violate physical laws, revealing a gap in their understanding of the physical world. Inspired by human cognition—where perception is fundamental to reasoning—we explore augmenting LLMs with enhanced perception abilities using Internet of Things (IoT) sensor data and pertinent knowledge for IoT task reasoning in the physical world. In this work, we systematically study LLMs’ capability to address real-world IoT tasks by augmenting their perception and knowledge base, and then propose a unified framework, IoT-LLM, to enhance such capability. In IoT-LLM, we customize three steps for LLMs: preprocessing IoT data into formats amenable to LLMs, activating their commonsense knowledge through chain-of-thought prompting and specialized role definitions, and expanding their understanding via IoT-oriented retrieval-augmented generation based on in-context learning. To evaluate the performance, We design a new benchmark with five real-world IoT tasks with different data types and reasoning difficulties and provide the benchmarking results on six open-source and close-source LLMs. Experimental results demonstrate the limitations of existing LLMs with naive textual inputs that cannot perform these tasks effectively. We show that IoT-LLM significantly enhances the performance of IoT tasks reasoning of LLM, such as GPT-4, achieving an average improvement of 65% across various tasks against previous methods. The results also showcase LLMs’ ability to comprehend IoT data and the physical law behind data by providing a reasoning process. Limitations of our work are claimed to inspire future research in this new era.

## 1 INTRODUCTION

Recent advancements in large generative models have showcased their exceptional performance and versatility in handling complex tasks across textual and visual domains, as evidenced by the GPT series (Radford et al., 2018; 2019; Brown et al., 2020; Achiam et al., 2023; OpenAI, 2023) and visual generation models (Dosovitskiy et al., 2020; Liu et al., 2021; Ho et al., 2020; Peebles & Xie, 2023; Blattmann et al., 2023). However, these models could occasionally generate outputs that are physically implausible, often referred to as “hallucinations” (Alkaissi & McFarlane, 2023; Huang et al., 2023). Even advanced video generation models, e.g., Sora (Brooks et al., 2024), are susceptible to producing animations that contravene fundamental physical laws, such as a video clip containing a tipping water glass that appears to defy gravity. These observations suggest that generative models may not really comprehend and apply physical laws of the physical world as accurately as humans when acting as world simulators. This has renewed interest in research on the *World Model* that focuses on understanding and modeling the physical world in a brain-like manner (Dawid & LeCun, 2023; Garrido et al., 2024; Mendonca et al., 2023; Liu et al., 2024).

Unlike Large Language Models (LLMs) that map descriptions of the physical world to a latent space and perform reasoning by predicting the text sequence according to the probability, research on human cognitive science illustrates a different mechanism. The human brain comprises multiple mutually-functional areas, of which the important components include the temporal and occipital lobes for perception, and the frontal cortex for reasoning (Churchland & Sejnowski, 1988; Saxe et al., 2009; Hobeika et al., 2016; Grèzes et al., 2001). Notably, perception is the primary mechanism through which information about the physical world is acquired, and then effective reasoning

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

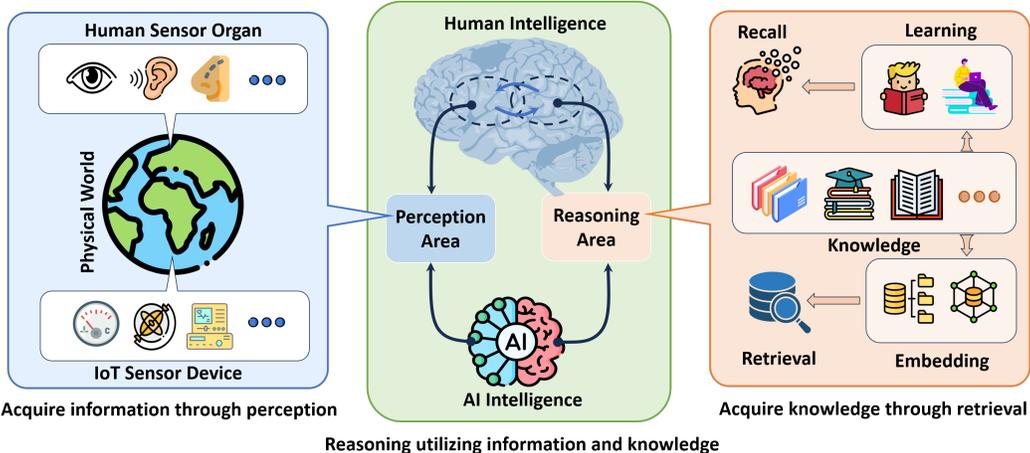


Figure 1: Inspired by human cognitive science, we augment LLMs with physical world perception from IoT data. Furthermore, by retrieving pertinent knowledge about IoT tasks, we enhance the reasoning capabilities of LLMs in executing real-world applications.

is inherently dependent on accurate perception. However, in LLMs, the physical world is only “perceived” through natural language, i.e., concepts and words in the semantic space, which denotes an indirect representation and abstraction of the physical world. A recent study in Nature shows language is primarily a tool for communication rather than thought (Fedorenko et al., 2024), so reasoning the physical-world problem with only language is limited. To enable LLMs with better reasoning capability in the real world, perception is highly demanded. Recent research on Vision Language Models (VLMs) builds the connection between visual perception and languages (Zhang et al., 2024a), yet the vision is only one of the various perceptual modalities. Many aspects of the physical world are still not perceived by existing LLMs.

We draw inspiration from how humans understand the physical world: perception to acquire information and reasoning with relevant domain knowledge. Firstly, humans perceive the world via a multitude of sensory organs, such as eyes for sight and ears for hearing. To empower machine perception, Internet of Things (IoT) sensors are developed. Since the first IoT sensor was designed for Coke machines to count the number of bottles in the 1980s (Madakam et al., 2015), IoT sensors become the “sensory organs” of machines, modeling the physical world for machine automation. Secondly, humans understand the world via the perception data with domain knowledge gained from experience and education. Similarly, LLMs can learn domain knowledge of both the physical world and sensors from the context to have stronger reasoning capabilities by in-context learning. In this manner, as shown in Fig. 1, we believe perception data with pertinent knowledge can enable LLMs to address complex problems with IoT-enabled perception in the real world. In this work, we aim to explore the following questions: (1) What types of real-world tasks can LLMs perform via the IoT perception of the physical world? (2) How can we enhance the LLM capability to deal with real-world tasks? (3) Do LLMs truly understand perception data and apply knowledge to realize real-world tasks?

Previous studies have primarily shown the viability of using LLMs for IoT task reasoning Xu et al. (2024b); Ji et al. (2024), but we find that these studies are not carefully scrutinized. (1) These studies only focus on specific tasks, such as R-peak identification and action recognition. The choices of tasks are not comprehensive, and thus they lack a benchmark to evaluate the performances of the methods. (2) They directly input raw IoT data into LLMs for reasoning, but LLMs are not good at dense numerical data and calculation (Zhou et al., 2024; Gruver et al., 2024). (3) They only evaluate their effectiveness on close-source LLMs, and lack a comprehensive study of benchmarking open-source LLMs with different parameter size.

To bridge this gap and answer the questions we proposed, we conduct an in-depth investigation of how to utilize LLMs to perform various tasks in the physical world using IoT data. Firstly, we explore whether LLMs can solve IoT classification and regression problems by setting a new bench-

108 mark with five classic IoT tasks with different data and levels of difficulties, including human activity  
109 recognition, industrial anomaly detection, heartbeat anomaly detection, WiFi-based human sensing,  
110 and indoor localization. The benchmark covers scenarios of daily life, industrial applications, and  
111 medical care, which will be detailed in the experiments. Secondly, we enhance LLMs’ reasoning  
112 capabilities with IoT data through three novel steps and consolidate three steps into IoT-LLM, a  
113 unified framework for IoT task reasoning. It is composed of three steps tailored for IoT reasoning:  
114 designing an LLM-friendly data format, activating knowledge by chain-of-thought prompting, and  
115 automatic IoT-oriented Retrieval-Augmented Generation (RAG) based on LLMs’ in-context learn-  
116 ing capability. Thirdly, to determine whether LLMs truly understand and then solve the task, we  
117 have LLMs generate analytical processes and analyze the reasonableness of the analytics. The anal-  
118 ysis generated by IoT-LLM indicates that LLMs can provide a reasonable process of solving simple  
119 tasks, but their efficacy diminishes in more specialized domains like heartbeat anomaly detection.  
120 This performance disparity is attributable to the complexity of data and limited domain-specific  
121 knowledge inherent in LLMs.

122 In summary, our contributions are as follows:

- 123 • We systematically study how Large Language Models (LLMs) can address real-world prob-  
124 lems by perceiving the physical world via IoT sensor data.
- 125 • We propose a unified framework to address IoT-related real-world problems, which en-  
126 hances the capability of LLMs through three steps: IoT data simplification and enrichment,  
127 IoT-oriented knowledge retrieval, and prompt configuration. To the best of our knowledge,  
128 this is the first unified framework for IoT-related tasks in the physical world.
- 129 • We establish the first benchmark for IoT task reasoning, including five real-world tasks with  
130 various types of IoT data. We benchmark both open-source and close-source LLMs with  
131 different parameter size. Empirical results show that our IoT-LLM significantly improves  
132 the performances of all base LLMs on IoT tasks.

