

Analyzing Prompt Design Choices in Biomedical Information Extraction for Low-Resource Languages

Ayesha Khatun and Bulut Ozler and Steven Bethard and Egoitz Laparra

College of Information Science

University of Arizona

Tucson, AZ, USA

{ayeshakhatun,kbozler,bethard,egoitz}@arizona.edu

Abstract

Despite the multilingual abilities of modern LLMs, biomedical information extraction remains challenging for low-resource, morphologically rich languages such as Bangla/Bengali and Basque. Prior investigations of prompt design and output schemas have focused on high-resource languages. To bridge this gap, we systematically analyze biomedical named entity recognition (NER) with open LLMs under multiple prompting settings on Bangla, Basque, Spanish, and English. We find that span-based extraction is more effective than BIO tagging for LLM prompting across all languages, while moving from statement-based prompting to question-based prompting has a stronger effect on the low-resource languages (e.g., +57% for Bangla and +109% for Basque) than on the high-resource languages (+28% for English and +22% for Spanish). Our breakdowns by error type show that translation-based prompting cuts Bangla hallucinations by 64% and QA-style prompting lowers Basque empty prediction errors by 61%. Our results offer practical guidance for building reliable multilingual biomedical NER systems in Bangla and Basque. We release our code at https://github.com/clulab/LRL_IE/.

1 Introduction

Research on health and medical texts is important because it affects people’s lives, can reduce cost and workload, and supports better clinical decisions (Agrawal et al., 2022). However, most medical information appears in unstructured text, so medical named entity recognition (NER) is needed to identify and label key medical terms (Averly and Ning, 2025). Manual annotation, whether for direct use or for fine-tuning machine learning models, requires medical expertise, and is slow, expensive, and hard to scale. To reduce the need for manual annotation, prior work has explored

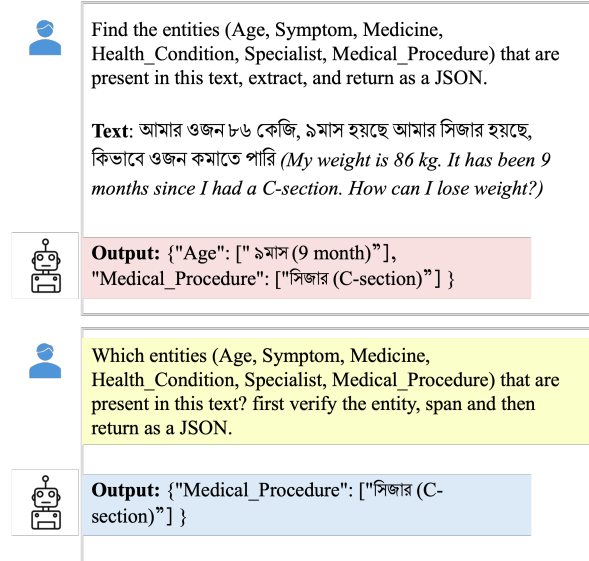


Figure 1: A statement-based prompt for biomedical NER in Bangla results in an incorrect extraction (red), while a question-based prompt results in a correct extraction (blue).

prompting and in-context learning with large language models (LLM) achieving great performance in high-resource languages like English (Ashok and Lipton, 2023; Li et al., 2023a; Li and Zhang, 2024; Wang et al., 2025; Averly and Ning, 2025). In contrast, for low-resource languages (languages with limited annotated datasets, smaller digital corpora, and fewer domain-specific resources), we still lack a systematic understanding of how different prompting strategies affect multilingual LLM performance in biomedical NER (Azime et al., 2025; Kumar et al., 2025).

We present a systematic analysis of LLM prompting strategies in Bangla/Bengali and Basque, two morphologically rich languages that pose challenging (Arnett and Bergen, 2025) low-resource settings (Bhattacharyya and Bhattacharya, 2025; López de Lacalle et al., 2020). Our analysis also includes Spanish and English as high-resource comparisons. We analyze different

prompted output representations (beginning–inside–outside tagging (BIO; [Ramshaw and Marcus, 1995](#)) and span-based extraction), prompt variations (statement, question, explanation, and translation), and common error types (entity hallucination, over-generation, boundary mistakes, etc.). We summarize these analyses in three research questions:

- Which output format is most effective for biomedical NER? Is it consistent across Bangla, Basque, Spanish, and English?
- Which prompt variation is most effective? Is it consistent across languages?
- Which errors are most common? Which prompt choices reduce hallucination, over-generation, and boundary errors?

We find that LLMs struggle to generate output in BIO format, but can succeed with span-based extraction. We find that although the best prompting strategy is often similar across languages, with question-style generally outperforming statement-style (as in [Figure 1](#)), selecting an optimal prompt is more critical in low-resource languages, where performance degrades much more under suboptimal prompting strategies. Our analyses of these four languages offer directions for future biomedical NER research in other low-resource and multilingual settings ([Pfeiffer et al., 2020](#); [Liu et al., 2021](#); [Choenni et al., 2023](#); [Pham et al., 2024](#)).

2 Related Work

Prior work has explored zero-shot and few-shot methods for Bangla across several NLP tasks (e.g., [Dementieva et al. 2025](#); [Adak et al. 2025](#); [Li et al. 2023b](#); [Shafayat et al. 2024](#)); for example, [Hasan et al. \(2024\)](#) use 3-shot and 5-shot prompting for sentiment analysis. Prompting has also been applied to general Bangla NER: [Mahtab et al. \(2025\)](#) provide an English instruction prompt describing the BIO tags and rules, include 10 in-context input–output examples, and then prompt the model to label a new Bangla sentence in BIO format. Early work on Bangla biomedical/telemedicine NER (e.g., [Islam et al. 2022](#); [Sazzed 2022](#)) primarily focused on dataset construction. More recently, [Khan et al. \(2023\)](#) introduce Bangla-HealthNER and evaluate fine-tuned BanglaBERT, Banglish-BERT, and mBERT models, and also report substantially lower performance for zero-shot ChatGPT than for supervised fine-tuning.

For Basque, there is prior work on biomedical NER, and most of it relies on supervised training or fine-tuning with BIO-style output formats. [Urbizu et al. \(2022\)](#) introduce the BasqueGLUE NLU benchmark; they use standard fine-tuning and report baseline results, and for the NER task they follow the BIO annotation scheme. [Zanoli et al. \(2024\)](#) examine whether a multilingual clinical corpus is effective for disorder NER; they train and fine-tune supervised NER models and use IBO/BIO-style tagging. Recent work has also proposed leaderboard-style benchmarking frameworks for evaluating LLMs across multiple tasks in high-resource settings, such as Italian ([Magnini et al., 2025](#)), but comparable systematic benchmarking remains limited for Bangla and Basque biomedical NER.

However, existing work for both Bangla and Basque largely emphasizes supervised BIO tagging, and we still lack a clear picture of how prompt design and output format affect LLM-based biomedical NER in these languages. To fill this gap, we systematically evaluate multiple prompting strategies and compare BIO tagging with span-based JSON extraction.

3 Methods

[Figure 2](#) details our NER evaluation pipeline and the different configurations we have explored in low-resource languages. We start from a base prompt that includes instructions for the task and some rules and hints to guide the LLM when annotating the entities. In all cases, the instructions are written in English ([Enomoto et al., 2025](#)) but specify the target language. Building on this base, we explore different configurations involving the output format of the annotations, whether or not to include demonstration examples, the format of these examples, or how the model should address the task.

