

# Are LLMs Truly Multilingual?

## Exploring Zero-Shot Multilingual Capability of LLMs for Information Retrieval: An Italian Healthcare Use Case

Vignesh Kumar Kembu<sup>1,2</sup>[0009-0002-3782-8111],  
 Pierandrea Morandini<sup>2</sup>[0000-0002-8615-3766],  
 Marta Bianca Maria Ranzini<sup>2</sup>[0000-0001-8275-6028], and  
 Antonino Nocera<sup>1</sup>[0000-0003-2120-2341]

<sup>1</sup> Department of Electrical, Computer and Biomedical Engineering,  
 University of Pavia, Italy  
 vigneshkumar.kembu01@universitadipavia.it, antonino.nocera@unipv.it  
<sup>2</sup> IRCCS Humanitas Research Hospital, Milan, Italy  
 {vignesh.kembu,pierandrea.morandini,marta.ranzini}@humanitas.it

**Abstract.** Large Language Models (LLMs) have become a key topic in AI and NLP, transforming sectors like healthcare, finance, education, and marketing by improving customer service, automating tasks, providing insights, improving diagnostics, and personalizing learning experiences. Information extraction from clinical records is a crucial task in digital healthcare. Although traditional NLP techniques have been used for this in the past, they often fall short due to the complexity, variability of clinical language, and high inner semantics in the free clinical text. Recently, Large Language Models (LLMs) have become a powerful tool for better understanding and generating human-like text, making them highly effective in this area. In this paper, we explore the ability of open-source multilingual LLMs to understand EHRs (Electronic Health Records) in Italian and help extract information from them in real-time. Our detailed experimental campaign on comorbidity extraction from EHR reveals that some LLMs struggle in zero-shot, on-premises settings, and others show significant variation in performance, struggling to generalize across various diseases when compared to native pattern matching and manual annotations.

**Keywords:** LLMs · Multilingual · Information Retrieval · Healthcare · EHRs.

## 1 Introduction

Large Language Models (LLMs) have revolutionized the field of natural language processing, showcasing impressive capabilities in text generation, comprehension,

and conversational interaction. The models, such as GPT’s and Google’s Bard etc., are based on advanced neural networks with billions of parameters. They can grasp context, semantics, and intricate language nuances, which allows them to perform exceptionally across diverse applications—from chatbots and virtual assistants to content creation and programming assistance. However, despite their strengths, LLMs encounter several challenges [21,16]. They can generate incorrect information, may misinterpret subtle input variations and have the potential to produce biased or inappropriate content. Ongoing research seeks to address these issues through improved training techniques, fine-tuning processes and the establishment of ethical guidelines. Today, LLMs play a crucial role in various sectors, including education, healthcare and customer service, offering innovative solutions while also raising important questions about security, privacy and the future of human-computer interactions. As these models continue to evolve, their potential to transform communication and creativity appears huge, presenting both exciting opportunities and complex challenges for society [6,28].

In context of healthcare LLMs are revolutionizing healthcare by enhancing clinical decision-making, automating administrative processes and improving patient engagement. These advanced AI systems, trained on large datasets, can interpret and generate human-like text, making them valuable for various applications in healthcare. One significant area of impact is in clinical decision support. LLMs can analyze medical literature, patient records and research data to provide evidence-based recommendations to healthcare professionals. By synthesizing complex information quickly, they aid clinicians in diagnosing conditions and selecting appropriate treatments, ultimately improving patient outcomes and are streamlining administrative tasks such as scheduling, billing and documentation [17,11,24].

**Scenario.** In a clinical context, we consider a clinician aiming to extract comorbidities from EHRs using LLMs, focusing on a simple and direct approach, which does not need any manipulation of the LLM prompt (Zero-shot). Considering this, we formulate three key questions to be addressed in the evaluation of the scenario.

- **Q1.** Can we use LLMs in Zero-shot to extract comorbidities from EHRs?
- **Q2.** Can LLMs substitute a regular expression-based approach?
- **Q3.** Can we find a best LLM among the chosen ones?

Extending from the scenario, we evaluate LLMs (Large Language Models) in multilingual settings within digital health and clinical decision support. This research explores the use of different LLMs in the healthcare domain, specifically focusing on the classification tasks involving Italian patient EHRs. The objective was to design, develop a generalized data gathering ETL of patient electronic health records, LLM pipeline and framework which can handle different data formats and ranges of models with different parameters, capable of extracting diverse types of information from clinical records. In context with this, our current study leverages usage of 6 open-source models from three model families. The inference was carried out in Zero-shot setting on free EHR text of patients, which are discussed below in detail.

## 2 Preliminaries

**Large Language Models (LLMs)** are advanced models designed to process and generate human language, trained on vast amount of text. These models use network architecture called transformers [25], which help them manage and produce text as in human communication. These models are typically built with billions of parameters and more, allowing them to capture internal patterns in language, context and reasoning [15].

*Closed* source models, such as GPT [5,1], Gemini [23,7], and Claude [4,3], are proprietary and accessible via APIs, enabling businesses to integrate advanced AI without building models from scratch. But their closed nature raises security concerns, particularly in sensitive areas like healthcare. On the other hand, *Open* source models like Qwen’s [20,27], Bloom’s [26], Llama’s [2,10], and Mistral’s [13,14], offer transparency and customization, allowing fine-tuning for specific tasks, such as healthcare applications, while providing better control over data privacy and regulatory compliance.