## 135 2 RELATED WORK

136 **ML/DL methods in IoT tasks.** The Internet of Things (IoT) sensors gather diverse data from  
137 the real world, such as tri-axial acceleration, electrocardiogram readings, WiFi signals, and pres-  
138 sure (Sehrawat & Gill, 2019). These data have empowered various human sensing tasks, including  
139 Human Activity Recognition (HAR) (Lara & Labrador, 2013), health monitoring like heartbeat and  
140 respiration anomaly detection (Mousavi & Afghah, 2018; Aytakin et al., 2022), and industrial appli-  
141 cations such as machine operational state monitoring (Kong et al., 2023). Currently, these IoT data  
142 are primarily processed using traditional machine learning techniques, such as Support Vector Ma-  
143 chines (SVM) and K-Nearest Neighbors (KNN) Algorithm (Alam et al., 2016; Luo et al., 2021), or  
144 deep learning methods (Li et al., 2021; Njima et al., 2019). These approaches build black-box pre-  
145 dictors for specific tasks, yet each predictor only supports one task, and the task cannot be addressed  
146 with reasoning analysis, which motivates us to explore LLM for IoT tasks.

147 **LLMs in IoT tasks.** Existing literature on Large Language Models (LLMs) in IoT mainly regards  
148 LLM as a user interface or as coordinators in smart machines (Li et al., 2023; Cui et al., 2023;  
149 Du et al., 2023). However, in these studies, LLMs function as intermediaries and do not directly  
150 interpret IoT data to perform real-world tasks. Recent studies, such as Penetrative AI (Xu et al.,  
151 2024a) and HarGPT (Ji et al., 2024), have begun integrating IoT data into LLMs for specific tasks,  
152 leveraging their inherent knowledge bases. Despite these advancements, the exploration of LLMs  
153 processing IoT data remains nascent. Penetrative AI converts IoT data into textual and numerical  
154 formats for basic tasks like R-peak identification in ECG data, heavily relying on manually crafted  
155 expert knowledge, which limits automation and scalability. Similarly, HarGPT processes raw IMU  
156 data to recognize human activities using a chain of thought technique but is restricted to this specific  
157 data type and task, not demonstrating the broader applicability of LLMs. While these studies provide  
158 initial insights into using LLMs in the IoT domain, they do not offer a comprehensive framework  
159 that fully exploits LLM capabilities or systematically explores the interaction between LLMs and  
160 the physical world through IoT devices, which is the primary focus of our work.

## 3 METHODOLOGY

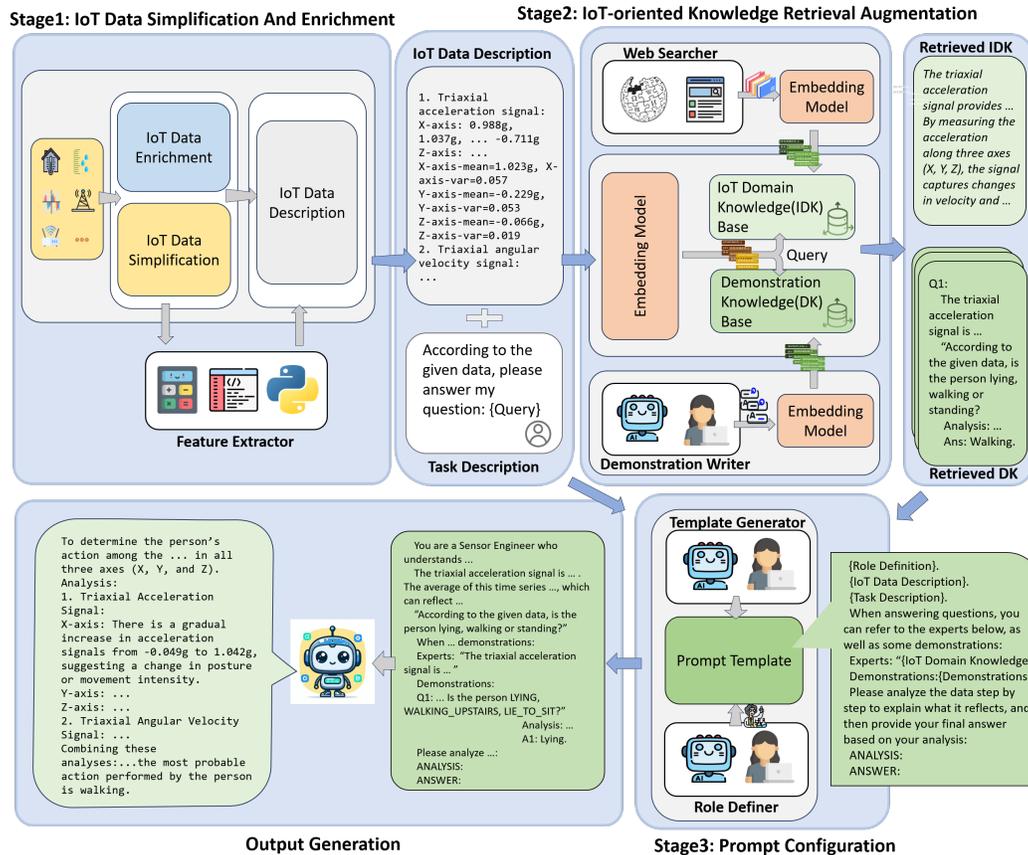


Figure 2: In our framework, IoT data is initially preprocessed to create a data description. Next, relevant IoT domain knowledge and task-specific demonstrations are retrieved. These elements are then combined into a prompt, which is input into a LLM to generate the final output.

In this section, we define the problem for IoT task reasoning with LLM and introduce our research methodology. The formulated research problem is how to leverage LLM and in-context learning for task reasoning for IoT data, termed as *IoT task reasoning*, e.g., using accelerators data for activity recognition or machine sensor for anomaly detection. The prompt for LLM should include two parts: data, as a way to perceive the physical world, and the task description, such as "Is it a Normal heartbeat (N) or Premature ventricular contraction beat (V)?" serves as the query. To evaluate the performance of IoT reasoning task, we build a new benchmark including 5 real-world tasks with different IoT data types and difficulty levels, encompassing both classification and regression problems.

At first, we employ LLMs to execute IoT tasks in a basic setting, similar to the existing approaches (Ji et al., 2024; Xu et al., 2024b), where the prompt provided to the LLMs includes only raw IoT data and the associated query. However, the performance of LLMs remains suboptimal. As shown by the baseline results in Table 2, even GPT-4 only achieves an accuracy of 43% for 3-way activity recognition and 50% for machine diagnosis based on their approach. These preliminary results akin to near-random guessing suggest a lack of comprehension of IoT data and tasks by this naive prompting way. Upon analyzing the characteristics of IoT data and real-world tasks, we identify that the challenges stem from the abstraction of dense numeric data and the lack of domain knowledge within LLMs. To address these challenges, we propose a unified framework (Fig. 2) consisting of three key stages: (1) IoT data simplification and enrichment, (2) IoT-oriented knowledge

216 augmentation, and (3) prompt configuration. Each stage addresses specific difficulties encountered  
217 by LLMs for IoT task reasoning, and we introduce each stage one by one.  
218

### 219 3.1 IOT DATA SIMPLIFICATION AND ENRICHMENT 220

221 Unlike textual human tasks that have been learned by LLMs, IoT data for IoT task reasoning presents  
222 unique challenges that hinder LLMs’ comprehension. Firstly, IoT data encompasses a diverse range  
223 of types and forms, many of which are complex time-series data (e.g., electrocardiogram read-  
224 ings) (Goldberger et al., 2000) or multi-variant data (e.g., WiFi CSI) (Yang et al., 2024). LLMs  
225 often struggle with accurately interpreting dense numerical data, especially when it involves long-  
226 sequence time-series data (Zhang et al., 2024b). Secondly, IoT data is typically composed of raw  
227 numerical values. This data often lacks essential textual annotations, such as units of measurement  
228 and metadata about the data collection process, which are critical for LLMs to interpret effectively  
229 in real-world applications. In summary, raw IoT data requires (1) appropriate simplification and  
230 (2) information enrichment. Previous studies (Xu et al., 2024b) have employed down-sampling  
231 techniques for time-series data but they only achieve coarse-grained simplification at a length level  
232 without enhancing the informational content of the IoT data. In contrast, we not only simplify IoT  
233 data at the token level but also enrich the IoT data by providing additional information to facilitate  
234 better understanding by LLMs (as illustrated in Fig. 10 in Appendix B). In this way, we transform  
235 complex raw IoT data into an LLM-friendly format for IoT task reasoning.

236 **IoT data simplification.** To achieve effective simplification, it is crucial to understand why LLMs  
237 struggle with dense numeric data. Firstly, according to recent research (Gruver et al., 2024; Spathis  
238 & Kawsar, 2023), tokenization methods, such as Byte Pair Encoding (BPE) often fragment numbers  
239 into tokens that do not align with their digits, resulting in inconsistent tokenization of floating-  
240 point numbers and complicating arithmetic operations. Therefore, in addition to down-sampling  
241 and keeping fixed precision (e.g., two digits of precision) to efficiently manage context length, we  
242 propose to insert spaces between digits to ensure distinct tokenization of each digit and use a comma  
243 (“;”) to separate each time step in a time series. Secondly, the complexity of long-sequence IoT data  
244 poses significant challenges for LLMs in analysis. To assist LLMs in processing this data, we extract  
245 essential statistical features, e.g., mean, variance, and FFT mean, utilizing external tools, such as  
246 Python scripts. We find that these fundamental features are strong enough for IoT task reasoning in  
247 classic IoT tasks. By doing so, we not only simplify IoT data at both length and token levels but also  
248 transform it into a format that is more suitable for tokenization and processing by LLMs.