3.1 Datasets

For our Bangla experiments, we use Bangla-HealthNER ([Khan et al., 2023](#)), a large Bengali biomedical NER dataset built from consumer health Q&A collected from a public online health platform in Bangladesh. It contains 31,783 samples and 144,136 sentences, annotated in token-level (BIO-style) format. The dataset includes seven entity types: Symptom, Health_Condition, Medicine, Specialist, Age, Dosage, and Medical_

Algorithm 1 BIO \rightarrow span-level JSON

Require: Tokens $x_{1:n}$, BIO tags $y_{1:n}$, entity types \mathcal{T}
Ensure: \mathcal{J} : map $t \in \mathcal{T} \mapsto$ list of extracted span strings

- 1: $\mathcal{J}(t) \leftarrow [] \quad \forall t \in \mathcal{T}$
- 2: $i \leftarrow 1$
- 3: **while** $i \leq n$ **do**
- 4: **if** $y_i = \text{B-}t$ for some $t \in \mathcal{T}$ **then**
- 5: $s \leftarrow i$
- 6: $i \leftarrow i + 1$
- 7: **while** $i \leq n$ **and** $y_i = \text{I-}t$ **do**
- 8: $i \leftarrow i + 1$
- 9: $e \leftarrow i - 1$
- 10: $\text{span} \leftarrow \text{CLEAN}(\text{DETOKEN}(x_{s:e}))$
- 11: **if** $\text{span} \neq \emptyset$ **then**
- 12: $\mathcal{J}(t) \leftarrow \mathcal{J}(t) \parallel [\text{span}]$
- 13: **else**
- 14: $i \leftarrow i + 1$
- 15: **return** \mathcal{J}

```
SYSTEM_PROMPT = (  
You are a biomedical NER assistant that performs medical  
Named Entity Recognition .....)   
Statement-style PROMPT = (  
"Find the entities (Age, Symptom, Medicine,  
Health_Condition, Specialist, Medical_Procedure) that are  
present in this text, extract, and return as a JSON."  
"Text: '\আমার বয়স ২৪ ওজন ৫৮ কেজি আমার কিছুখন পর  
পর প্রসাবে হয় খুবকম এখন আমি কি করতে পারি'\n"  
"Answer: {\n"Age": [], "  
"\n"Symptom": [{"প্রসাবের রাস্তায়  
জ্বালাপোড়া", "\n"ব্যাথা", "\n"প্রসাবে মাঝে মাঝে গন্ধ", "  
"\n"Medicine": [], "  
"\n"Health_Condition": [], "  
"\n"Specialist": [], "  
"\n"Medical_Procedure": []}\n\n")
```

Figure 4: A statement-style prompt for Bangla biomedical NER, where the yellow-highlighted text tells the model to find predefined medical entity types in the input text.

yields a set of extracted spans:

$$\hat{\mathcal{S}} = \left\{ (t, x_{s:e}) \mid \begin{array}{l} y_s = \text{B-}t, \forall i \in (s+1, \dots, e): y_i = \text{I-}t \\ y_{e+1} \neq \text{I-}t \end{array} \right\}$$

where \mathcal{T} is the set of entity types and $x_{s:e}$ is the token span reconstructed into text. We used [Algorithm 1](#) to detokenize each extracted span and append it to $\mathcal{J}(t)$. Code is available at https://anonymous.4open.science/r/LRL_IE-1959/.

3.3 Demonstration Examples

Zero-shot prompting (zero) is a strategy where we provide only the task instruction and the list of entity types, without any labeled examples (Liu et al., 2023). We provided detailed guidelines to the model, e.g., instructing it to skip negated symptoms, not confuse duration with age, not treat lab/test names as symptoms, and to select short, clean spans with no duplicates. This setting tests whether the model can perform biomedical extraction with only instructions and no demonstrations.

Few-shot prompting (few) is a way to use an LLM without training it, where we show the model a small number of labeled examples (Pan et al., 2023; Cheng et al., 2025) inside the prompt and then ask it to do the same task for a new input. Such labeled examples are known as the *support set*, \mathcal{S}_K , where K labeled examples are included in the prompt. Formally, we combine the base prompt P_{base} with this support set, and find the LLM’s most probable output, \hat{o}_{FS} , given this combined prompt:

$$P_{\text{FS}} = P_{\text{base}} \oplus \mathcal{S}_K$$
$$\hat{o}_{\text{FS}} = \arg \max_o p(o \mid x, P_{\text{FS}})$$

The few-shot prompt contains $K = 9$ in-context examples in the target language. Since there are nine entity types in total, we set $K = 9$ and designed the examples so that each entity type appears at least once. The selected set includes both short and long examples.

3.4 Prompt Variations

Statement-style prompting (stmt) appends a statement of the task s to the prompt while using the same support set \mathcal{S}_K :

$$P_{\text{QA}} = P_{\text{base}} \oplus \mathcal{S}_K \oplus s$$
$$\hat{o}_{\text{QA}} = \arg \max_o p(o \mid x, P_{\text{QA}})$$

This approach presents the biomedical entity extraction task as a statement, such as “Find the entities (AGE, SYMPTOM, ...) that are present in this text.”, as shown in [Figure 4](#).

Question-style prompting (qa) appends an explicit question q to the prompt while using the same support set \mathcal{S}_K :

$$P_{\text{QA}} = P_{\text{base}} \oplus \mathcal{S}_K \oplus q$$
$$\hat{o}_{\text{QA}} = \arg \max_o p(o \mid x, P_{\text{QA}})$$

This approach presents the biomedical entity extraction task as a direct question, such as “Which entities of types (AGE, SYMPTOM, ...) appear in this text?”, as shown in [Figure 5](#).

Explanation-based prompting (expl) is a prompting strategy ([Figure 6](#)) that provides a brief description of what entities are likely present in the current text. First, the prompt explains the label boundaries (e.g., HEALTH_CONDITION is a diagnosis or stated condition, while SYMPTOM is a complaint). Then, for each example, we add a short natural-language description in English that highlights the relevant cues in the input. We have

```

SYSTEM_PROMPT = (
You are a question-answering assistant that performs
medical Named Entity Recognition ..... )

QA-Style_PROMPT = (
"Question: Which entities (Age, Symptom, Medicine,
Health_Condition, Specialist, Medical_Procedure)
are present in this text?\n"

"Text:\nআমার বয়স ২৪ ওজন ৫৮ কেজি আমার কিছুখন
পর পর প্রসাব হয় খুবকম এখন আমি কি করতে
পারি!\n\n"
"Answer: {\nAge\":[], "
"\nSymptom\":[\nপ্রসাবের রাস্তায়
স্বালাপোড়া\n,\nব্যথা\n,\nপ্রসাবে মাঝে মাঝে গন্ধ\n,"
"\nMedicine\":[], "
"\nHealth_Condition\":[], "
"\nSpecialist\":[], "
"\nMedical_Procedure\":[]\n\n")

```

Figure 5: A question-style prompt for Bangla biomedical NER, where the yellow-highlighted text explicitly asks the model to identify predefined medical entity types in the input text.