**Information Extraction (IE)** is a crucial task in natural language processing (NLP) that involves automatically extracting structured information from unstructured text. This process is key for converting large volumes of textual data, such as news articles, medical text (EHRs) and social media posts, into usable information for further analysis or decision-making[18]. Large Language Models (LLMs), such as GPT and BERT, have significantly advanced the field of Information Extraction (IE). These models, trained on massive corpora, excel at identifying entities, relationships and other structured information in diverse and unstructured text without requiring task-specific training data [8]. By leveraging their deep contextual understanding, LLMs have shown superior performance in extracting nuanced and complex information, enabling more accurate and adaptable IE systems across various domains [5].

**Multilingual** understanding and generation have made notable progress through models trained on large and diverse multilingual data combined with advanced training techniques. Large language models demonstrate impressive robustness in English, leveraging abundant data and resources. Evidence on their performance and reliability in other languages remains limited and underexplored [19]. LLMs show potential in healthcare, helping in medical Q&A, diagnosis, counseling, and EHR data extraction, even in multilingual settings [29].

**Electronic Health Records (EHRs)** are digital information that contains patients medical history, such as diagnoses, treatments, medications, laboratory results and documentation notes from clinicians. EHRs improve continuity of care and efficiency by allowing the update of patient information in real time [12]. However, most EHR data is unstructured, making it difficult to extract meaningful insights. Information extraction on EHR data was traditionally addressed by regular expressions, however building the pattern is complicated, not generalizable and requires deep domain knowledge with all linguistics nuances of the domain. Currently, Large Language Models (LLMs) have great potential to process both structured and unstructured data, enabling them to extract valuable insights for clinical decision-making and research [17]. LLMs can extract

data from clinical notes with high accuracy, outperforming traditional pattern-matching methods. They are particularly effective in retrieving data, such as lab results and vital signs, which are crucial for clinical analysis. [11].

### 3 Methodology

The methodology comprises of EHR data gathering pipeline, first we applied regular expressions to automatically annotate the texts, establishing a measure for comparison. Then we leveraged a LLMs to extract comorbidities from the same set of EHRs, comparing the LLMs outputs to the Regexp annotations. To manage inaccuracies in the regexp annotations, we chose 100 regexp-false classified EHRs and proceeded with manual annotation by clinicians, creating a ground truth. Finally, further comparison of the LLMs performance against these manually annotated labels was carried out. Figure 1 presents the layout of the methodological approach.

#### 3.1 Data Pipeline & description

The protocol required collection of different data from the patient EHR from different hospital areas. For this, a large ETL using Oracle and Python has been implemented for the collection of the cardiological risk factors, presence of previous cardiac interventions, treatment received during the hospitalization, labs, discharge diagnosis, at-home treatment, and different dataframes have been created: Anamnesis, Interventions, Labs, Procedures, Procedures ICD9, Therapies. We will focus solely on *Anamnesis* data which has been obtained by performing a large ETL process because it represents unstructured text with the comorbidity information of interest. The dataframe consists of three columns, which are: Nosologico (hospital admission ID), Ente (hospital identifier), Anamnesis (EHR free text). With a total of 8223 patient records.

To emphasize the IR part, we focused on the following comorbidities which are clinically relevant for patient evaluation in the cardiac domain: Fibrillazione atriale (Atrial fibrillation), Insufficienza Renale (Kidney failure), BPCO - Broncopneumopatia cronica ostruttiva (COPD-Chronic obstructive pulmonary disease), Diabete mellito (Diabetes mellitus) and Ipertensione arteriosa (Hypertension) will be the key focus of our research. Since the EHR is in the Italian language, all the keys shall be used as per it.

#### 3.2 Automated Data Annotation using Regular expression

Building the regexp for the key comorbidities was done in collaboration with the help of clinicians, since the electronic health records are created and updated by them for each patient. So, to create the best regexp classifier, domain experts have been involved in developing the patterns. An example of why we need domain experts for this regexp creation is as follows, Different clinical mentions for

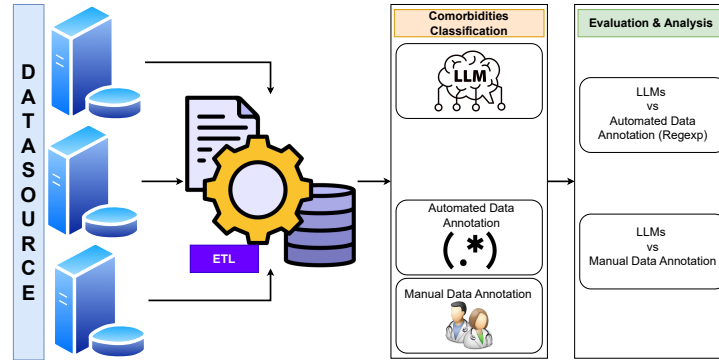


Fig. 1: Methodology - Data gathering pipeline, data annotation and comparison: from regex based automatic classification and clinicians validated ground truth to LLM extraction.

the same term: Term: “Diabetes” - Diabetes mellitus - refers to the general disease entity, DM - Abbreviation for Diabetes Mellitus, often used in medical notes and prescriptions, Type 1 diabetes mellitus, Type 2 diabetes mellitus, Insulin-dependent diabetes mellitus (IDDM), Non-insulin-dependent diabetes mellitus (NIDDM), Diabetes mellitus with nephropathy, Diabetes mellitus with retinopathy, Diabetes mellitus with neuropathy, Gestational diabetes mellitus (GDM). A domain expert is essential for identifying key comorbidities and creating accurate regex patterns. These patterns are then applied to EHRs to classify comorbidities, producing baseline data for comparison with LLM results.