249 **IoT data enrichment.** As previously noted, IoT data alone is insufficient for LLMs to effectively  
250 perform real-world tasks. To address this, we enrich the data by incorporating contextual information  
251 about the physical world. Specifically, we provide a comprehensive overview of IoT data collection  
252 and the integration of physical information. For instance, in human activity recognition (HAR) tasks  
253 where we employ inertial measurement unit (IMU) data including triaxial acceleration and angular  
254 velocity from accelerometers and gyroscopes, we meticulously outline the data collection process,  
255 incorporating the metadata such as sampling frequency (e.g., 10 Hz), device placement on the body,  
256 and units of measurement (e.g., gravitational acceleration and radians per second). This approach  
257 enables LLMs to not only align the three-axis IMU data with the corresponding three-dimensional  
258 spatial orientations of the human body but also to understand the physical significance of these  
259 numerical values, thereby enhancing the comprehension of LLMs for the task in the physical world.

### 260 3.2 IOT-ORIENTED KNOWLEDGE RETRIEVAL AUGMENTATION 261

262 In IoT task reasoning, the knowledge of LLMs to perform IoT tasks is significant. For example, de-  
263 tecting abnormal heartbeats from electrocardiogram (ECG) data requires interpreting ECG signals  
264 and associating them with specific heartbeat states (e.g., normal, premature ventricular contraction),  
265 necessitating specialized domain knowledge. Although previous research (Xu et al., 2024b) pro-  
266 poses to include specific expert knowledge for specific tasks, the augmentation is task-specific and  
267 added manually, which is time-consuming and not scalable. To address this, we enable LLMs with  
268 IoT knowledge in an automatic fashion. Inspired by the in-context learning capability of LLMs, we  
269 also retrieve task-specific demonstrations, such as question-answer pairs, to guide LLMs in effec-  
tively utilizing IoT data for analyzing IoT tasks.

We first construct an IoT domain knowledge base and a demonstration knowledge base, which will be utilized for retrieving domain knowledge and task-specific demonstrations. To ensure comprehensive coverage of knowledge about IoT data and tasks within the IoT domain knowledge base, we gather relevant documents (e.g., Wikipedia articles, research papers) through web searches encompassing the following themes: (1) IoT data domain knowledge, (2) IoT task domain knowledge, and (3) expert insights on leveraging IoT data for task execution. For the demonstration knowledge base, we create task-specific demonstrations (i.e., question-answer pairs) authored by human or AI models (e.g., ChatGPT). We then employ an embedding model (e.g., text-embedding-ada-002<sup>1</sup> by OpenAI) to embed texts into vectors and store the text chunks and corresponding embeddings as key-value pairs, which allows for efficient and scalable search capabilities. To improve the quality of retrieved contents, we also store metadata (e.g., IoT data type for IoT domain knowledge base and task type for demonstration knowledge base) alongside the vector embeddings within the vector database. This approach allows for advanced post-processing techniques, such as metadata filtering (Poliakov & Shvai, 2024), to refine search results and improve task-specific retrieval accuracy. Secondly, we retrieve relevant knowledge using both IoT data description and task description as query. We adopt a hybrid search method, which means utilizing both keyword-based retrievers and embedding-based retrievers to harness their unique strengths, ensuring the consistent retrieval of highly relevant and context-rich information. Finally, after applying a re-ranking technique to recalibrate the similarity between the query and retrieved texts using ranker models (e.g. bge-reranker-base<sup>2</sup>), we filter out the top-m most relevant pieces, thus obtaining pertinent knowledge, encompassing documents with specific domain knowledge and task demonstrations relevant to the task at hand.

### 3.3 PROMPT CONFIGURATION

In addition to augmenting LLMs’ knowledge by providing external documents in the context utilizing the in-context learning capability of LLMs, we further invoke LLMs’ internal knowledge by carefully configuring the prompt. Recent studies demonstrate that LLMs possess strong role-playing capabilities (Park et al., 2023). To leverage this, we assign specific roles to LLMs for particular tasks. For instance, we have LLMs assume the role of a professional doctor when performing heartbeat anomaly detection, thereby activating their internal domain knowledge. What’s more, since LLMs’ reasoning capability can be improved a lot by decomposing the whole problem into several parts (Wei et al., 2022), we decompose the reasoning procedure into two steps, prompting LLMs to analyze the IoT data and task first, and then provide the final answer based on this analysis. By doing so, we can also evaluate the extent to which the LLM understands IoT data and its capability to perform IoT tasks through the generated analysis. In the end, we employ a prompt template (refer to Fig. 8 in Appendix B) to structure the content discussed previously. The ultimate prompt is crafted based on the template and subsequently fed into a downstream LLM. The LLM then produces the final output, encompassing both analysis and answer to the specified task.

## 4 EXPERIMENTS

### 4.1 A BENCHMARK ON IOT TASK REASONING

#### 4.1.1 IOT TASKS.

To comprehensively assess the capability boundaries of LLMs for IoT task reasoning, we develop a new benchmark comprising five real-world tasks with diverse IoT data types and difficulty levels: (1) Human Activity Recognition (HAR) using Inertial Measurement Unit (IMU) data, (2) Industrial anomaly detection using metrics such as temperature, cooling power, and cooling efficiency, (3) Heartbeat anomaly detection using Electrocardiogram (ECG) data, (4) Human sensing using WiFi Channel State Information (CSI), and (5) Indoor localization based on WiFi signal strength. It is important to note that we don’t need to construct a knowledge base for each task especially, instead, we just need to construct two knowledge bases (i.e., one IoT domain knowledge base and one demonstration knowledge base), each of which contains all the domain/demonstration knowledge about the total five tasks. During the retrieval phase, we can easily fetch pertinent knowledge precisely corresponding to the task utilizing metadata (e.g., IoT data type and task type) stored within

<sup>1</sup><https://platform.openai.com/docs/guides/embeddings>

<sup>2</sup><https://huggingface.co/BAAI/bge-reranker-base>

the bases. For demonstrations, we utilize the one-shot setting, which means we retrieve one example for each category in classification tasks.

Table 1: **Performance of LLMs on WiFi-based Indoor Localization task.** Since this is a regression task, we choose the Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and standard deviation (STD) of the RMSE as the main performance metrics.

Method		Model					
		Llama2-7B	Mistral-7B	Claude-3.5	Gemini-pro	GPT-3.5	GPT-4
Base-line	RMSE (m)	0.374	11.570	0.829	2.318	2.598	0.741
	MAE (m)	0.313	9.347	0.696	1.814	1.937	0.581
	STD	0.903	6.856	1.607	5.999	6.715	1.502
Ours	RMSE (m)	0.355	9.995	0.404	<b>0.313</b>	0.719	0.402
	MAE (m)	0.295	7.980	0.341	<b>0.265</b>	0.592	0.341
	STD	0.852	11.146	<b>0.706</b>	0.763	1.765	0.697
Improvement	RMSE (m)	<b>+5.1%</b>	<b>+13.6%</b>	<b>+51.3%</b>	<b>+86.5%</b>	<b>+72.3%</b>	<b>+45.7%</b>
	MAE (m)	<b>+5.8%</b>	<b>+14.6%</b>	<b>+51.0%</b>	<b>+85.4%</b>	<b>+69.4%</b>	<b>+41.3%</b>

Table 2: **Overall performance of LLMs on IoT tasks.** **HAR-2cls** stands for classifying walking and standing activities. **HAR-3cls** stands for classifying lying, walking upstairs, and transitioning from lying to sitting activities. **Heartbeat** stands for classifying normal and abnormal heartbeats. **Machine** stands for determining whether the coolers work properly or not. **Occupancy** stands for detecting the presence of a person in a room.

Model		IoT tasks (Accuracy $\uparrow$ )				
		HAR-2cls	HAR-3cls	Heartbeat	Machine	Occupancy
Llama2-7B	Baseline	50.0%	32.8%	50.0%	35.0%	48.4%
	Ours	57.2%	38.0%	54.5%	56.4%	82.5%
	Improvement	<b>+14.4%</b>	<b>+15.9%</b>	<b>+9.0%</b>	<b>+61.1%</b>	<b>+70.5%</b>
Mistral-7B	Baseline	61.5%	26.0%	44.0%	31.5%	50.0%
	Ours	84.9%	42.7%	60.5%	<u>92.1%</u>	61.1%
	Improvement	<b>+38.0%</b>	<b>+64.2%</b>	<b>+37.5%</b>	<b>+192.4%</b>	<b>+22.2%</b>
Claude-3.5	Baseline	98.3%	80.1%	52.4%	51.0%	50.0%
	Ours	<b>100.0%</b>	<b>95.3%</b>	<b>81.0%</b>	86.3%	82.5%
	Improvement	<b>+1.7%</b>	<b>+19.0%</b>	<b>+54.6%</b>	<b>+69.2%</b>	<b>+65.0%</b>
Gemini-pro	Baseline	39.3%	34.0%	52.0%	49.0%	55.9%
	Ours	88.4%	<u>82.8%</u>	51.5%	70.1%	66.2%
	Improvement	<b>+124.9%</b>	<b>+143.5%</b>	<b>-1.0%</b>	<b>+43.1%</b>	<b>+18.4%</b>
GPT-3.5	Baseline	91.5%	33.3%	35.3%	51.5%	50.0%
	Ours	92.1%	45.8%	51.0%	61.5%	<b>92.1%</b>
	Improvement	<b>+0.7%</b>	<b>+37.5%</b>	<b>+44.5%</b>	<b>+19.4%</b>	<b>+84.2%</b>
GPT-4	Baseline	77.3%	43.3%	54.0%	49.5%	43.7%
	Ours	<b>100.0%</b>	87.8%	<u>69.8%</u>	<b>92.4%</b>	<u>86.6%</u>
	Improvement	<b>+29.4%</b>	<b>+102.8%</b>	<b>+29.3%</b>	<b>+86.7%</b>	<b>+98.2%</b>

#### 4.1.2 IOT DATASETS.

In our benchmark, we choose public IoT datasets on the five tasks to ensure fairness. Since some datasets are too challenging for LLMs with many classes, we simplify some datasets by only using a subset, which is also employed in previous works (Ji et al., 2024).