used explanation only for examples, not the entire prompt. This prompt uses the same support set \mathcal{S}_K and adds label definitions $D = \{d_\ell\}_{\ell \in \mathcal{L}}$:

$$P_{\text{Desc}} = P_{\text{base}} \oplus \mathcal{S}_K \oplus D$$

$$\hat{o}_{\text{Desc}} = \arg \max_o p(o | x, P_{\text{Desc}})$$

Translation-based prompting (trans) is a process of prepending an English rendering (Figure 7) of the input while also keeping the original Bangla/Basque text in the prompt. This encourages the model to reason in English, which can help when the model is stronger in English than in the target language. However, this approach can fail when translation paraphrases the meaning, drops details, mistranslates medical terms, or changes span boundaries. Therefore, we treat translation-based prompting as a practical baseline rather than a guaranteed improvement. Let $T(\cdot)$ be a translation function (e.g., Bangla/Basque \rightarrow English). We use the same support set \mathcal{S}_K , but the input includes the translated text:

$$x' = T(x)$$

$$P_{\text{Trans}} = P_{\text{base}} \oplus \mathcal{S}_K$$

$$\hat{o}_{\text{Trans}} = \arg \max_o p(o | x', P_{\text{Trans}})$$

3.5 Language Models

We primarily conduct our experiments using **Meta-Llama-3-8B** (Grattafiori et al., 2024), a widely adopted open-weights model with strong general performance. However, to test the effect of different model families, we also evaluate the best-performing prompt strategy under Llama-3 on two additional models: **Qwen3-8B** (Yang et al.,

```

SYSTEM_PROMPT = (
You are a biomedical NER assistant that performs
medical Named Entity Recognition ..... )

Explanation-based_PROMPT = (
"Text:\nআমার পিতা খলিতে পাথর আছে। গত ২.৫ বছর
যাবৎ এটা হয়েছে। যার আকার ৪ সেমি। এতদিন তীব্র
কোন ব্যাথা ছিল না কিন্তু গত ৮ দিনের মধ্যে ৩ দিন
পেটে তীব্র ব্যাথা হয়েছিল এবং পেটের ডানদিকে পাজরের
নিচে চাপ দিলে ব্যাথা অনুভূত হয়। পিতা খলিতে পাথর
হলে কী পেট ব্যাথার সাথে সাথে পেটের ডানদিকে পাজরের
নিচে চাপ দিলে ব্যাথা অনুভূত হয়? এবং কেন এই চাপ
দিলে ব্যাথা অনুভূত হয়?? অভিজ্ঞ ডাক্তারের পরামর্শ
চাচ্ছি!\n\n"

"Description of text: The Bangla text says the
patient has gallbladder stones for about 2.5 years
and the stone size is 4 cm. The phrase "পিতা খলিতে
পাথর" is a Health_Condition entity (a
diagnosis/condition). Because of this condition,
the patient reports pain-related symptoms, such
as severe abdominal pain and pain when pressing
under the right ribs (these are Symptom
entities).\n\n"

"Answer: {\nAge\":[], "
"\nSymptom\":[\nব্যথা\n,\nপেটে তীব্র ব্যাথা\n,\nচাপ
দিলে ব্যাথা\n,\nপেটে ব্যাথার\n], "
"\nMedicine\":[], "
"\nHealth_Condition\":[\nপিতা খলিতে পাথর\n], "
"\nSpecialist\":[], "
"\nMedical_Procedure\":[]\n\n")

```

Figure 6: An explanation-based prompt for Bangla biomedical NER, where the yellow-highlighted text provides a natural-language explanation of the clinical context and entity semantics.

2025) and **Aya** (Aryabumi et al., 2024). We include Qwen3 because it is a competitive recent model that often performs particularly well on high-resource languages such as English, and Aya because it is designed for multilingual use and thus provides a useful contrast for low-resource settings. To ensure a fair comparison, we select similar size models (all in the $\sim 8B$ parameter range) and keep inference settings fixed.

4 Experimental Results

Our analyses are presented below. As datasets for different languages are by necessity different, we focus on relative gains, not absolute values, when looking across languages.

4.1 Output Format Analysis

Table 1 shows a large gap between BIO tagging and span-based extraction, indicating that the output schema strongly affects LLM-based NER. Across languages, switching from *bio_few* to *json_few_stmt* yields large F1 gains: +538% (Bangla; 0.054 \rightarrow 0.345), +200% (Basque; 0.104 \rightarrow 0.312), +2242% (Spanish; 0.019 \rightarrow 0.445), and +339% (English; 0.120 \rightarrow 0.527). BIO tagging is likely harder for LLMs because the model must label

<p>SYSTEM_PROMPT = (</p> <p>You are a biomedical NER assistant that performs medical Named Entity Recognition</p> <p>Translation-based_PROMPT = (</p> <p>"Text: \”আমার পিত খলিতে পাথর আছে। গত ২.৫ বছর যাবৎ এটা হয়েছে। যার আকার ৪ সেমি। এতদিন তীব্র কোল ব্যাথা ছিল না কিন্তু গত ৮ দিনের মধ্যে ৩ দিন পেটে তীব্র ব্যাথা হয়েছিল এবং পেটের ডানদিকে পাজরের নিচে চাপ দিলে ব্যাথা অনুভূত হয়। পিত খলিতে পাথর হলে কী পেট ব্যাথার সাথে সাথে পেটের ডানদিকে পাজরের নিচে চাপ দিলে ব্যাথা অনুভূত হয়? এবং কেন এই চাপ দিলে ব্যাথা অনুভূত হয়?? অভিজ্ঞ ডাক্তারের পরামর্শ চাচ্ছি।\”\n”</p> <p>"Translation of the text: I have gallstones in my gallbladder. This has been happening for the last 2.5 years. The stone size is 4 cm. I did not have any severe pain for a long time, but within the last 8 days, I had severe stomach pain on 3 days, and I feel pain when I press under the right rib on the right side of my abdomen. If there are gallstones, does abdominal pain occur along with pain when pressing under the right rib on the right side? And why do I feel pain when I press there?? I want advice from an experienced doctor.\n”</p> <p>"Answer: {\”Age\":[], ”</p> <p>”\”Symptom\”:[\”ব্যাথা\”,\”পেটে তীব্র ব্যাথা\”,\”চাপ দিলে ব্যাথা\”,\”পেট ব্যাথার\”], ”</p> <p>”\”Medicine\”:[], ”</p> <p>”\”Health_Condition\”:[\”পিত খলিতে পাথর\”], ”</p> <p>”\”Specialist\”:[], ”</p> <p>”\”Medical_Procedure\”:[]\n\n”</p>

Figure 7: A translation-based prompt for Bangla biomedical NER, where the yellow-highlighted text shows the English translation of the original Bangla clinical text used for entity extraction.

Prompt	Bangla	Basque	Spanish	English
<i>bio_few</i>	0.054	0.104	0.019	0.120
<i>json_few</i>	0.345	0.312	0.445	0.527

Table 1: Biomedical NER Performance across output schemas with Llama3-8B. In both cases, we apply few-shot prompting.

every token in the exact order, with even one missing/extra token breaking the whole alignment. Span-based extraction is easier because the model can simply copy the entity text spans and return them as clean JSON, which matches how LLMs naturally answer.

4.2 Prompt Variation Analysis

Table 2 shows that QA-style prompting performs best across all languages. It even outperforms detailed label explanations (*expl*) and translation of the input to English (*trans*), despite using only simple WH-questions. The strong performance of QA-style prompting is likely because modern LLMs are instruction tuned on many question-answer pairs. Prior work reformulating NER as machine reading comprehension similarly finds that QA-style formulations improve NER, attributing gains

Prompt	Bangla	Basque	Spanish	English
<i>json_zero</i>	0.291	0.241	0.403	0.422
<i>json_few_stmt</i>	0.345	0.312	0.496	0.527
<i>json_few_qa</i>	0.458	0.503	0.528	0.541
<i>json_few_trans</i>	0.347	0.441	0.494	0.313
<i>json_few_expl</i>	0.417	0.386	0.477	0.304

Table 2: Prompting Strategy Comparison for Multilingual Biomedical NER Using Meta-Llama-3-8B

to query conditioning and (when available) semantically informative queries that encode entity-type knowledge (Li et al., 2020).

Prompting strategy has a stronger effect on low-resource languages than on high-resource languages: moving from zero-shot to question-style prompting improves Bangla by 57.4% and Basque by 108.7%, while Spanish improves by 21.8% and English by only 28.2%.