### 3.3 Manual Data Annotation

To be more precise on the created pattern for data classification using Automated Data Annotation using Regular expression with the help of clinicians, 100 “false” classified records (i.e., regex classified as negative(0) for the presence of the comorbidities) of the patient on all five key comorbidities have been manually annotated with help of clinicians (Doctors). Although EHRs aim to improve patient care, they also pose several challenges. One major concern is data quality, as EHRs are only as good as the data entered by clinicians, which can be prone to errors and inaccuracies. EHRs can lack standardization, resulting in varying formats and terminology across different systems and organizations, making it challenging for healthcare providers to share information and coordinate care. Considering this, the regex created with the help of clinicians information might not be able to capture all the patterns of the comorbidities. So, manual annotation was carried out to double-check the “false” classified records of the patient. In the process of manual annotations, two clinicians annotated all five key comorbidities for each selected EHR. They annotated all five key comorbidities discussing case by case to reach a agreement of the class.

## 4 Experimental Setup

There are various factors to consider when designing the research setup. For our experiment, the primary criteria we have considered are discussed below.

**Privacy concerns & Licensing.** Data privacy is a concern with Large Language Models (LLMs), and organizations must address these issues when deploying them. Humanitas AI Center and Humanitas Research Hospital prioritize data privacy, healthcare patient data should be securely managed and used for research only on-premises. Recalling the difference between “Closed-source” and “Open-Source” models, only “Open-Source” ones have been selected as they allow for deployment in the “on-premises” environment.

**Language Support.** Language support in large language models (LLMs) refers to the number of languages a model can understand and generate text, as well as its proficiency in those languages. The models capabilities can vary significantly depending on their training data, architecture and intended use cases. They can be broadly classified as “Multilingual Models” and “Language-Specific Models”. Multilingual models are explicitly designed to handle multiple languages and often support dozens or even more languages. In our case, patients EHR is in the Italian language, thus only "Multilingual Models" have been chosen and the experiments have been carried out.

**Size and Resource Requirements.** The size and resource requirements of LLMs significantly influence their accessibility and usability. Large models typically yield impressive capabilities but come with high demands for computational power, memory and storage, making them more costly and complex to operate. Prior understanding of these requirements was crucial for us to adopt LLMs for the proposed information classification task. High-Performance Computing (HPC) with Multiple GPUs was used to carry out the proposed research.

**Selected LLMs.** Considering the points stated before, chosen models show great level of accuracy in the leaderboard [9], especially in consideration with the Italian language score. From OpenLLaMA family 3B & 7B models, from Mistral family 7B & 8x7B models and from Qwen2.5 family 3B & 7B models have been selected for this study.

## 5 Result Analysis

This section examines and compares the performance of cutting-edge large language models (LLMs) towards regular expression and human annotation in the IR task discussed in section 3.

### 5.1 Performance Comparison: LLMs vs Regexp

To evaluate the performance of various large language models (LLMs), it is essential to establish a reference dataset for comparison. With the availability of annotated data through automated methods, this dataset can serve as the

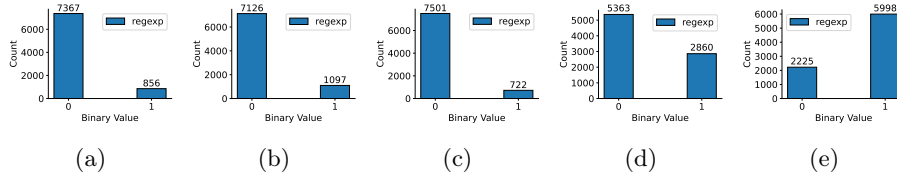


Fig. 2: Classification using regular expressions for the chosen comorbidities - a) Fibrillazione atriale, b) Insufficienza Renale, c) BPCO-Broncopneumopatia cronica ostruttiva, d) Diabete mellito and e) Ipertensione arteriosa.

reference. Consequently, the results generated by different LLMs can be compared against these actual values. The automated annotation is carried out by applying a series of different combinations of match patterns for the discussed comorbidities, utilizing regular expressions. From the Figure 2 we could see the classification, in which 0 class represents that the comorbidities is not found in the EHRs and 1 class represents the availability of the comorbidities in the EHRs. Diabete mellito and ipertensione arteriosa are most positive classified field by the regexp, comorbidities like Fibrillazione atriale, Insufficienza Renale and BPCO were the most negative classified.

**LLMs.** A standard prompt in a zero-shot setting has been used across all the LLMs in context with the scenario 1. To avoid multiple classifications in a single inference, each comorbidity has been classified in an individual inference for each EHR. Assigning one comorbidity per inference per task maximizes the accuracy of large language models (LLMs) and leads to more precise predictions. This will avoid the misinterpretation of data, especially medical information, which is prone to bias when several tasks are handled by one inference. In executing one task at a time, the model is able to give complete attention to that one task, thereby optimizing the performance. In the following, the classification results of different LLM family models against the regular expressions have been discussed.