**Human Activity Recognition.** We employ the Smartphone-Based Recognition of Human Activities and Postural Transitions Dataset (Reyes-Ortiz et al., 2015). This dataset comprises raw IMU data, specifically 3-axial linear acceleration, and 3-axial angular velocity, captured at a sampling rate of 50Hz by the smartphone’s accelerometer and gyroscope. The data encompasses twelve distinct activities. To reduce both the sequence length and data complexity, we down-sample the data to 10Hz. Given the challenges associated with multi-class classification for LLM, instead of utilizing all twelve activity categories, we conduct a binary classification task involving the WALKING and STANDING labels, and a ternary classification task with the LYING, WALKING UPSTAIRS, and LIE TO SIT labels.

**Industrial anomaly detection.** We employ the Condition Monitoring of Hydraulic Systems Dataset (Helwig et al., 2018), which facilitates the assessment of a hydraulic test rig’s condition using multi-sensor data, including temperature, cooling power, and efficiency factor series, all experimentally derived from the rig. The dataset categorizes cooler conditions into three severity grades: (1) close to failure; (2) reduced efficiency; and (3) full efficiency. For simplicity, we focus on a binary classification task using only “close to failure” and “full efficiency” categories.

**Heartbeat anomaly detection.** We employ the MIT-BIH Arrhythmia Database (Goldberger et al., 2000). This dataset comprises ECG recordings from 48 subjects, each sampled at 360Hz, and categorizes heartbeats into several types, including Normal beat (N), Atrial premature beat (A), and Premature ventricular contraction (V), among others. To reduce the difficulty of the task, we down-sample the signals to 72Hz and focus on a binary classification task using only the Normal beat (N) and Premature ventricular contraction (V) categories.

**Human sensing task.** We utilize a dataset collected using a TP-Link TL-WDR4300 WiFi router operating at 5 GHz with a 40 MHz bandwidth (Zhuravchak et al., 2022). The dataset specifically captures the absence of human presence across three different rooms. Each room’s environment is carefully monitored to record Channel State Information (CSI) that reflects the presence or absence of occupants, providing a robust basis for occupancy detection tasks.

**Indoor localization task.** We utilize a dataset collected in a laboratory environment using an IoT system developed in (Huang et al., 2022). The dataset consists of RSSI signals, the basis for determining human positions within the space. By collecting RSS fingerprints at various reference points, a signal radio map is constructed using a modified Gaussian Process Regression (GPR) method. This approach allows us to estimate the RSS distribution at any given location, providing a reliable means of localizing human presence in the environment.

#### 4.1.3 LLM BASELINES.

In the conducted experiments, we utilize a combination of proprietary and open-source LLMs, including gpt-3.5-turbo, gpt-4-turbo, claude-3-5-sonnet, gemini-pro, Mistral-7B<sup>3</sup>, and LLama2-7B<sup>4</sup>. This diverse selection of models enables a comprehensive evaluation of the LLMs’ capabilities in executing IoT tasks and provides insights into their respective strengths and limitations in real-world applications. The code implementations of IoT-LLM have been attached in the supplementary materials and will be made public after publication.

## 4.2 RESULTS AND ANALYSIS

To evaluate the efficacy of our proposed framework in enhancing the capabilities of IoT task reasoning for LLMs, we use HarGPT (Ji et al., 2024) as the baseline, of which the prompts only contain

<sup>3</sup><https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

<sup>4</sup><https://huggingface.co/togethercomputer/LLaMA-2-7B-32K>

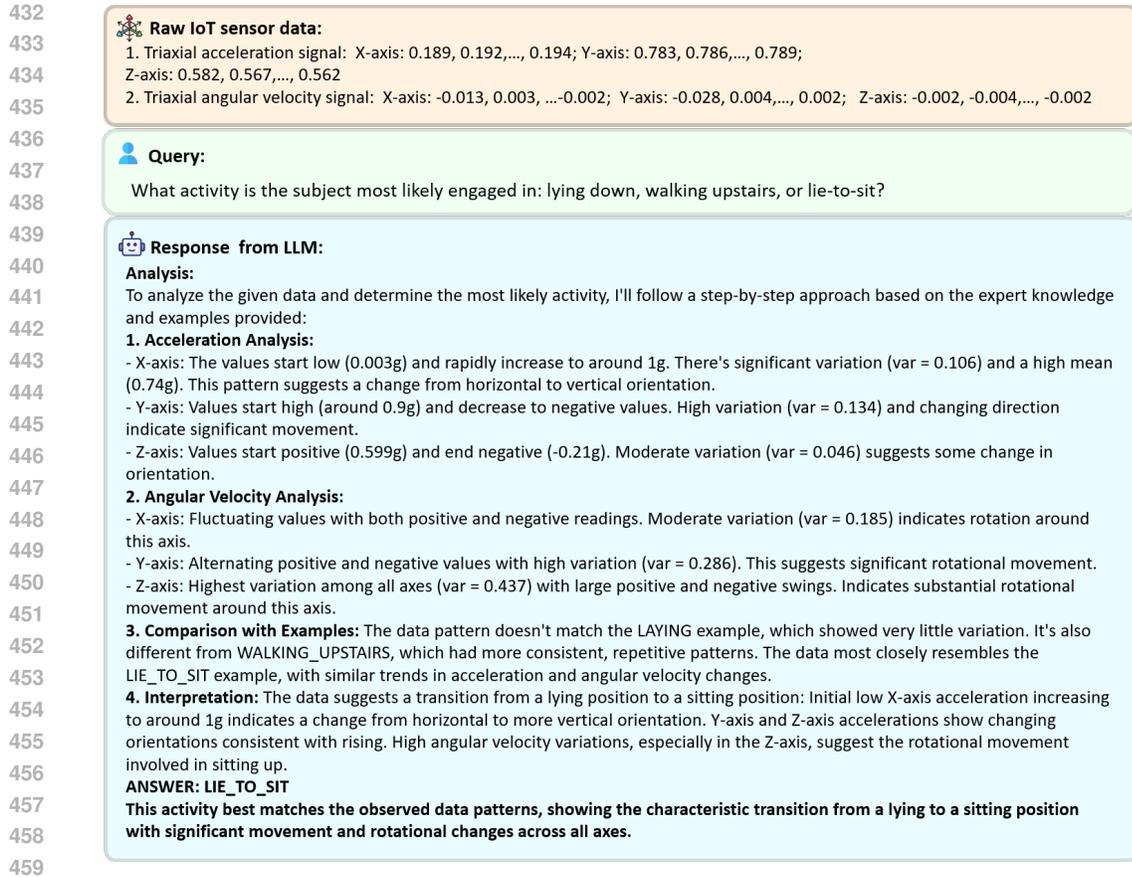


Figure 3: Response example of LLM for human activity recognition (HAR). For more comprehensive examples, please refer to Appendix A.

raw IoT data and corresponding task descriptions, without any data preprocessing, domain knowledge, and demonstrations. The overall performance of LLMs on IoT tasks is shown in Table 1 and Table 2. The results show that our proposed method consistently boosts the performance of all the LLMs to complete IoT tasks in real-world scenarios. Notably, advanced LLMs such as Claude-3.5, Gemini-pro, and GPT-4 have demonstrated significant performance improvements, evolving from near-random guessing to effectively solving certain tasks. After analyzing the overall performance and outputs of LLMs in IoT task reasoning, we can answer the questions we proposed in the introduction now. Here is a summary of our arguments regarding the IoT task reasoning with LLMs.

**LLMs excel in various IoT tasks but struggle with complex data challenge.** Based on the experimental results, we observe that when provided with perception data (i.e., IoT data collected by sensors) and external knowledge, advanced LLMs like GPT-4 and Claude-3.5 can effectively perform various IoT tasks in the physical world, particularly excelling in HAR using IMU data. However, LLMs' performance is limited by their intrinsic lack of domain-specific knowledge and difficulty in comprehending numerical data. For instance, in the task of heartbeat anomaly detection, even provided with external knowledge, LLMs perform sub-optimally. This is because the time-series nature of ECG data presents significant challenges for LLMs due to its numerical density and length. Although we have mitigated some of these challenges by simplifying the data, this approach only addresses the issue at the data level without fundamentally resolving it at the model level. Additionally, LLMs inherently lack the extensive medical knowledge required for comprehensive analysis. While retrieved knowledge can suffice for simpler tasks, more complex problems may necessitate further model fine-tuning to incorporate deeper and broader medical expertise.