4.3 Entity Type Analysis

Table 3 and Table 4 break down performance of prompting strategies by named entity types for Bangla and Basque, respectively. In both languages, question-style prompting is best for most types, with the largest absolute gains for MEDICAL_PROCEDURE in Bangla (37.8%) and H-PROFESSIONAL in Basque (37.4%), indicating that such prompting is especially helpful for rarer or harder-to-extract categories.

Occasionally, explanation or translation-based prompting outperforms question-based prompting. Explanation-based prompting helps Bangla SYMPTOMS and translation-based prompting helps Bangla SPECIALISTS. But the gains are modest over question-based prompting, and translation-based prompting fails badly for the AGE category.

4.4 Analysis Under the Best Prompt

Table 5 compares different LLMs under the best-performing (question-style) prompt. Overall, Llama-3-8B performs best on Bangla (F1 = 0.494) and Spanish (F1 = 0.529), while Qwen3-8B achieves the highest score on English (F1 = 0.541) and Basque (F1 = 0.550). Aya performs worst on all languages. This suggests that Qwen may have stronger English and Basque knowledge from pre-training, while Llama may have stronger Spanish and Bangla knowledge, and both Qwen and Llama likely have better multilingual and/or biomedical domain knowledge than Aya. These results reveal the importance of LLM selection when working with low-resource languages.

	Age	Symptom	Medicine	Health_Condition	Specialist	Medical_Procedure
<i>json_zero</i>	0.421	0.067	0.400	0.400	0.514	0.133
<i>json_few_stmt</i>	0.200	0.198	0.490	0.400	0.621	0.167
<i>json_few_qa</i>	0.471	0.215	0.600	0.533	0.600	0.546
<i>json_few_trans</i>	0.010	0.196	0.560	0.333	0.625	0.222
<i>json_few_expl</i>	0.235	0.253	0.600	0.519	0.500	0.364

Table 3: Entity-type-wise F1 scores for Bangla biomedical NER under five prompting strategies. Bold indicates the best F1 per entity type.

	Disorder	Patient	H-Professional
<i>json_zero</i>	0.215	0.368	0.140
<i>json_few_stmt</i>	0.314	0.441	0.183
<i>json_few_qa</i>	0.427	0.515	0.558
<i>json_few_trans</i>	0.380	0.410	0.320
<i>json_few_expl</i>	0.411	0.409	0.328

Table 4: Entity-type-wise F1 scores for Basque biomedical NER under five prompting strategies. Bold indicates the best F1 per entity type.

Model	Bangla	Basque	Spanish	English
<i>Llama-3-8B</i>	0.494	0.503	0.529	0.457
<i>Qwen3-8B</i>	0.399	0.550	0.514	0.541
<i>Aya</i>	0.229	0.282	0.349	0.318

Table 5: Prompting Strategy Comparison for Multilingual Biomedical NER Using Meta-Llama-3-8B and question-style prompting (*json_few_qa*).

5 Error Analysis

We inspected the errors of the best performing models in Bangla and Basque. Below are the main failure modes, with the operational definitions used in our evaluation. Examples of errors are given in [Appendix D](#) and [Appendix E](#).

Hallucination The gold annotation contains no entities, but the model outputs one or more spans.

All-Missed The gold contains entities, but the model returns an empty JSON. This is the most severe recall failure.

Missed Entities The model extracts some entities but misses others, without adding extra entities.

Extra Entities The model predicts additional entities that are not in the gold but does not miss gold entities.

Type Confusion The predicted entity span matches the gold text, but the assigned type is wrong.

Mixed Errors The model both adds and misses entities. Extra-dominant cases indicate over-extraction, while missed-dominant cases indicate under-extraction.

Boundary Mismatch The model captures the right concept but with different boundaries, becoming both a false positive and a false negative under exact-match scoring.

Figure 8 (top) shows clear differences across prompting methods for Bangla biomedical NER errors. Although we have seen that QA-style prompts perform better on average, here we see that translation-based prompting reduces over-generation: it has the lowest HALLUCINATION (2.6%) and comparatively low BOUNDARY MISMATCH (13%). For span-level precision, question-style prompting yields the largest BOUNDARY MISMATCH (29%), suggesting it identifies the right concept but struggles to copy the exact multiword span. Thus, question-style prompting helps the model extract more entities, while translation-based prompting helps reduce hallucination, extra entities, and boundary mismatch errors.

For Basque (Figure 8, bottom), across all prompts, HALLUCINATION dominates (38–52% of errors), indicating that over-generation is the primary failure mode in Basque. However, the prompts shift the secondary errors: statement-based prompting (*json_few_stmt*) reduces HALLUCINATION (37.5%) compared to the other methods, but it increases ALL-MISSED (14.8%), suggesting a more conservative behavior that sometimes fails to extract anything. In contrast, question-style prompting extracts more aggressively and keeps ALL-MISSED low (5.8%), but this comes with higher EXTRA ENTITIES (8.6%) and substantial BOUNDARY MISMATCH (18.5%).

6 Conclusion and Future Work

We systematically analyzed prompt design choices and output schemas for multilingual biomedical NER, focusing on the low-resource, morphologically rich languages Bangla and Basque while comparing them against the high-resource languages Spanish and English. We find that question-style prompting is the most consistent way to improve

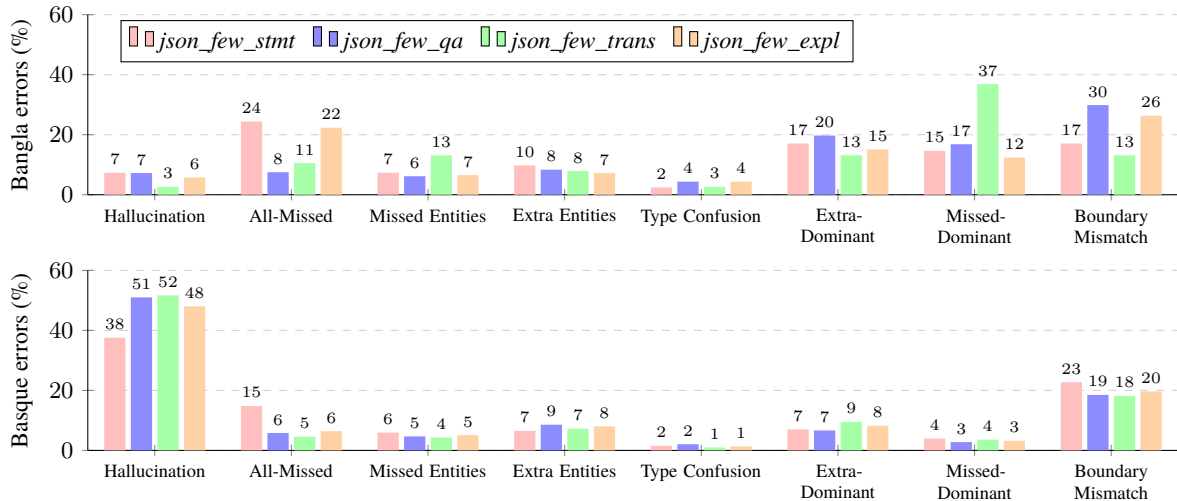


Figure 8: Error type distribution for LLM-extracted entities for Bangla (top) and Basque (bottom).

exact-match entity-level F1 across all languages. We also find that selecting an optimal prompt is more critical in low-resource languages, where performance degrades much more under suboptimal prompting strategies. Our error analysis shows that question-style prompts result in fewer empty outputs, while translation-based prompting can reduce hallucinations and boundary errors in some settings.