*OpenLLaMA* model family generally show a increase in the classification accuracy w.r.t regular expression results of the comorbidities as the model size increases, which can be seen from the Figure 3. In particular with *OpenLLaMA* 3B, the model tends to perform well in terms of comorbidities like Diabete mellito and Ipertensione arteriosa with an accuracy of 34.78 % and 72.86 % compared to the other comorbidities which have an classification accuracy of less than 15 %. In contrast *OpenLLaMA* 7B has an accuracy of 50 % and above across all the comorbidities, Insufficienza Renale and BPCO have been among the top in classification with 83.73 % and 81.95 %. *Mistral* model family generally show a decrease in the classification accuracy w.r.t regular expression results of the comorbidities as the model size increases, which is shown in Figure 3. *Mistral* 7B shows a accuracy of 75 % above across all the comorbidities except Ipertensione arteriosa which is at 57.33 %, Fibrillazione atriale and BPCO have shown a accuracy of 90 % above which is the highest among all the models. *Mixtral* 8x7B shows a varied performance, comorbidities like BPCO, Diabete mellito

and Ipertensione arteriosa classification metrics are more equivalent to the performance OpenLLaMA 3B. Fibrillazione atriale and Insufficienza Renale show a slight increase in accuracy when compared to OpenLLaMA 3B. In contrast *Qwen2.5* model family does not show much of a increase or decrease with increase in the model parameters, which is shown in the Figure 3. *Qwen2.5* 7B which shows a very small marginal difference w.r.t to the 3B model, especially only for comorbidities like Fibrillazione atriale, Insufficienza Renale and BPCO. This shows that the domain-specific features in the dataset are being handled similarly by both models and the differences in model size is not substantial enough to affect their performance.

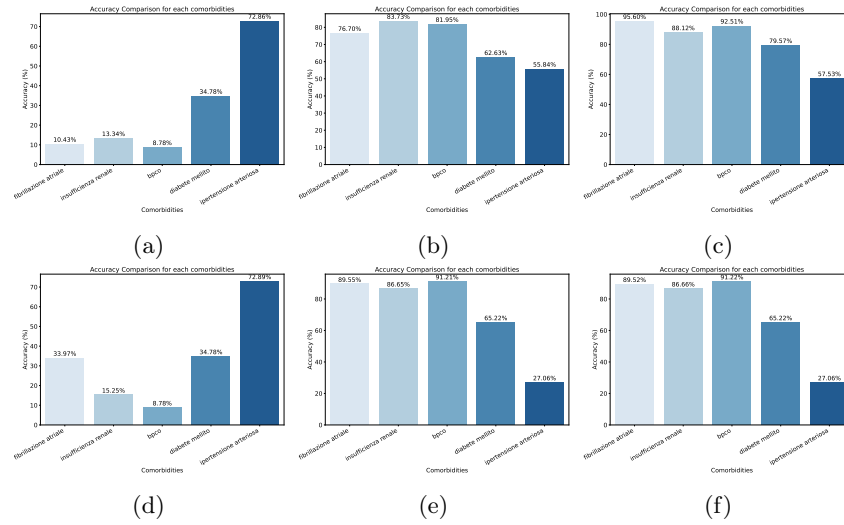


Fig. 3: LLMs accuracy compared to regular expression -a)OpenLLaMA 3B, b)OpenLLaMA 7B, c)Mistral 7B, d)Mixtral 8x7B, e)Qwen2.5 3B and e)Qwen2.5 7B

From the Figure 4, we could compare the overall accuracy of the selected models to the a) automated and b) manual annotated data as the reference values. It shows the mean accuracy of the models across all the defined comorbidities. From a) of the Figure 4 compared to the automated annotation, we could see that OpenLlama 3B and Mixtral 8x7B overall accuracies were lower than 35 %. In contrast, OpenLlama 7B, Mistral 7B, Qwen2.5 3B and 7B both had an overall accuracy of more than 70 %, which is nearly two times the performance of the other two models. Mistral 7B performed very well when compared to the other models, with an overall accuracy of 82.67 % when compared to regular expression (automated annotation).

While the classification accuracy provides a general insight into the LLMs performance, it is important to examine the classification reports for a deeper



understanding. Metrics such as F1-score, precision and recall help to assess how well the model performs across different comorbidities especially in the presence of data imbalance, and highlight potential areas where it may be struggling. This is needed for evaluating the LLMs true effectiveness.

LLMs and Metrics		Automated (8223 data records)				
		Fibrillazione Atriale	Insufficienza Renale	BPCO	Diabete Mellito	Ipertensione Arteriosa
OpenLLaMA 3B	Precision	0.1	0.13	0.09	0.35	0.73
	Recall	1	1	1	1	1
	F1-score	0.19	0.24	0.16	0.52	0.84
OpenLLaMA 7B	Precision	0.23	0.19	0.22	0.46	0.8
	Recall	0.52	0.07	0.42	0.47	0.53
	F1-score	0.31	0.1	0.29	0.47	0.64
Mistral 7B	Precision	0.86	0.96	1	0.99	0.99
	Recall	0.69	0.11	0.15	0.42	0.42
	F1-score	0.77	0.2	0.26	0.59	0.59
Mixtral 8x7B	Precision	0.10	0.13	0.09	0.35	0.73
	Recall	0.69	0.98	1	1	1
	F1-score	0.18	0.24	0.16	0.52	0.84
Qwen2.5 3B	Precision	0	0	0	0	0
	Recall	0	0	0	0	0
	F1-score	0	0	0	0	0
Qwen2.5 7B	Precision	0.25	0	0	0	0
	Recall	0	0	0	0	0
	F1-score	0.1	0	0	0	0

Table 1: Precision, Recall and F1 Score for LLMs vs. Regular Expressions in Identifying Comorbidities (Class 1).