**LLMs are excellent learners in IoT task reasoning.** Without domain-specific knowledge and relevant demonstrations, LLMs face significant challenges in performing IoT tasks, often resorting to near-random guessing, especially in tasks such as heartbeat anomaly detection. This indicates that real-world tasks remain challenging for LLMs to execute directly. However, LLMs are excellent learners, and their capabilities can be significantly enhanced through data simplification & enrichment and knowledge retrieval augmentation. Specifically, the LLama2-7B, Mistral-7B, Claude-3.5, Gemini-pro, GPT-3.5, and GPT-4 models exhibit average performance improvements of 30%, 62%, 44%, 69%, 43%, and 65% respectively across various tasks, underscoring the effectiveness of our methodology.

**LLMs can act as experts, not just classifiers or predictors.** In our study, we prompt LLMs to generate both an analysis of the task and the final answer. Based on this analysis, we demonstrate that LLMs can fully comprehend preprocessed IoT data and effectively utilize the provided knowledge to perform IoT tasks. Unlike traditional DL/ML methods, which are trained end-to-end to produce only the final answer, LLMs offer more explainable results. Specifically, LLMs not only provide the final answer but also the reasoning behind it, akin to expert suggestions in daily life. For instance, when tasked with human activity recognition (as illustrated in Fig.3), the LLM delivers a detailed step-by-step analysis before presenting the final answer.

### 4.3 ABLATION STUDY

To evaluate the impact of different components within our framework, we performed an ablation study using GPT-4 on HAR and industrial anomaly detection tasks. We tested the following configurations: (1) IoT data simplification and enrichment, (2) addition of retrieved domain knowledge based on (1), (3) inclusion of retrieved demonstrations based on (2), and (4) the full configuration, which incorporates role descriptions and chain-of-thought techniques as outlined in the Prompt Configuration stage. The results, presented in Table 3, reveal that for straightforward tasks such as classifying walking and standing activities, IoT data simplification and enrichment and domain knowledge retrieval are sufficient. However, for more complex tasks, the inclusion of additional modules significantly boosts performance. Overall, our findings indicate that each module in our framework progressively enhances the ability of LLMs to perform IoT-related tasks using IoT data.

Table 3: Ablation study of different modules within our framework on three tasks.

Method	IoT tasks (Accuracy↑)		
	HAR-2cls	HAR-3cls	Machine
Baseline	77.3%	43.3%	49.5%
+ IoT data simplification and enrichment	96.0%	47.3%	62.7%
+ retrieved domain knowledge	100.0%	78.7%	78.0%
+ retrieved demonstrations	100.0%	86.7%	83.3%
Full setting	<b>100.0%</b>	<b>87.8%</b>	<b>92.4%</b>

## 5 CONCLUSION

LLMs often struggle with tasks requiring an understanding of physical laws. To address this, we propose IoT-LLM, a framework that integrates IoT sensor data with LLMs to enhance their perception and reasoning abilities in the physical world. Evaluated on tasks like human activity recognition and industrial anomaly detection, IoT-LLM improves LLM performance by approximately 65%, though challenges remain in specialized domains. This approach systematically enhances LLM capabilities for real-world applications by leveraging IoT data.

**Limitations.** While LLMs can manage low-dimensional time-series data, they face significant challenges with higher-dimensional data, such as audio and 3D point cloud data, due to their extensive length and complexity. Integrating such data into the context of LLMs is both difficult and impractical. Instead, directly fine-tuning LLMs with data specific to these modalities may be a more effective approach. This strategy could be explored in future research to extend the capabilities of LLMs to handle more complex data types within the IoT domain.

## REFERENCES

- 540  
541  
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Furqan Alam, Rashid Mehmood, Iyad Katib, and Aiiad Albeshri. Analysis of eight data mining  
546 algorithms for smarter internet of things (iot). *Procedia Computer Science*, 98:437–442, 2016.
- 547 Hussam Alkaissi and Samy I McFarlane. Artificial hallucinations in chatgpt: implications in scien-  
548 tific writing. *Cureus*, 15(2), 2023.
- 550 Idil Aytekin, Onat Dalmaz, Kaan Gonc, Haydar Ankishan, Emine Ulku Saritas, Ulas Bagci, Haydar  
551 Celik, and Tolga Çukur. Covid-19 detection from respiratory sounds with hierarchical spectro-  
552 gram transformers. *IEEE Journal of Biomedical and Health Informatics*, 28:1273–1284, 2022.  
553 URL <https://api.semanticscholar.org/CorpusID:250699108>.
- 554 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
555 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
556 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 558 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
559 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
560 generation models as world simulators. 2024. URL [https://openai.com/research/  
561 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 562 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
563 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
564 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 565 Patricia S. Churchland and Terrence J. Sejnowski. Perspectives on cognitive neuroscience. *Sci-*  
566 *ence*, 242 4879:741–5, 1988. URL [https://api.semanticscholar.org/CorpusID:  
567 11737656](https://api.semanticscholar.org/CorpusID:11737656).
- 568 Hongwei Cui, Yuyang Du, Qun Yang, Yulin Shao, and Soung Chang Liew. Llmind: Orchestrating  
569 ai and iot with llms for complex task execution. *arXiv preprint arXiv:2312.09007*, 2023.
- 571 Anna Dawid and Yann LeCun. Introduction to latent variable energy-based models: A path towards  
572 autonomous machine intelligence. *arXiv preprint arXiv:2306.02572*, 2023.
- 573 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
574 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-  
575 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition  
576 at scale. *ArXiv*, abs/2010.11929, 2020. URL [https://api.semanticscholar.org/  
577 CorpusID:225039882](https://api.semanticscholar.org/CorpusID:225039882).
- 578 Jiachen Du, Boyang Jia, and Xinyi Fu. Space brain: An ai autonomous spatial decision system. In  
579 *CAAI International Conference on Artificial Intelligence*, pp. 61–67. Springer, 2023.
- 581 Evelina Fedorenko, Steven T Piantadosi, and Edward AF Gibson. Language is primarily a tool for  
582 communication rather than thought. *Nature*, 630(8017):575–586, 2024.
- 584 Quentin Garrido, Mahmoud Assran, Nicolas Ballas, Adrien Bardes, Laurent Najman, and Yann  
585 LeCun. Learning and leveraging world models in visual representation learning. *arXiv preprint  
586 arXiv:2403.00504*, 2024.
- 587 Ary Goldberger, Luís Amaral, L. Glass, Shlomo Havlin, J. Hausdorg, Plamen Ivanov, R. Mark,  
588 J. Mietus, G. Moody, Chung-Kang Peng, H. Stanley, and Physiokit Physiobank. Components  
589 of a new research resource for complex physiologic signals. *PhysioNet*, 101, 01 2000.
- 590 Julie Grèzes, Pierre Fonlupt, Bennett I. Bertenthal, Chantal Delon-Martin, Christoph Segebarth,  
591 and Jean Decety. Does perception of biological motion rely on specific brain regions? *Neu-  
592 roImage*, 13:775–785, 2001. URL [https://api.semanticscholar.org/CorpusID:  
593 2899041](https://api.semanticscholar.org/CorpusID:2899041).