We also find that BIO-tagging output is less suitable for LLMs than text-span based output. This contrasts with available datasets: not only biomedical NER, but nearly all Bangla NER datasets are annotated in the BIO format (e.g., Islam et al., 2022; Sazed, 2022; Khan et al., 2023; Mahtab et al., 2025). BIO format is fragile: the model must keep the exact token order and token count, and even one extra or missing token breaks alignment. This becomes worse when tokenization is difficult, as in morphologically rich languages like Bangla and Basque, where suffixes and punctuation frequently change token boundaries. In contrast, text-span-based extraction asks the model to list the entity strings in the sentence, which matches how LLMs naturally respond and is easier to parse and store as JSON.

In the future, we plan to test richer question prompts that include short label definitions and examples (not only WH-questions) to see whether they further reduce boundary mismatches and type confusion. We also want to study better post-processing and evaluation for near-miss spans and explore lightweight adaptation methods that choose the best prompt per language and entity

type. Finally, expanding to more low-resource languages and more medical datasets will help confirm how general these findings are.

Limitations

A key limitation of this study is that we evaluate only three open LLMs around 8B parameters and two datasets, Bangla HealthNER and Basque E3C, so further work is needed to determine whether the findings generalize to other model sizes, closed models, or biomedical text styles. We also do not cover additional low-resource languages because validating prompts, examples, and error cases reliably requires language expert support. In addition, we did not evaluate paid or proprietary GPT models, and our strict exact match scoring can penalize near-correct span boundaries in morphologically rich languages. Finally, since different datasets are used for each language, cross-language comparisons should be interpreted cautiously, as proportional improvements may still reflect dataset-specific characteristics rather than true language differences.

References

- Sayantana Adak, Pauras Mangesh Meher, Paramita Das, and Animesh Mukherjee. 2025. [REVerSum: A multi-staged retrieval-augmented generation method to enhance Wikipedia tail biographies through personal narratives](#). In *Proceedings of the 31st International Conference on Computational Linguistics: Industry Track*, pages 732–750, Abu Dhabi, UAE. Association for Computational Linguistics.
- Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. [Large language](#)

- models are few-shot clinical information extractors. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1998–2022, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Catherine Arnett and Benjamin Bergen. 2025. [Why do language models perform worse for morphologically complex languages?](#) In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6607–6623, Abu Dhabi, UAE. Association for Computational Linguistics.
- Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh Fadaee, Ahmet Üstün, and Sara Hooker. 2024. [Aya 23: Open weight releases to further multilingual progress.](#) *Preprint*, arXiv:2405.15032.
- Dhananjay Ashok and Zachary C. Lipton. 2023. [PromptNER: Prompting For Named Entity Recognition.](#) *arXiv preprint*. ArXiv:2305.15444 [cs].
- Reza Averly and Xia Ning. 2025. [Entity decomposition with filtering: A zero-shot clinical named entity recognition framework.](#) In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2935–2951, Albuquerque, New Mexico. Association for Computational Linguistics.
- Israel Abebe Azime, Atnafu Lambebo Tonja, Tadesse Destaw Belay, Yonas Chanbe, Bontu Fufa Balcha, Negasi Haile Abadi, Henok Biadglign Ademteu, Mulubrhan Abebe Nerea, Debela Desalegn Yadeta, Derartu Dagne Geremew, Assefa Atsbiha Tesfu, Philipp Slusallek, Tamar Solorio, and Dietrich Klakow. 2025. [ProverbEval: Exploring LLM evaluation challenges for low-resource language understanding.](#) In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 6250–6266, Albuquerque, New Mexico. Association for Computational Linguistics.
- Pramit Bhattacharyya and Arnab Bhattacharya. 2025. [BanglaByT5: Byte-level modelling for Bangla.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 5551–5560, Suzhou, China. Association for Computational Linguistics.
- Xiang Cheng, Chengyan Pan, Minjun Zhao, Deyang Li, Fangchao Liu, Xinyu Zhang, Xiao Zhang, and Yong Liu. 2025. [Revisiting chain-of-thought prompting: Zero-shot can be stronger than few-shot.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 13533–13554, Suzhou, China. Association for Computational Linguistics.
- Rochelle Choenni, Dan Garrette, and Ekaterina Shutova. 2023. [Cross-lingual transfer with language-specific subnetworks for low-resource dependency parsing.](#) *Computational Linguistics*, 49(3):613–641.
- Daryna Dementieva, Nikolay Babakov, Amit Ronen, Abinew Ali Ayele, Naqee Rizwan, Florian Schneider, Xintong Wang, Seid Muhie Yimam, Daniil Moskovskiy, Elisei Stakovskii, Eran Kaufman, Ashraf Elnagar, Animesh Mukherjee, and Alexander Panchenko. 2025. [Multilingual and explainable text detoxification with parallel corpora.](#) In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7998–8025, Abu Dhabi, UAE. Association for Computational Linguistics.
- Taisei Enomoto, Hwichan Kim, Zhousi Chen, and Mamoru Komachi. 2025. [A fair comparison without translationese: English vs. target-language instructions for multilingual LLMs.](#) In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 649–670, Albuquerque, New Mexico. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models.](#) *Preprint*, arXiv:2407.21783.
- Md. Arid Hasan, Shudipta Das, Afiyat Anjum, Firoj Alam, Anika Anjum, Avijit Sarker, and Sheak Rashed Haider Noori. 2024. [Zero- and few-shot prompting with LLMs: A comparative study with fine-tuned models for Bangla sentiment analysis.](#) In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 17808–17818, Torino, Italia. ELRA and ICCL.
- Tanvir Islam, Sakila Mahbin Zinat, Shamima Sukhi, Zakir Hossain Zamil, Aynur Nahar, and M. F. Mridha. 2022. [An attention-based medical NER in the Bengali language.](#) In G. Mathur and 1 others, editors, *Proceedings of 2nd International Conference on Artificial Intelligence: Advances and Applications*, Algorithms for Intelligent Systems, pages 131–140. Springer Nature Singapore.
- Alvi Khan, Fida Kamal, Nuzhat Nower, Tasnim Ahmed, Sabbir Ahmed, and Tareque Chowdhury. 2023. [NERvous about my health: Constructing a Bengali medical named entity recognition dataset.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5768–5774, Singapore. Association for Computational Linguistics.
- Somnath Kumar, Vaibhav Balloli, Mercy Ranjit, Kabir Ahuja, Sunayana Sitaram, Kalika Bali, Tanuja Ganu,