From the Table 1 we could compare comorbidity-wise LLMs performance. OpenLLama 3B has perfect recall (1.0) but struggles with precision, leading to many false positives. OpenLLaMA 7B model shows low precision in most of the conditions i.e, ranging from 0.19 for Insufficienza Renale to 0.8 for Ipertensione Arteriosa. OpenLLama 3B with high recall, low precision and low F1, the model struggles to generalization across comorbidities. OpenLLama 7B shows better generalization than 3B, but still has low scores across various comorbidities.

Mistral 7B demonstrates outstanding precision across all comorbidities, which shows it is very effective at avoiding false positives. But recall is low compared to high precision, indicating the model often misses many cases of the comorbidities. In terms of F1-score, the model shows better balanced results across all comorbidities, compared to other LLMs. Mixtral 8x7B model exhibits low precision and high recall across all comorbidities. Due to this, the F1-scores are low for all except for Ipertensione Arteriosa is 0.84, showing the model achieves a better balance between precision and recall for this comorbidity. Mistral 7B shows better generalization, but still has low scores across a few comorbidities. Mixtral 8x7B has low scores in most cases and struggles to generalize across comorbidities, even though it is the most advanced among all the chosen models.

Both Qwen2.5 3B and Qwen2.5 7B models, despite their parameter differences, seem to produce classes that do not match the true positive class at all in most of the cases. The precision is 0 in most of the cases because there are no positive predictions. Similarly, recall is also 0 because the model does not correctly identify any of the actual positive samples in the data. In context with

these, F1-score is also 0 for most of the comorbidities. Qwen2.5 7B exhibits a slight difference for Fibrillazione Atriale with precision(0.25) and F1-score(0.1). Overall, Qwen2.5 family could not produce true positive class when automated data was compared.

LLMs in a zero-shot setting struggle to extract comorbidities from EHRs and do not match the accuracy of regular expression-based extractions, making them unsuitable as a substitute. These results are linked to re-search questions **Q1**) and **Q2**).

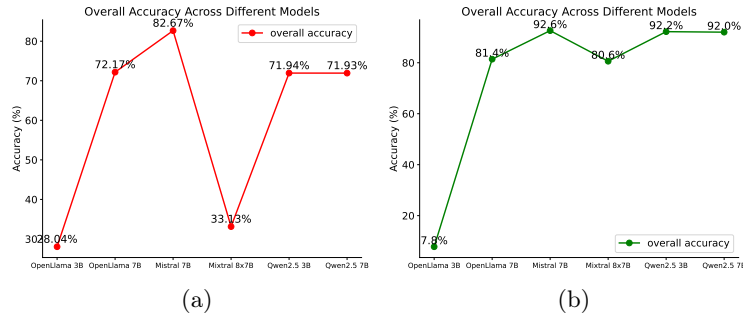


Fig. 4: Overall accuracy across different models compared to a) regular expression annotation and b) manual annotation.

## 5.2 Performance Comparison: LLMs vs Humans

As discussed in section 3 100 *false* classified comorbidities has been manually annotated by clinical experts. This is because positive classifications contain comorbidity information in the EHRs. False negatives in regex are critical in healthcare, as missed diagnoses can delay treatment, progression of disease and lead to poor patient outcomes. By annotating these false negatives, we can ensure that critical misclassifications are corrected, improving the ability to identify important conditions.

During this process, clinicians were not informed about the nature of the dataset provided i.e., with only the false class from the regex annotations provided. This was done to avoid bias or prejudgment from influencing the annotation. The guidance provided to the annotators is same as before, 0 class to represent the comorbidities is not found and 1 class represent the availability the comorbidities in the EHRs. From the Figure 5 we could see that Ipertensione arteriosa is the one with more false negatives, followed by the Fibrillazione atriale. Both of these comorbidities have 10 % or more false classification by the regular expression. In the other hand Insufficienza Renale, BPCO and Diabete

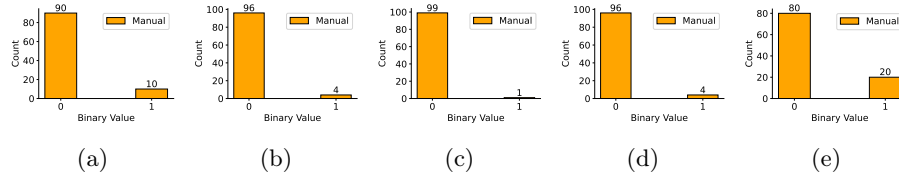


Fig 5: Manual annotation classification of the chosen comorbidities - a)Fibrillazione atriale, b)Insufficienza Renale, c)BPCO-Broncopneumopatia cronica ostruttiva, d)Diabete mellito and e)Ipertensione arteriosa.

mellito have 4 % or less false negative compared to the counterpart. In addition to this we can derive the performance of the regular expression created for each comorbidities when it is compared to the manual annotation, in particular pattern created for BPCO has the higher accuracy of 99 % and 80 % for Ipertensione arteriosa being the least accuracy. Overall comorbidities classification accuracy using regular expression when compared to the manual annotation was 92.2 %.