- 594 Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. Large language models are zero-shot  
595 time series forecasters. *Advances in Neural Information Processing Systems*, 36, 2024.  
596
- 597 Nikolai Helwig, Eliseo Pignanelli, and Andreas Schtze. Condition monitoring of hydraulic systems.  
598 UCI Machine Learning Repository, 2018. DOI: <https://doi.org/10.24432/C5CW21>.
- 599 Jonathan Ho, Ajay Jain, and P. Abbeel. Denoising diffusion probabilistic models. *ArXiv*,  
600 abs/2006.11239, 2020. URL [https://api.semanticscholar.org/CorpusID:  
601 219955663](https://api.semanticscholar.org/CorpusID:219955663).  
602
- 603 Lucie Hobeika, C. Diard-Detoef, Béatrice Garcin, Richard Levy, and Emmanuelle Volle. General  
604 and specialized brain correlates for analogical reasoning: A meta-analysis of functional imaging  
605 studies. *Human Brain Mapping*, 37, 2016. URL [https://api.semanticscholar.org/  
606 CorpusID:205849677](https://api.semanticscholar.org/CorpusID:205849677).
- 607 He Huang, Jianfei Yang, Xu Fang, Hao Jiang, and Lihua Xie. Varifi: Variational inference for indoor  
608 pedestrian localization and tracking using imu and wifi rss. *IEEE Internet of Things Journal*, 10  
609 (10):9049–9061, 2022.  
610
- 611 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong  
612 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language  
613 models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*,  
614 2023.
- 615 Sijie Ji, Xinzhe Zheng, and Chenshu Wu. Hargpt: Are llms zero-shot human activity recognizers?  
616 *arXiv preprint arXiv:2403.02727*, 2024.  
617
- 618 Fanhui Kong, Jianqiang Li, Bin Jiang, Huihui Wang, and Houbing Herbert Song. Integrated  
619 generative model for industrial anomaly detection via bidirectional lstm and attention mecha-  
620 nism. *IEEE Transactions on Industrial Informatics*, 19:541–550, 2023. URL [https://api.  
621 semanticscholar.org/CorpusID:236732804](https://api.semanticscholar.org/CorpusID:236732804).
- 622 Oscar D. Lara and Miguel A. Labrador. A survey on human activity recognition using wearable  
623 sensors. *IEEE Communications Surveys & Tutorials*, 15:1192–1209, 2013. URL [https://  
624 api.semanticscholar.org/CorpusID:206583924](https://api.semanticscholar.org/CorpusID:206583924).
- 625 Bing Li, Wei Cui, Wen Wang, Le Zhang, Zhenghua Chen, and Min Wu. Two-stream convolution  
626 augmented transformer for human activity recognition. In *AAAI Conference on Artificial Intelli-  
627 gence*, 2021. URL <https://api.semanticscholar.org/CorpusID:235192773>.
- 628 Fu Li, Jiaming Huang, Yi Gao, and Wei Dong. Chatiot: Zero-code generation of trigger-action based  
629 iot programs with chatgpt. In *Proceedings of the 7th Asia-Pacific Workshop on Networking*, pp.  
630 219–220, 2023.
- 631 Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and  
632 language with ringattention. *arXiv preprint arXiv:2402.08268*, 2024.  
633
- 634 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Bain-  
635 ing Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *2021  
636 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9992–10002, 2021. URL  
637 <https://api.semanticscholar.org/CorpusID:232352874>.
- 638 Minhui Luo, Jin Zheng, Wei Sun, and Xing Zhang. Wifi-based indoor localization using clustering  
639 and fusion fingerprint. *2021 40th Chinese Control Conference (CCC)*, pp. 3480–3485, 2021. URL  
640 <https://api.semanticscholar.org/CorpusID:238413227>.
- 641 Somayya Madakam, Ramya Ramaswamy, and Siddharth Tripathi. Internet of things (iot): A litera-  
642 ture review. *Journal of Computer and Communications*, 3(5):164–173, 2015.  
643
- 644 Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured world models from human videos.  
645 *arXiv preprint arXiv:2308.10901*, 2023.  
646  
647

- 648 Sajad Mousavi and Fatemeh Afghah. Inter- and intra- patient ecg heartbeat classification for ar-  
649 rhythmia detection: A sequence to sequence deep learning approach. *ICASSP 2019 - 2019 IEEE*  
650 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1308–1312,  
651 2018. URL <https://api.semanticscholar.org/CorpusID:56173441>.
- 652 Wafa Njima, Iness Ahriz, Rafik Zayani, Michel Terre, and Ridha Bouallegue. Deep cnn for indoor  
653 localization in iot-sensor systems. *Sensors*, 19(14):3127, 2019.
- 654 R OpenAI. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2(5), 2023.
- 655 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and  
656 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings*  
657 *of the 36th annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
- 658 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
659 *the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 660 Mykhailo Poliakov and Nadiya Shvai. Multi-meta-rag: Improving rag for multi-hop queries using  
661 database filtering with llm-extracted metadata. *ArXiv*, abs/2406.13213, 2024. URL <https://api.semanticscholar.org/CorpusID:270620041>.
- 662 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language under-  
663 standing by generative pre-training. 2018.
- 664 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
665 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 666 Jorge Reyes-Ortiz, Davide Anguita, Luca Oneto, and Xavier Parra. Smartphone-based recognition  
667 of human activities and postural transitions. UCI Machine Learning Repository, 2015. DOI:  
668 <https://doi.org/10.24432/C54G7M>.
- 669 Rebecca Saxe, Susan L. Whitfield-Gabrieli, Jonathan Scholz, and Kevin A. Pelphrey. Brain re-  
670 gions for perceiving and reasoning about other people in school-aged children. *Child develop-*  
671 *ment*, 80 4:1197–209, 2009. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:13911130)  
672 [13911130](https://api.semanticscholar.org/CorpusID:13911130).
- 673 Deepti Sehrawat and Nasib Singh Gill. Smart sensors: Analysis of different types of iot sensors. In  
674 *2019 3rd International Conference on Trends in Electronics and Informatics (ICOEI)*, pp. 523–  
675 528. IEEE, 2019.
- 676 Dimitris Spathis and Fahim Kawsar. The first step is the hardest: Pitfalls of representing and tok-  
677 enizing temporal data for large language models. *Journal of the American Medical Informatics*  
678 *Association : JAMIA*, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:261697232)  
679 [261697232](https://api.semanticscholar.org/CorpusID:261697232).
- 680 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc  
681 Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language mod-  
682 els. *ArXiv*, abs/2201.11903, 2022. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:246411621)  
683 [CorpusID:246411621](https://api.semanticscholar.org/CorpusID:246411621).
- 684 Huatao Xu, Liying Han, Qirui Yang, Mo Li, and Mani Srivastava. Penetrative AI: Making LLMs  
685 comprehend the physical world. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.),  
686 *Findings of the Association for Computational Linguistics ACL 2024*, pp. 7324–7341, Bangkok,  
687 Thailand and virtual meeting, August 2024a. Association for Computational Linguistics. URL  
688 <https://aclanthology.org/2024.findings-acl.437>.
- 689 Huatao Xu, Liying Han, Qirui Yang, Mo Li, and Mani Srivastava. Penetrative ai: Making llms  
690 comprehend the physical world. In *Proceedings of the 25th International Workshop on Mobile*  
691 *Computing Systems and Applications*, pp. 1–7, 2024b.
- 692 Jianfei Yang, He Huang, Yunjiao Zhou, Xinyan Chen, Yuecong Xu, Shenghai Yuan, Han Zou,  
693 Chris Xiaoxuan Lu, and Lihua Xie. Mm-fi: Multi-modal non-intrusive 4d human dataset for  
694 versatile wireless sensing. *Advances in Neural Information Processing Systems*, 36, 2024.

702 Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks:  
703 A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024a.  
704

705 Xiyuan Zhang, Ranak Roy Chowdhury, Rajesh K. Gupta, and Jingbo Shang. Large language  
706 models for time series: A survey. *ArXiv*, abs/2402.01801, 2024b. URL [https://api.  
707 semanticscholar.org/CorpusID:267411923](https://api.semanticscholar.org/CorpusID:267411923).

708 Lexin Zhou, Wout Schellaert, Fernando Martínez-Plumed, Yael Moros-Daval, Cèsar Ferri, and José  
709 Hernández-Orallo. Larger and more instructable language models become less reliable. *Nature*,  
710 pp. 1–8, 2024.

711 Andrii Zhuravchak, Oleg Kapshii, and Evangelos Pournaras. Human activity recognition based on  
712 wi-fi csi data-a deep neural network approach. *Procedia Computer Science*, 198:59–66, 2022.  
713  
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 RESPONSE EXAMPLES OF LLMs FOR IOT TASKS

In this section, we present examples of responses generated by LLM for various applications, including industrial anomaly detection, heartbeat anomaly detection, WiFi-based human sensing, and indoor localization. To emphasize the output of the LLMs, we do not provide detailed prompts; instead, we simply display the raw IoT sensor data and the corresponding user queries.

 **Raw IoT sensor data:**

1. Temperature Change Sequence: 57.277, ..., 57.301
2. Cooling Power Change Sequence: 1.488, 1.497, ..., 1.495
3. Cooling Efficiency Change Sequence: 18.704, ..., 18.821

 **Query**

Is the machine's cooling system functioning properly?

 **Response from LLM:**

To analyze the operational status of the machine's cooling system, we'll examine the Temperature Change Sequence, Cooling Power Change Sequence, and Cooling Efficiency Change Sequence step by step.

**1. Temperature Change Sequence Analysis:**  
The sequence starts at 55.836°C and fluctuates throughout the 60-second period, with a slight increasing trend towards the end, closing at 55.75°C. The temperature shows a range of fluctuations but remains within a relatively narrow band (~55.5°C to ~56.0°C). This indicates that while there is a slight increase, the temperature does not significantly spike or drop, suggesting a potentially stable but slightly elevated operational temperature.

**2. Cooling Power Change Sequence Analysis:**  
The cooling power starts at 1.568KW and shows variability throughout the sequence, peaking at 1.636KW and then decreasing towards the end, closing at 1.568KW. The cooling power's fluctuations are relatively minor, indicating that the cooling system is adjusting its output in response to the system's needs. However, the fact that it increases and then returns to the starting value suggests that the cooling system is actively working to manage the temperature, albeit with varying degrees of success.

**3. Cooling Efficiency Change Sequence Analysis:**  
The cooling efficiency begins at 20.09% and exhibits some variability, peaking at 20.908% and then showing a decrease in efficiency towards the end, closing at 20.125%. The efficiency values indicate that the cooling system's effectiveness in removing heat from the system varies but reaches relatively high efficiency levels during the sequence. The decrease towards the end suggests a reduction in cooling performance over time.

**Final Analysis and Conclusion:**  
The data reflects a cooling system that is actively working to manage the machine's temperature, as evidenced by the fluctuations in cooling power and efficiency. However, several indicators suggest that the cooling system may not be functioning optimally:

- The temperature, while relatively stable, is on the higher side throughout the sequence and shows a slight increasing trend. This suggests that the cooling system might be struggling to maintain or lower the temperature effectively.
- The cooling power shows that the system is adjusting its output, but the return to the initial value suggests that these adjustments may not be fully effective over time.
- The cooling efficiency, despite reaching relatively high levels, decreases towards the end of the sequence, indicating a potential drop in the cooling system's performance.