- and Akshay Nambi. 2025. [Bridging the language gap: Dynamic learning strategies for improving multilingual performance in LLMs](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 9209–9223, Abu Dhabi, UAE. Association for Computational Linguistics.
- Mingchen Li, Yang Ye, Jeremy Yeung, Huixue Zhou, Huaiyuan Chu, and Rui Zhang. 2023a. [W-procer: Weighted Prototypical Contrastive Learning for Medical Few-Shot Named Entity Recognition](#). *arXiv preprint*. ArXiv:2305.18624 [cs].
- Mingchen Li and Rui Zhang. 2024. [How far is language model from 100% few-shot named entity recognition in medical domain](#). *arXiv preprint arXiv:2307.00186*.
- Xiaoqian Li, Ercong Nie, and Sheng Liang. 2023b. [Crosslingual retrieval augmented in-context learning for Bangla](#). In *Proceedings of the First Workshop on Bangla Language Processing (BLP-2023)*, pages 136–151, Singapore. Association for Computational Linguistics.
- Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020. [A unified MRC framework for named entity recognition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5849–5859, Online. Association for Computational Linguistics.
- Chaoqun Liu, Wenxuan Zhang, Guizhen Chen, Xiaobao Wu, Anh Tuan Luu, Chip Hong Chang, and Lidong Bing. 2023. [Zero-shot text classification via self-supervised tuning](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1743–1761, Toronto, Canada. Association for Computational Linguistics.
- Linlin Liu, Bosheng Ding, Lidong Bing, Shafiq Joty, Luo Si, and Chunyan Miao. 2021. [MulDA: A multilingual data augmentation framework for low-resource cross-lingual NER](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5834–5846, Online. Association for Computational Linguistics.
- Maddalen López de Lacalle, Xabier Saralegi, and Iñaki San Vicente. 2020. [Building a task-oriented dialog system for languages with no training data: the case for Basque](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2796–2802, Marseille, France. European Language Resources Association.
- Bernardo Magnini, Begoña Altuna, Alberto Lavelli, Manuela Speranza, and Roberto Zanolli. 2020. [The E3C project: collection and annotation of a multilingual corpus of clinical cases](#). In *Proceedings of the Seventh Italian Conference on Computational Linguistics (CLiC-it 2020)*, pages 190–196, Bologna, Italy. CEUR Workshop Proceedings.
- Bernardo Magnini, Marco Madeddu, Michele Resta, Roberto Zanolli, Martin Cimmino, Paolo Albano, and Viviana Patti. 2025. [A leaderboard for benchmarking LLMs on Italian](#). In *Proceedings of the Eleventh Italian Conference on Computational Linguistics (CLiC-it 2025)*, pages 636–646, Cagliari, Italy. CEUR Workshop Proceedings.
- Md. Motahar Mahtab, Faisal Ahamed Khan, Md. Ekramul Islam, Md. Shahad Mahmud Chowdhury, Labib Imam Chowdhury, Sadia Afrin, Hazrat Ali, Mohammad Mamun Or Rashid, Nabeel Mohammed, and Mohammad Ruhul Amin. 2025. [BanNERD: A benchmark dataset and context-driven approach for Bangla named entity recognition](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 6807–6828, Albuquerque, New Mexico. Association for Computational Linguistics.
- Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. 2023. [What in-context learning “learns” in-context: Disentangling task recognition and task learning](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8298–8319, Toronto, Canada. Association for Computational Linguistics.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020. [MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7654–7673, Online. Association for Computational Linguistics.
- Trinh Pham, Khoi Le, and Anh Tuan Luu. 2024. [UniB-ridge: A unified approach to cross-lingual transfer learning for low-resource languages](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3168–3184, Bangkok, Thailand. Association for Computational Linguistics.
- Lance Ramshaw and Mitch Marcus. 1995. [Text chunking using transformation-based learning](#). In *Third Workshop on Very Large Corpora*.
- Matias Rojas, Felipe Bravo-Marquez, and Jocelyn Dunstan. 2022. [Simple yet powerful: An overlooked architecture for nested named entity recognition](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2108–2117, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Salim Sazzed. 2022. [BanglaBioMed: A biomedical named-entity annotated corpus for Bangla \(Bengali\)](#). In *Proceedings of the 21st Workshop on Biomedical Language Processing*, pages 323–329, Dublin, Ireland. Association for Computational Linguistics.
- Sheikh Shafayat, H M Quamran Hasan, Minhajur Rahman Chowdhury Mahim, Rifki Afina Putri, James Thorne, and Alice Oh. 2024. [BEnQA: A question answering benchmark for Bengali and English](#). In *Findings of the Association for Computational Linguistics*.

tics: *ACL 2024*, pages 1158–1177, Bangkok, Thailand. Association for Computational Linguistics.

Mohammad Golam Sohrab and Makoto Miwa. 2018. [Deep exhaustive model for nested named entity recognition](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2843–2849, Brussels, Belgium. Association for Computational Linguistics.

IV Styler, William F., Steven Bethard, Sean Finan, Martha Palmer, Sameer Pradhan, Piet C de Groen, Brad Erickson, Timothy Miller, Chen Lin, Guergana Savova, and James Pustejovsky. 2014. [Temporal annotation in the clinical domain](#). *Transactions of the Association for Computational Linguistics*, 2:143–154.

Gorka Urbizu, Iñaki San Vicente, Xabier Saralegi, Rodrigo Agerri, and Aitor Soroa. 2022. [BasqueGLUE: A natural language understanding benchmark for Basque](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1603–1612, Marseille, France. European Language Resources Association.

Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, Guoyin Wang, and Chen Guo. 2025. [GPT-NER: Named entity recognition via large language models](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 4257–4275, Albuquerque, New Mexico. Association for Computational Linguistics.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.

Roberto Zanoli, Alberto Lavelli, Daniel Verdi do Amarante, and Daniele Toti. 2024. [Assessment of the E3C corpus for the recognition of disorders in clinical texts](#). *Natural Language Engineering*, 30(4):851–869.

A BIO Tagging vs Span-Based

The left panel of [Figure 9](#) shows span-based extraction, where entities are directly returned as text spans grouped by type in a JSON structure, while the right panel shows token-level BIO tagging, which requires assigning a label to every token and maintaining strict alignment. This example illustrates how span-based formats provide a simpler and more robust representation for LLM-based NER, especially for morphologically rich languages like Bangla.

Span-level JSON Format	Token-level BIO Format
<pre>{ "ID": 5, "Input": "বেশ কয়েকদিন যাবত গলায় খুব ব্যাথা সেই সাথে কাশি হয়। কাশির সাথে ঘন কফ বের হয় আর বুকে খুব ব্যাথা হয়। জর ও আছে", "Output": { "Symptom": ["গলায় খুব ব্যাথা", "কাশি", "কাশির সাথে ঘন কফ", "বুকে খুব ব্যাথা", "জর"] } }</pre>	<pre>{ "ID": 5, "TOKEN": ["বেশ", "কয়েকদিন", "যাবত", "গলায়", "খুব", "ব্যাথা", "সেই", "সাথে", "কাশি", "হয়", "।", "কাশির", "সাথে", "ঘন", "কফ", "বের হয়", "আর", "বুকে", "খুব", "ব্যাথা", "হয়", "।", "জর", "ও", "আছে", "।"], "NER_TAG": ["O", "O", "O", "B-Symptom", "I-Symptom", "I-Symptom", "O", "O", "B-Symptom", "O", "O", "B-Symptom", "I-Symptom", "I-Symptom", "I-Symptom", "O", "O", "B-Symptom", "I-Symptom", "I-Symptom", "O", "O", "B-Symptom", "O", "O", "O"] }</pre>

Figure 9: Comparison of span-level JSON and token-level BIO formats for Bangla biomedical NER.

B Prompting Behavior and Qualitative Analysis for Bangla

[Table 6](#) illustrates qualitative differences across three prompting strategies for Bangla biomedical NER. In the question-style prompt, the model correctly identifies most entities but misses the HEALTH_CONDITION (হীনফেকশন), showing that QA prompting can still suffer from recall gaps. In contrast, translation-based prompting leads to over-generation: the model incorrectly treats common home remedies such as গরম পানি (hot water) and মধু as MEDICINE, likely due to reasoning in English and broader semantic interpretation. Explanation-based prompting improves recall but introduces type confusion, as the model incorrectly labels the polite address স্যার (“Sir”) as a SPECIALIST. These examples highlight that while extended prompting strategies can help extraction, they also introduce distinct error patterns depending on how semantic cues are framed.