**LLMs.** With this new set of data, a zero-shot setting with a standard prompt is used across all the selected LLMs. As per previous experiment, to prevent multiple classifications in a single inference, each comorbidity is classified individually for each EHR. In the following we discuss the classification results of various LLM families against manual annotations.

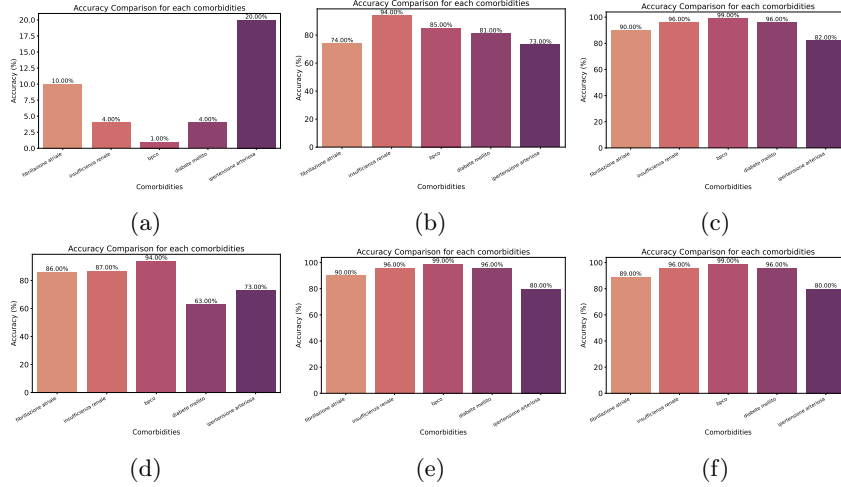


Fig 6: LLMs accuracy compared to manual annotation. (a)OpenLLaMA 3B, (b)OpenLLaMA 7B, (c)Mistral 7B, (d)Mixtral 8x7B, (e)Qwen2.5 3B and (f)Qwen2.5 7B

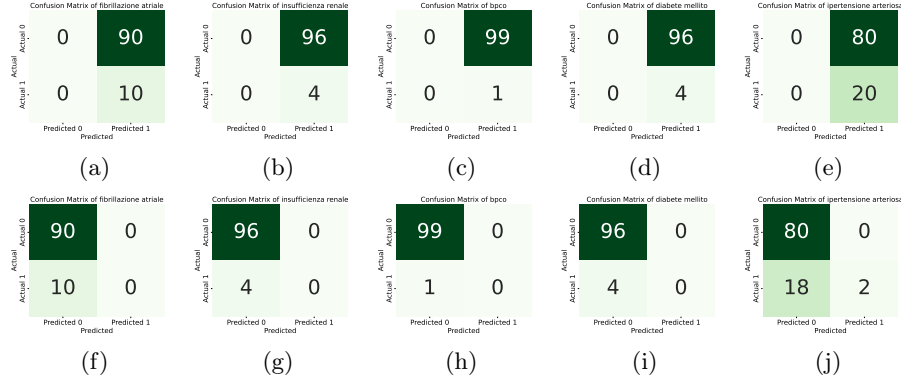


Fig. 7: Confusion matrix for LLMs when compared to manual annotation. The first row represents OpenLLaMA 3B and the second, Mistral 7B. Comorbidities are (a,f) Fibrillazione atriale, (b,g) Insufficienza Renale, (c,h) BPCO-Broncopneumopatia cronica ostruttiva, (d,i) Diabete mellito and (e,j) Ipertensione arteriosa.

*OpenLLaMA* family of models shows nearly the same kind of behavior as seen with the regular expression, increase in the model accuracy w.r.t the model size which can be seen in the Figure 6. OpenLLaMA 3B has varying range of classification accuracy when compared to the manual annotation set, starting from BPCO as low as 1 % to 20 % for Ipertensione arteriosa. In contrast OpenLLaMA 7B shows a extreme level of accuracy increase compared to the 3B version. In particular the model shows classification accuracy of greater than 70 % across all the comorbidities, notably 94 % accuracy for Insufficienza Renale and 85 % for BPCO. These results show that the model are mostly in alignment with the human annotation. *Mistral* family shows increased level of accuracy w.r.t to the previous regular expression case, the model Mixtral 8x7B showed a drastic performance increase in manual annotation. Mistral 7B shows 90 % and above classification accuracy in most of the comorbidities except for Ipertensione arteriosa which is 82 % and reaching 99 % for BPCO. Mistral 8x7B has a 94 % for BPCO and 63 % for Diabete mellito being the least, the model shows a overall performance increase across all the comorbidities when compared to the manual annotation. *Qwen2.5* family shows increased level of accuracy w.r.t to the regular expression case, notably a huge increase 52.94% for Ipertensione arteriosa and 30.78% for Diabete mellito. Qwen2.5 3B and 7B differ slightly in the Fibrillazione atriale and other comorbidities classification accuracy being identical.