Based on this analysis, while the cooling system is operational and making adjustments in response to the machine's cooling needs, the slight increase in temperature, along with the variability in cooling power and the drop in efficiency towards the end, suggests that the cooling system may not be functioning at its best. It appears to be underperforming, possibly due to issues like partial blockages, reduced coolant flow, or other inefficiencies within the system.

**ANSWER:** The machine's cooling system is not functioning properly and may require maintenance or further inspection to address potential issues affecting its performance.

Figure 4: Response example of LLM for industrial anomaly detection.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

 **Raw IoT sensor data:**

-0.205, -0.195, -0.155, -0.05, -0.015, 0.02, 0.125, 0.21, , ..., 0.055

 **Query:**

Is it a Normal heartbeat(N) or Premature ventricular contraction beat(V)?

 **Response from LLM:**

To analyze the given ECG data, we will follow a systematic approach, examining the amplitude, intervals, and progression of the electrical signal that represents a heartbeat. This will involve looking for patterns that match either a normal heartbeat or a premature ventricular contraction (PVC).

**Step 1: Baseline and P-wave Identification**

- The initial segment shows a gradual increase in amplitude from -0.22mV to 0.125mV, which could represent the P-wave and the beginning of the QRS complex. In a normal heartbeat, the P-wave represents atrial depolarization. The gradual increase is consistent with the expected pattern, although the amplitude is somewhat low for a typical P-wave, which might not be clearly defined in this case.

**Step 2: QRS Complex Analysis**

- The amplitude sharply rises to 1.15mV, which is a characteristic feature of the QRS complex, representing ventricular depolarization. This sharp rise is expected in both normal and PVC beats. However, the amplitude of 1.15mV is within a normal range for the QRS complex, not indicating an overtly abnormal morphology based on amplitude alone.

**Step 3: Examination of the ST Segment and T-wave**

- After the peak, there's a dramatic drop to -1.395mV, which is highly unusual for a normal T-wave following the QRS complex. In a normal ECG, the T-wave represents ventricular repolarization and should be a positive deflection following the QRS complex, not a sharp negative drop. This significant negative deflection is not characteristic of a normal T-wave and indicates a significant abnormality.

**Step 4: Checking for Premature Beat**

- The dramatic changes in amplitude, especially the sharp drop after the peak, could suggest a PVC. PVCs are characterized by an abnormal QRS complex, which is wide and has an unusual shape. The drop to -1.395mV could be interpreted as part of an abnormally wide and bizarre QRS complex, which is a hallmark of PVCs. However, the presence of a preceding P-wave-like activity and the recovery back towards a more normal rhythm pattern complicate the interpretation slightly.

**Step 5: Final Analysis and Conclusion**

- The data shows a significant deviation from the normal ECG pattern, especially with the dramatic drop after the peak, which is not characteristic of a normal T-wave. The overall pattern suggests a significant abnormality in the ventricular repolarization process. Given the sharp rise indicative of the QRS complex followed by an abnormal repolarization phase, the pattern is more suggestive of a PVC due to the abnormal and dramatic changes in voltage that do not follow the typical progression of a normal heartbeat.

**ANSWER: Premature ventricular contraction beat (V)**

**This conclusion is based on the abnormal morphology of the waveform, particularly the unusual negative deflection following what appears to be the QRS complex, which does not align with the expected pattern of a normal T-wave. This analysis highlights the importance of examining the entire ECG waveform, including the amplitude, intervals, and progression of electrical activity, to differentiate between normal heartbeats and arrhythmias such as PVCs.**

Figure 5: Response example of LLM for heartbeat anomaly detection.

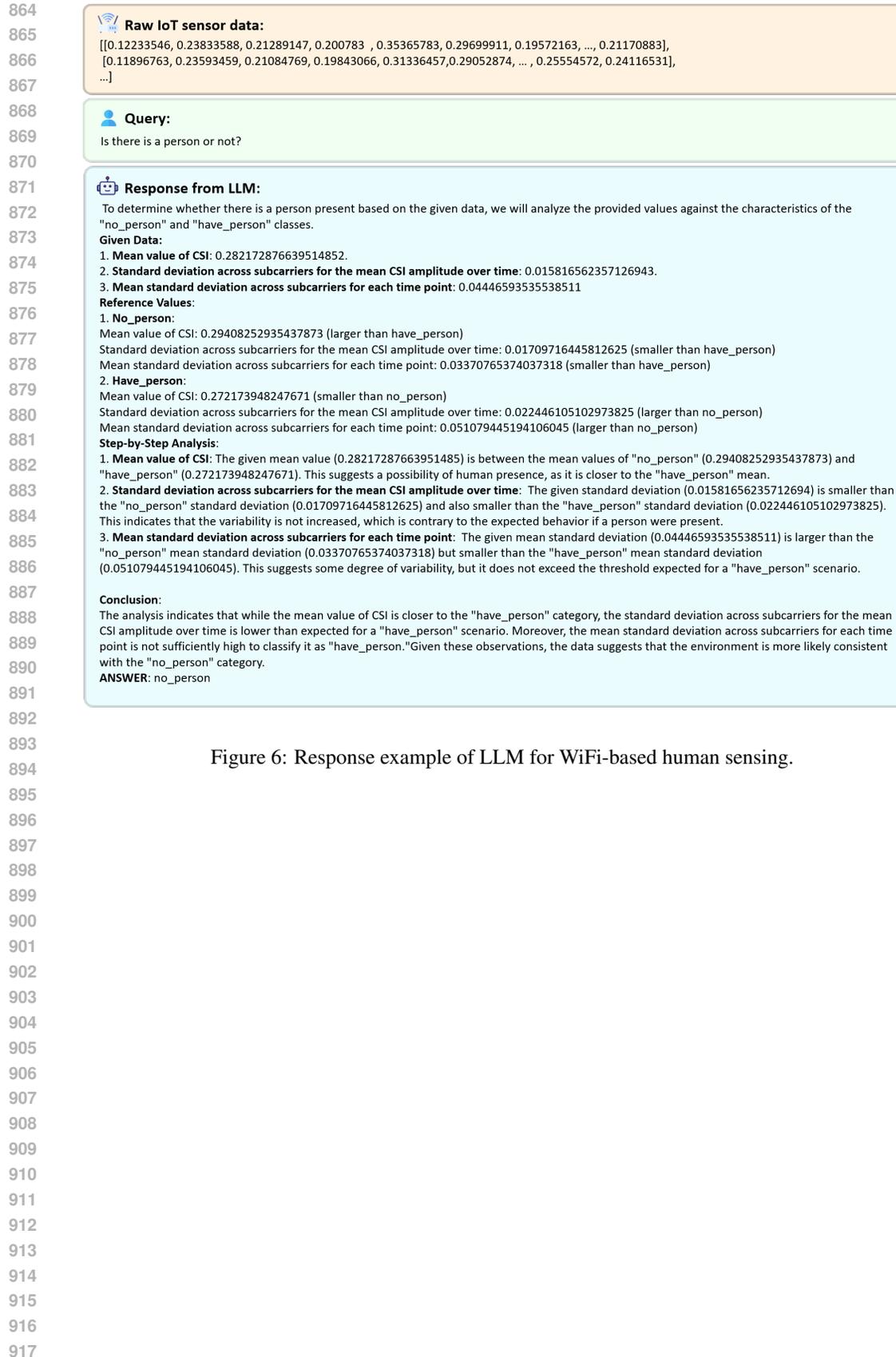


Figure 6: Response example of LLM for WiFi-based human sensing.

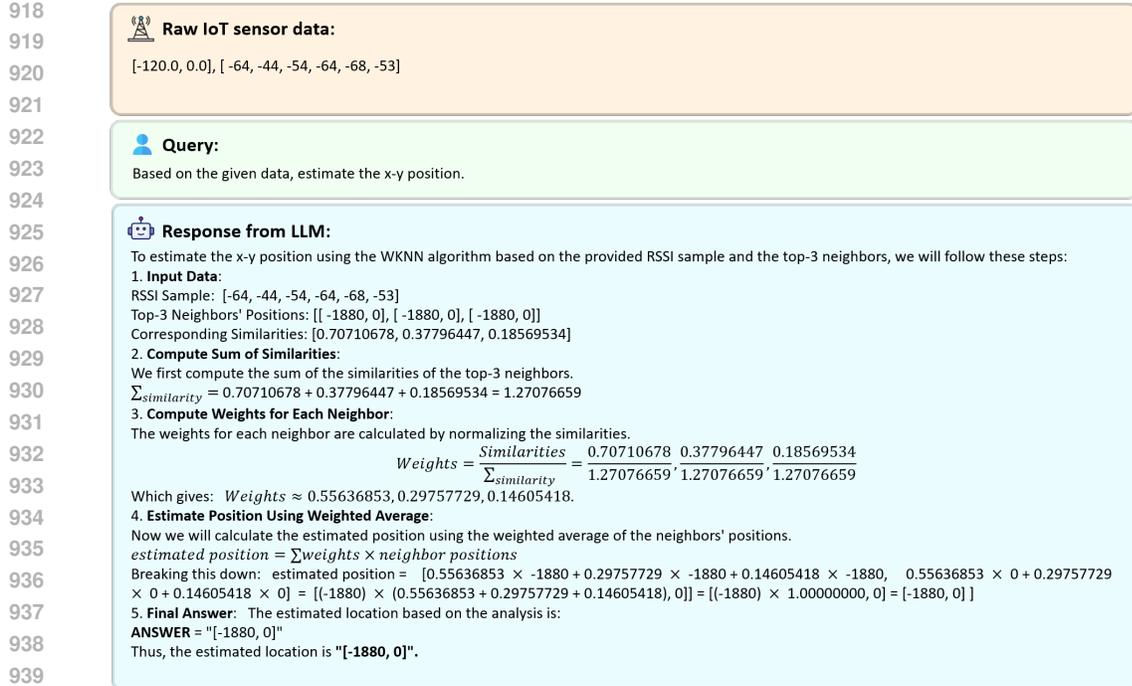


Figure 7: Response example of LLM for WiFi-based indoor localization.