C Prompting Behavior and Qualitative Analysis for Basque

[Table 7](#) presents qualitative differences in model behavior across three prompting strategies for Basque biomedical NER. Under the question-style prompt, the model successfully identifies major disorders but exhibits boundary and normalization issues, such as over-specifying severity or partially mismatching gold entity forms. Translation-based prompting generally improves coverage by leveraging English semantics, but it can introduce label drift, as closely related clinical concepts are normalized differently from the gold labels. Explanation-based prompting further increases re-

E.g.	Prompt Type	Prompt	Llama3-8B	Gold Label
1	json_few_- qa	<p>You are a question-answering assistant that performs medical Named Entity Recognition (NER). For each question, identify ONLY the entities that are explicitly present in the provided text and answer STRICTLY in JSON.</p> <p>Question: Which entities (Age, Symptom, Medicine, Health_Condition, Specialist, Medical_Procedure) are present in this text?</p> <p>Example 1: এলাজির সমস্যার কারণে রোদে গেলে গা চিটমিট করে, মাথার ভিতরে কিলবিল করে। সকালে ঘুম থেকে উঠলে অনবরত হাঁচি হয়। কখনো নিয়মিত কোনো এলাজির ওষুধ খাইনি। এক্ষেত্রে আমি কি করতে পারি? এলাজির কারণে অনেক দৈনন্দিন কাজ করতে পারি না। Answer: Symptom – রোদে গেলে গা চিটমিট করে মাথার ভিতরে কিলবিল করে হাঁচি; Medicine – এলাজির ওষুধ. Health_Condition – এলাজির.</p> <p>Now extract all entities for this text: ধন্যবাদ আপনার প্রশ্নের জন্য। ফটিয়ে দিবেন না। ইনফেকশন হয়ে যাবে। আপাতত এলাজি উদ্বোধকারী খাবার, ধূলাবালি এড়িয়ে চলুন। এলাজির খান। না দেখে সমাধান দেয়া যাচ্ছে না। সরাসরি চর্মরোগ বিশেষজ্ঞ দেখিয়ে পরামর্শ নিতে হবে ধন্যবাদ</p>	Medicine: এলাজি Specialist: চর্মরোগ বিশেষজ্ঞ	Medicine: এলাজি Health_Condition: ইনফেকশন Specialist: চর্মরোগ বিশেষজ্ঞ
2	json_few_- trans	<p>You are a NER assistant that performs medical Named Entity Recognition (NER). For each question, identify ONLY the entities that are explicitly present in the provided text and answer STRICTLY in JSON.</p> <p>Example 1: এলাজির সমস্যার কারণে রোদে গেলে গা চিটমিট করে, মাথার ভিতরে কিলবিল করে। সকালে ঘুম থেকে উঠলে অনবরত হাঁচি হয়। কখনো নিয়মিত কোনো এলাজির ওষুধ খাইনি। এক্ষেত্রে আমি কি করতে পারি? এলাজির কারণে অনেক দৈনন্দিন কাজ করতে পারি না।</p> <p>English translation of the Example: Because of allergy problems, when I go out in the sun my body tingles, and I feel a crawling sensation inside my head. In the morning after waking up I sneeze continuously. I have never taken any allergy medicine regularly. In this case, what can I do? Because of this allergy, I cannot do many daily tasks. Now find out the entities (Age, Symptom, Medicine, Health_Condition, Specialist, Medical_Procedure) are present in this text.</p> <p>Answer: Symptom – রোদে গেলে গা চিটমিট করে মাথার ভিতরে কিলবিল করে হাঁচি; Medicine – এলাজির ওষুধ. Health_Condition – এলাজির.</p> <p>Now extract all entities for this text: আপনার বাচাকে আদা, মধু, গরম পানি খাওয়ান। সিরাপ এমব্রোজ খাওয়াতে পারেন আধা চামচ করে তিন বার একজন শিশু বিশেষজ্ঞ কে দেখিয়ে নিন।</p>	Medicine: আদা গরম পানি মধু সিরাপ এমব্রোজ ; Specialist: শিশু বিশেষজ্ঞ	Medicine: সিরাপ এমব্রোজ ; Specialist: শিশু বিশেষজ্ঞ
3	json_few_- expl	<p>You are a NER assistant that performs medical Named Entity Recognition (NER). For each question, identify ONLY the entities that are explicitly present in the provided text and answer STRICTLY in JSON.</p> <p>Example 1: এলাজির সমস্যার কারণে রোদে গেলে গা চিটমিট করে, মাথার ভিতরে কিলবিল করে। সকালে ঘুম থেকে উঠলে অনবরত হাঁচি হয়। কখনো নিয়মিত কোনো এলাজির ওষুধ খাইনি। এক্ষেত্রে আমি কি করতে পারি? এলাজির কারণে অনেক দৈনন্দিন কাজ করতে পারি না।</p> <p>Explanation of this Example: The Bangla text clearly describes an allergy problem, which is a Health_Condition entity (“এলাজির”). It also describes allergy-related symptoms, such as skin discomfort in the sun, crawling sensation in the head, and continuous sneezing (Symptom). It mentions “এলাজির ওষুধ” only in a generic way, but it still refers to medicine use (Medicine). No age, specialist, or procedure is stated. Now read the Bangla text and check the Answer, then identify which entities (Age, Symptom, Medicine, Health_Condition, Specialist, Medical_Procedure) are present in the text.</p> <p>Answer: Symptom – রোদে গেলে গা চিটমিট করে মাথার ভিতরে কিলবিল করে হাঁচি; Medicine – এলাজির ওষুধ. Health_Condition – এলাজির.</p> <p>Example 2: রোগীর বয়স ৫৫ বছর। তিনি ডায়াবেটিস ও উচ্চ রক্তচাপের রোগী এবং মেটফরমিন ও লোসারটান খাচ্ছেন। Answer: Age – ৫৫ বছর. Health_Condition – ডায়াবেটিস উচ্চ রক্তচাপ.</p> <p>Now extract all entities for this text: আমার বয়স ১৩ আমার পচ পায়ের হার এর বেথা। আমি কি ঔষধ খাব বলেন স্যার। Calcium BP 500mg কিনছি খাব কি স্যার।</p>	Age: ১৩ বছর ; Symptom: পচ পায়ের হার এর বেথা ; Medicine: Calcium BP 500mg ; Specialist: স্যার	Age: ১৩ ; Symptom: পচ পায়ের হার এর বেথা ; Medicine: Calcium BP 500mg

Table 6: Example prompt templates for Bangla biomedical NER. We show three prompt variants (question-style, translation-based, and explanation-based) with one in-context demonstration and one target Bangla input. For each example, we report the Meta-Llama-3-8B output and the corresponding gold labels for comparison.

E.g.	Prompt Type	Prompt	Llama3-8B	Gold Label
1	json_few_- qa	<p>You are an assistant that performs medical Named Entity Recognition (NER) in Basque. For each input text, identify ONLY the entities that are explicitly present in the provided text and answer STRICTLY in JSON.</p> <p>Question: Which entities (Disorder, Patient, H-Professional) are present in this text?</p> <p>Example 1: Text: “-Eskuineko besapea:aurkikuntza patologikorik ez.” Answer: {'Disorder': ['aurkikuntza'], 'Patient': [], 'H-Professional': []}</p> <p>Now extract all entities for this text: <i>Estenosi aortiko larriari bihotz-gutxiegitasuna eta sinkopea gehitu zaizkionez, ebakuntza egiteko baldintzak betetzen ditu, eta Bihotz-hodietako Kirurgiako Zerbitzuarekin harremanetan jarri gara.</i></p>	<p>Disorder: bihotz-gutxiegitasuna estenosi aortiko larria sinkopea ; H-Professional: Bihotz-hodietako Kirurgiako Zerbitzuarekin</p>	<p>Disorder: Estenosi aortiko bihotz-gutxiegitasuna sinkopea]</p>
2	json_few_- trans	<p>You are an assistant that performs medical Named Entity Recognition (NER) in Basque. For each input text, identify ONLY the entities that are explicitly present in the provided text and answer STRICTLY in JSON.</p> <p>Example 1: Text: “-Eskuineko besapea:aurkikuntza patologikorik ez.” English translation of the text: “-Right armpit: no pathological findings.” Answer: {'Disorder': ['aurkikuntza'], 'Patient': [], 'H-Professional': []}</p> <p>Now extract all entities for this text: <i>Txankro sifilitikoa/Sifilis primarioa.</i></p>	<p>Disorder: Sifilis primarioa Txankro sifilitikoa</p>	<p>Disorder: Sifilis Txankro sifilitikoa</p>
3	json_few_- expl	<p>You are an assistant that performs medical Named Entity Recognition (NER) in Basque. For each input text, identify ONLY the entities that are explicitly present in the provided text and answer STRICTLY in JSON.</p> <p>Example 1: Text: “-Eskuineko besapea:aurkikuntza patologikorik ez.” Description of text: A clinical examination of a right armpit with no health disorders or pathological findings identified. You have to find out which entities (Disorder, Patient, H-Professional) are present in the text. Answer: {'Disorder': ['aurkikuntza'], 'Patient': [], 'H-Professional': []}</p> <p>Now extract all entities for this text: <i>Edemarik ez; ez dago bena sakonetako tronbosi-zeinurik.</i></p>	<p>Disorder: tronbosi-zeinurik</p>	<p>Disorder: Edemarik bena sakonetako tronbosi-zeinurik</p>