From b) of the Figure 4 compared to the manual annotation, we could see that OpenLlama 3B shows a very low overall performance among all the models with a overall accuracy of less than 10 % when compared to the manual annotation. OpenLlama 7B, Mistral 7B, Mixtral 8x7B, Qwen2.5 3B and 7B showed a overall accuracy of more than 80 %. With mistral 7B and openllama 7B we see an

accuracy increase of about  $\approx 10\%$  when compared to the automated annotation, Qwen2.5 3B and 7B we see an accuracy increase of about  $\approx 20\%$  when compared to the automated annotation and Mixtral 8x7B we see a whopping 47.48 % of accuracy increase in the manual annotation when compared to the automated annotation. Regular expression had a accuracy of 92 % when compared to manual annotation, llama and mixtral 8x7B are below that and the others three reports a similar performance that of regexp. In consideration of the huge unbalanced class distribution post the manual annotation, we use confusion matrix to compare the results with the LLMs instead of classification report. From b) of the Figure 4 we take the best and least accuracy model and matrix is created and shown in Figure 7. This is done to evaluate whether the model is generalizing effectively or simply making random predictions. Openllama 3B being we could clearly see a pattern that the model always classifies the comorbidities as positives on gathered subset of data for manual annotation. Mistral 7B shows a pattern, by classifying most of comorbidities as negatives on the gathered subset. This shows how the LLMs classify the selected comorbidities from the EHRs, showcasing a NO better performance compared to the regular expression approach.

Even though we see varied performance of LLMs across the comorbidities in this setting, we could conclude that Mistral 7B seems to generalize considerably well across three out of five comorbidities compared to the other LLMs when regular expression is considered. However, when direct comparison with manual annotation is considered the model always predicts class 0 (the majority class), thus not showing real understanding on the semantics of the comorbidities (**Q3**).

## 6 Discussion

The results of the scenario 1 considered in our study, and the extended technical evaluation of the same, provide valuable insights into the usage of LLMs in zero-shot, on-premises settings for comorbidities extraction from Italian EHRs. The initial accuracy results of the LLMs when compared to the regular expression or the manual annotation seem to be good, but when a classification report and confusion matrix are constructed and analyzed, they show how these LLMs struggle to generalize across all the chosen comorbidities.

The results from our evaluation demonstrate that Mistral 7B can achieve performance in extracting a few comorbidities from EHRs among the selected LLMs when regular expression is considered for comparison, but cannot generalize across all the chosen comorbidities in accordance with manual annotation. Adding to this, neither of the chosen LLMs could reach the level of the regular expression-based pattern matching approach. It is evident that the selected multilingual LLMs face serious challenges and trustworthiness related issues when used in a zero-shot setting with no technical expertise in high-risk domains as Healthcare. In this state of usage, considering LLMs in place of a

pattern-matching approach in extraction pipelines would be problematic, since it would create uncertainty in the extracted information.

In context with the proposed scenario, when a clinician without much technical knowledge of prompt tuning techniques adheres to a simple and direct LLM prompt (zero-shot), aiming to extract comorbidities from EHRs, the overall results show that this approach would not be advisable in this setting. Although different prompt engineering approaches [22] could be employed to increase the accuracy of LLMs usage and making them more adaptive to the task intended, we on purpose did not attempt to optimize or engineer the use of LLMs, thus replicating a simulated environment, where a clinician or doctor would be expected to use these models directly.

Overall, we want to highlight that the increasing use of LLMs without technical understanding in high-risk domains like healthcare, law, finance, and autonomous systems etc., brings concern. In these high-stakes environments, misinformation, biased decision-making, lack of accountability and security vulnerabilities are significant risks when deploying LLMs in critical fields, as mistakes could lead to dangerous outcomes. These analyses showed that LLMs must be properly tested before deployment into production, global metrics should be accurately selected to avoid misleading conclusions and continuous close monitoring should be done for hallucinations and other deviations.

## 7 Conclusion

In this paper, we explored the potential of six multilingual general-purpose open-source large language models (LLMs) in understanding and extracting valuable information from Electronic Health Records (EHRs) in Italian. Our study, focused on the real-time retrieval of comorbidities from EHRs, shows important insights into the LLMs ability to handle non-English language understanding and processing. Our findings show how LLMs struggle in extracting comorbidities and their difficulty in understanding Italian language EHR when used in a zero-shot approach. LLMs in the zero-shot setting cannot be used as a substitute for traditional pattern matching in the extraction pipelines. In the future, we will explore In-Context learning (ICL) approaches and will also consider fine-tuning an LLM model for IR from EHRs to improve its language processing capabilities and enhance trustworthiness in handling healthcare-related tasks.

## References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
2. AI, M.: Llama 3.2: Connect 2024 vision for edge and mobile devices. <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/> (2024)
3. Anthropic: Introducing claude 4, <https://www.anthropic.com/news/claude-4>