## B PROMPT TEMPLATE

In the Prompt Configuration stage within our framework, we systematically organize IoT data description, task description, retrieved pertinent knowledge (including IoT domain knowledge and task-specific demonstrations), and role description to generate the final prompt according to the prompt template, as shown in Fig. 8. For example, based on the final prompt template, we obtain the final prompt (as shown in Fig. 9) for heartbeat anomaly detection.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

**{Role definition}**  
The combined application of these knowledge and skills would enable you to accurately perform the task and provide relevant analysis and interpretation

EXPERT:  
**{Retrieved domain knowledge}**

EXAMPLES:  
**{Retrieved task-specific demonstrations}**

THE GIVEN DATA:  
**{IoT data description}**

QUESTION:  
**{Task description}**  
Please analyze the data step by step to explain what it reflects, and then provide your answer based on your analysis.

ANALYSIS:  
ANSWER:

Figure 8: Final prompt template.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

You are an experienced physician who is familiar with various types of electrocardiogram (ECG) data. You can easily make preliminary judgments on whether heartbeats are abnormal based on ECG data. You possess the following medical and domain knowledge:

1. ECG Interpretation: You understand the basic principles of electrocardiography and know how to interpret ECG waveforms, including identifying different phases of the cardiac cycle and recognizing abnormalities.
2. Cardiac Physiology: You are familiar with the physiological functions of the heart, the generation and propagation of cardiac electrical signals, and the characteristics and manifestations of various cardiac arrhythmias.
3. Recognition of ECG Abnormalities: You are able to identify abnormal waveforms in ECG data, such as arrhythmias, myocardial ischemia, myocardial infarction, etc., and differentiate them from normal ECG patterns.
4. Medical Statistics: You are proficient in statistical analysis of ECG data, identification of outliers, and quantitative assessment of abnormalities.
5. Clinical Experience: You have extensive clinical experience to integrate ECG data with patient symptoms and medical history for accurate diagnosis and evaluation.
6. Medical Ethics and Legal Knowledge: You understand medical ethics and legal regulations to ensure confidentiality and lawful use of patient data.

The combined application of these domain knowledge and skills would enable you to accurately assess whether there are any abnormalities in the ECG data and provide relevant analysis and interpretation.

EXPERT:

Electrocardiography is the process of producing an electrocardiogram (ECG or EKG[a]), a recording of the heart's electrical activity through repeated cardiac cycles.[4] It is an electrogram of the heart which is a graph of voltage versus time of the electrical activity of the heart[5] using electrodes placed on the skin. In clinical applications, labeled ECG data are used to build a heartbeat classification system. Then this system is used to determine the types of heartbeats in unknown patients' ECG recordings. The overall magnitude of the heart's electrical potential is then measured from twelve different angles ("leads") and is recorded over a period of time (usually ten seconds). In this way, the overall magnitude and direction of the heart's electrical depolarization is captured at each moment throughout the cardiac cycle. A premature heart beat or extrasystole[1] is a heart rhythm disorder corresponding to a premature contraction of one of the chambers of the heart. Premature heart beats come in two different types: premature atrial contractions and premature ventricular contractions. Diagnosis Normal sinus rhythm and ectopic beats - premature ventricular contractions (PVC) and premature atrial contractions (PAC) shown on an EKG PVCs may be found incidentally on cardiac tests such as a 12-lead electrocardiogram (ECG/EKG) performed for another reason. In those with symptoms suggestive of premature ventricular complexes, the ECG/EKG is the first investigation that may identify PVCs as well as other cardiac rhythm issues that may cause similar symptoms. You can analyze whether the heartbeat is normal by considering a combination of factors such as the amplitude of peaks or valleys appearing in the electrocardiogram (ECG) time series, the time intervals between adjacent peaks or valleys, and the fluctuations in voltage values within the ECG data.

EXAMPLE1:

THE GIVEN ECG DATA:

-0.39mV, -0.38mV, -0.36mV, -0.355mV, -0.35mV, -0.37mV, -0.365mV, -0.35mV, -0.335mV, -0.35mV, -0.345mV, -0.355mV, -0.35mV, -0.33mV, -0.32mV, -0.295mV, -0.29mV, -0.295mV, -0.285mV, -0.23mV, -0.165mV, -0.08mV, -0.03mV, -0.1mV, -0.15mV, -0.185mV, -0.185mV, -0.16mV, -0.095mV, 0.325mV, 1.02mV, 0.53mV, -0.15mV, -0.22mV, -0.26mV, -0.305mV, -0.335mV, -0.31mV, -0.295mV, -0.275mV, -0.27mV, -0.245mV, -0.21mV, -0.145mV, -0.09mV, -0.03mV, 0.04mV, 0.105mV, 0.23mV, 0.365mV, 0.475mV, 0.52mV, 0.48mV, 0.375mV, 0.275mV, 0.16mV, 0.075mV, 0.0mV, -0.08mV, -0.135mV

ANSWER: Normal heartbeat (N)

EXAMPLE2:

THE GIVEN ECG DATA:

-0.55mV, -0.41mV, -0.29mV, -0.24mV, -0.16mV, -0.105mV, -0.08mV, -0.04mV, -0.055mV, -0.005mV, 0.085mV, 0.21mV, 0.42mV, 0.63mV, 0.785mV, 0.875mV, 0.9mV, 0.83mV, 0.705mV, 0.56mV, 0.405mV, 0.21mV, 0.125mV, 0.12mV, 0.08mV, 0.065mV, 0.06mV, 0.18mV, 0.48mV, 1.05mV, 1.57mV, 1.25mV, 0.81mV, 0.9mV, 0.05mV, -0.365mV, -0.525mV, -0.69mV, -0.76mV, -0.62mV, -0.79mV, -0.68mV, -0.685mV, -0.735mV, -0.785mV, -0.795mV, -0.82mV, -0.775mV, -0.7mV, -0.6mV, -0.485mV, -0.355mV, -0.24mV, -0.14mV, -0.115mV, -0.11mV, -0.115mV, -0.11mV, -0.085mV, -0.095mV

ANSWER: Premature ventricular contraction (V)

THE GIVEN DATA:

-0.205mV, -0.195mV, -0.155mV, -0.05mV, -0.015mV, 0.02mV, 0.125mV, 0.21mV, 0.345mV, 0.47mV, 0.615mV, 0.675mV, 0.705mV, 0.655mV, 0.6mV, 0.47mV, 0.36mV, 0.2mV, 0.135mV, 0.095mV, 0.045mV, 0.09mV, 0.025mV, 0.055mV, 0.05mV, 0.065mV, 0.14mV, 0.27mV, 0.29mV, 0.825mV, 1.35mV, 0.655mV, -1.18mV, -0.96mV, -0.89mV, -0.67mV, -0.47mV, -0.335mV, -0.27mV, -0.145mV, -0.12mV, -0.08mV, -0.1mV, -0.07mV, -0.005mV, 0.055mV, 0.155mV, 0.31mV, 0.52mV, 0.705mV, 0.875mV, 0.92mV, 0.865mV, 0.75mV, 0.6mV, 0.44mV, 0.245mV, 0.19mV, 0.125mV, 0.055mV

The ECG data is collected from a patient's heart. The data consists of a series of electrical signals that represent the heart's electrical activity. The signals are measured in millivolts (mV) and are recorded over a period of time at the sampling frequency of 60Hz. This means there is an interval of 0.017 seconds between the two voltage values. The data is divided into two categories: normal heartbeats (N) and ventricular ectopic beats (V). The normal heartbeats represent the regular electrical activity of the heart, while the ventricular ectopic beats represent abnormal electrical activity. The data is collected using a single-channel ECG device.

QUESTION:

Is the ECG heartbeat signal normal or abnormal?

Please analyze the data step by step to explain what it reflects, and then provide your answer based on your analysis.

ANALYSIS:

ANSWER:

Figure 9: Final prompt for heartbeat anomaly detection. Note that role description is generated automatically by AI models (e.g., ChatGPT).

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133



Figure 10: During IoT data simplification and enrichment stage, raw IoT data is transformed into IoT data description, which is easier to understand by LLMs. Raw IoT data is enriched with descriptive metadata, including natural language expressions of implicit physical information like units. Specialized tokenization techniques and extraction of temporal or frequency domain features further enhance LLMs' understanding of numerical and time-series data. These improvements make IoT data more accessible and interpretable for LLMs, facilitating its use in real-world applications.