Table 7: Example prompt templates for Basque biomedical NER. We show three prompt variants (question-style, translation-based, and explanation-based) with one in-context demonstration and one target Basque input. For each example, we report the Meta-Llama-3-8B output and the corresponding gold labels for comparison.

Error Type	Input Example	Gold Output	Pred Output
Hallucinations	ধন্যবাদ আপনাকে প্রশ্নের জন্য। জ্বর ৯৯ হলে ওষুধ খাওয়াবেন। নয়ত দরকার নেই। ধন্যবাদ	Empty	Symptom: জ্বর
All-Missed	আমি। আমার বয়স ২৭ বছর। সাম্প্রতিক সময়ে আমি আমার শরীরে অস্বস্তি অনুভব করছি, আগের থেকে শরীরের ওজন বৃদ্ধি পেয়েছে, পেটে প্র...	Age: ২৭ বছর ; Symptom: একাধারে মাটিতে বসে থাকলে মাজার ...। শরীরের ওজন বৃদ্ধি হাটলেই মনে হয় ক্লান্ত হয়ে	Empty
Missed Entities	total cholesterol 200 mg / dl serum triglycerides 553 mg / dl...how to lower triglycerides? am i prone to heart disease?	Symptom: HDL cholesterol 25 mg LDL cholesterol 93 mg cholesterol...। serum triglycerides 553 mg ; Health_Condition: heart disease	Health_Condition: heart disease
Extra Entities	আপনাকে ধন্যবাদ প্রশ্ন করার জন্য। আপনি গরম পানির খোয়া নাকে নিয়...রাতে একটা করে খান। একজন মেডিসিন বিশেষজ্ঞ পরামর্শ নিন। ধন্যবাদ	Medicine: এলাট্রিল ; Specialist: মেডিসিন বিশেষজ্ঞ	Medicine: এলাট্রিল গরম পানির খোয়া ; Specialist: মেডিসিন বিশেষজ্ঞ
Type Confusion	আমার স্ত্রী বয়স ২৪ ও ১৪ দিনের প্রেগন্যান্ট অবস্থা MM Kit খায় ...আল্ট্রাসোনোগ্রাফি করতে চাই, দয়া করে আমাকে পরিষ্কার নাম টা লিখে দিবেন	Age: ২৪ ; Symptom: bleeding আনে বেশি যাচ্ছে ; Medicine: MM Kit ...; Health_Condition: মাসিক ১৪ দিনের প্রেগন্যান্ট	Age: ১৪ দিন ২৪ ; Symptom: bleeding ; Medicine: MM Kit ; Health_Condition: প্রেগন্যান্ট অবস্থা ; Medical_Procedure: আল্ট্রাসোনোগ্রাফি
Mixed Errors	Thank you for your question. Your serum Triglyceride is sligh...with your reports. Avoid fatty food, carbohydrate. Thank you.	Symptom: serum Triglyceride is slightly raised	Medicine: drug ; Specialist: doctor
Boundary Mismatch	বেশ কয়েকদিন যাবত গলায় খুব ব্যাথা সেই সাথে কাশি হয়। কাশির সাথে ঘন কফ বেরহয় আর বুকে খুব ব্যাথা হয়। জর ও আছে।	Symptom: কাশি কাশির সাথে ঘন কফ গলায় খুব ব্যাথা জর বুকে খুব ব্যাথা	Symptom: কাশি খুব ব্যাথা ঘন কফ জর বুকে খুব ব্যাথা

Table 8: Representative QA-prompt error examples for Bangla biomedical NER (Llama3.1-8B). Each row shows one real instance of a major error category under exact-match span scoring.

Error Type	Input Example	Gold Output	Pred Output
Hallucinations	<i>Burua eta lepoa.</i>	<i>Empty</i>	Disorder: Burua eta lepoa
All-Missed	<i>Ez dago aurkikuntza patologikorik.</i>	Disorder: aurkikuntza	<i>Empty</i>
Missed Entities	<i>Estenosi aortiko larria eta bihotz-gutxiegitasuna.</i>	Disorder: Estenosi aortiko larria bihotz-gutxiegitasuna	Disorder: bihotz-gutxiegitasuna
Extra Entities	<i>Sabela.</i>	<i>Empty</i>	Patient: Sabela
Type Confusion	<i>Biguna eta zanpagarria.</i>	<i>Empty</i>	Disorder: Biguna Disorder: zanpagarria
Mixed Errors	<i>Biriketako murmurio normala.</i>	<i>Empty</i>	Disorder: Biriketako murmurio
Boundary Mismatch	<i>Ez dago aurkikuntza patologikorik.</i>	Disorder: aurkikuntza patologikorik	Disorder: aurkikuntza

Table 9: Representative QA-prompt error examples for Basque biomedical NER (Meta-Llama-3-8B). Each row illustrates a distinct error category under exact-match span evaluation.

call by encouraging semantic inference, yet it also leads to over-generation, with the model extracting implicit or negated findings (e.g., absence of edema or thrombosis) as entities. Overall, these examples demonstrate that richer prompting strategies can enhance extraction in Basque but also introduce systematic errors related to semantic inference, normalization, and negation handling.

D Bangla Errors (QA Prompt)

Table 8 presents representative Bangla error cases for the QA-style prompt. A recurring pattern is that the model often identifies the correct medical concept but fails to reproduce the exact mention boundaries, especially for multiword symptoms and colloquial expressions, resulting in **Boundary Mismatch**. We also observe **Hallucinations** when generic medical advice or common words are treated as symptoms/conditions even when the gold annotation is empty. In contrast, **All-Missed** cases reflect severe recall failures where the model returns an empty JSON despite clear entity cues (e.g., explicit age or specialist mentions). Finally, **Type Confusion** commonly arises between semantically adjacent categories (e.g., tests/procedures vs. conditions, or informal medication names vs. remedies), indicating that the model recognizes the span but struggles with fine-grained label assignment.

E Basque Errors (QA Prompt)

Table 9 shows representative Basque error cases for the same QA-style prompt. Across prompts, the dominant failure mode is **Hallucinations**, suggesting that the model is prone to over-generation in Basque medical text: short phrases and anatomical references are frequently over-labeled as medical entities even when the gold annotation contains none. We also see **Boundary Mismatch** in cases where the model selects a shorter head noun

rather than the full descriptive mention, which is penalized under exact-match evaluation. Compared with Bangla, Basque exhibits more frequent over-extraction and fewer purely recall-only failures, consistent with the overall error distribution where FP-heavy errors occupy a larger fraction of error-only cases.