4. Anthropic: Introducing the next generation of claude, <https://www.anthropic.com/news/claude-3-family>
5. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners (2020), <https://arxiv.org/abs/2005.14165>
6. Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang, C., Wang, Y., Ye, W., Zhang, Y., Chang, Y., Yu, P.S., Yang, Q., Xie, X.: A survey on evaluation of large language models. *ACM Trans. Intell. Syst. Technol.* **15**(3) (Mar 2024). <https://doi.org/10.1145/3641289>, <https://doi.org/10.1145/3641289>
7. Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang, D., Rosen, E., et al.: Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261* (2025)
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (eds.) *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423/>
9. EuroLingua-GPT: Internal european leaderboard - a hugging face space by eurolingua, <https://huggingface.co/spaces/Eurolingua/european-llm-leaderboard>
10. Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al.: The llama 3 herd of models (2024)
11. Gu, B., Shao, V., Liao, Z., Carducci, V., Romero-Brufau, S., Yang, J., Desai, R.: Scalable information extraction from free text electronic health records using large language models. *medrxiv* (2024)
12. Häyrynen, K., Saranto, K., Nykänen, P.: Definition, structure, content, use and impacts of electronic health records: a review of the research literature. *International journal of medical informatics* **77**(5), 291–304 (2008)
13. Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L.R., Lachaux, M.A., Stock, P., Scao, T.L., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E.: Mistral 7b (2023), <https://arxiv.org/abs/2310.06825>
14. Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D.S., de las Casas, D., Hanna, E.B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L.R., Saulnier, L., Lachaux, M.A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T.L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E.: Mixtral of experts (2024), <https://arxiv.org/abs/2401.04088>
15. Kaplan, J., McCandlish, S., Henighan, T., Brown, T.B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., Amodei, D.: Scaling laws for neural language models (2020), <https://arxiv.org/abs/2001.08361>
16. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L.: Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension (2019), <https://arxiv.org/abs/1910.13461>

17. Li, L., Zhou, J., Gao, Z., Hua, W., Fan, L., Yu, H., Hagen, L., Zhang, Y., Assimes, T.L., Hemphill, L., Ma, S.: A scoping review of using large language models (llms) to investigate electronic health records (ehrs) (2024), <https://arxiv.org/abs/2405.03066>
18. Nadeau, D., Sekine, S.: A survey of named entity recognition and classification. *Linguisticae Investigationes* **30**, 3–26 (2007), <https://api.semanticscholar.org/CorpusID:8310135>
19. Qin, L., Chen, Q., Zhou, Y., Chen, Z., Li, Y., Liao, L., Li, M., Che, W., Yu, P.S.: A survey of multilingual large language models. *Patterns* **6**(1), 101118 (2025). <https://doi.org/https://doi.org/10.1016/j.patter.2024.101118>, <https://www.sciencedirect.com/science/article/pii/S2666389924002903>
20. Qwen, :, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., Lin, H., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Lin, J., Dang, K., Lu, K., Bao, K., Yang, K., Yu, L., Li, M., Xue, M., Zhang, P., Zhu, Q., Men, R., Lin, R., Li, T., Tang, T., Xia, T., Ren, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Wan, Y., Liu, Y., Cui, Z., Zhang, Z., Qiu, Z.: Qwen2.5 technical report (2025), <https://arxiv.org/abs/2412.15115>
21. Radford, A., Narasimhan, K., Salimans, T., et al.: Improving language understanding by generative pretraining. OpenAI (2019), <https://openai.com/research/language-unsupervised>
22. Sahoo, P., Singh, A.K., Saha, S., Jain, V., Mondal, S., Chadha, A.: A systematic survey of prompt engineering in large language models: Techniques and applications (2025), <https://arxiv.org/abs/2402.07927>
23. Team, G., Anil, R., Borgeaud, S., Alayrac, J.B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., Millican, K., et al.: Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805 (2023)
24. Thirunavukarasu, A.J., Ting, D.S.J., Elangovan, K., Gutierrez, L., Tan, T.F., Ting, D.S.W.: Large language models in medicine. *Nature medicine* **29**(8), 1930–1940 (2023)
25. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need (2023), <https://arxiv.org/abs/1706.03762>
26. Workshop, B., Scao, T.L., Fan, A., Akiki, C., Pavlick, E., Ilić, S., Hesslow, D., Castagné, R., Luccioni, A.S., Yvon, F., et al.: Bloom: A 176b-parameter open-access multilingual language model. arXiv preprint arXiv:2211.05100 (2022)
27. Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., Zheng, C., Liu, D., Zhou, F., Huang, F., Hu, F., Ge, H., Wei, H., Lin, H., Tang, J., Yang, J., Tu, J., Zhang, J., Yang, J., Yang, J., Zhou, J., Zhou, J., Lin, J., Dang, K., Bao, K., Yang, K., Yu, L., Deng, L., Li, M., Xue, M., Li, M., Zhang, P., Wang, P., Zhu, Q., Men, R., Gao, R., Liu, S., Luo, S., Li, T., Tang, T., Yin, W., Ren, X., Wang, X., Zhang, X., Ren, X., Fan, Y., Su, Y., Zhang, Y., Zhang, Y., Wan, Y., Liu, Y., Wang, Z., Cui, Z., Zhang, Z., Zhou, Z., Qiu, Z.: Qwen3 technical report (2025), <https://arxiv.org/abs/2505.09388>
28. Zhao, W.X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., Du, Y., Yang, C., Chen, Y., Chen, Z., Jiang, J., Ren, R., Li, Y., Tang, X., Liu, Z., Liu, P., Nie, J.Y., Wen, J.R.: A survey of large language models (2025), <https://arxiv.org/abs/2303.18223>
29. Zhu, S., Supryadi, Xu, S., Sun, H., Pan, L., Cui, M., Du, J., Jin, R., Branco, A., Xiong, D.: Multilingual large language models: A systematic survey (2024), <https://arxiv.org/abs/2411.11072>