

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

Anonymous ACL submission

Abstract

Despite the strong multilingual abilities of modern LLMs, biomedical information extraction remains inconsistent for low-resource, morphologically rich languages such as Bangla/Bengali and Basque. Prior investigations of prompt design and output schemas have been focused on high-resource settings. To bridge this gap, we systematically evaluate biomedical named entity recognition (NER) with open LLMs under multiple prompting settings. We find that span-based extraction is much 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 low-resource languages than on high-resource languages (e.g., +57% for Bangla and +109% for Basque, but only +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 low-resource languages. We release our code at <https://github.com/<anonymized>>.

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; Averly and Ning, 2025). 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). Creating these labels requires medical expertise, which makes annotation slow, expensive, and hard to scale. Tra-

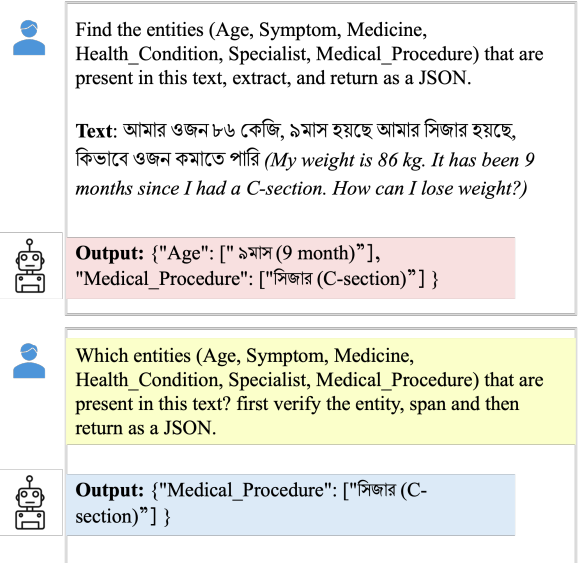


Figure 1: A statement-based prompt for biomedical NER in Bangla results in an incorrect extraction (red), while a question-based prompt results in a correct extraction (blue).

ditional fine-tuning for NER also needs many labeled examples. To reduce this dependence, prior work has explored 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, 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 study of how different prompting strategies behave in Bangla/Bengali and Basque, two morphologically rich languages (Bhattacharyya and Bhattacharya, 2025; López de Lacalle et al., 2020) with complex grammar that pose challenging low-resource settings. To test whether prompt-

ing effects generalize across languages, we also evaluate Spanish and English as high-resource comparisons, which lets us separate shared prompting behavior from language-specific effects under the same multilingual models. We compare two output representations, Beginning–Inside–Outside (BIO) tagging (Chen et al., 2022) and span-based extraction, and provide practical guidance for converting BIO outputs into span format. Finally, we analyze which prompting strategies help reduce common extraction errors, including entity hallucination, over-generation, and boundary mistakes. We summarize these goals in three research questions:

- Which prompting strategy most reliably improves biomedical NER for low-resource languages?
- Does the span-based format beat BIO tagging format, and why?
- Which errors occur most often across prompts, and which prompting choices reduce hallucination, over-generation, and boundary errors?

We show that LLMs struggle generating outputs in BIO format and that span-based extraction is a much better option. We also find that formulating NER as a question-style task (see Figure 1) obtains the best results, providing major improvements for low-resource languages like Bangla and Basque. Our outcomes can benefit biomedical NER not only in other low-resource languages but also in other multilingual settings. Since multilingual models share parameters across languages, methods designed for low-resource settings often strengthen the shared multilingual backbone (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, BanglishBERT, 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. Zanolli et al. (2024b) 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.

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 in English but indicate the target language. For example, “You are a biomedical NER assistant that performs Named Entity Recognition in Bangla.” 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 Ben-

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:     end while
10:     $e \leftarrow i - 1$ 
11:     $\text{span} \leftarrow \text{DETOK}(x_{s:e})$ 
12:     $\text{span} \leftarrow \text{CLEAN}(\text{span})$ 
13:    if  $\text{span} \neq \emptyset$  then
14:       $\mathcal{J}(t) \leftarrow \mathcal{J}(t) \parallel [\text{span}]$ 
15:    end if
16:  else
17:     $i \leftarrow i + 1$ 
18:  end if
19: end while
20: return  $\mathcal{J}$ 
```

Span-based format (json) means the model does not output a tag for every token. Instead, it directly outputs the entity mentions as text spans from the sentence, usually grouped by entity type in a JSON-like structure. Many prior works adopt span-based NER models (Sohrab and Miwa, 2018; Rojas et al., 2022). In contrast, our approach predicts BIO tags and then performs BIO decoding (chunk extraction) to convert token-level predictions into span strings. Formally, BIO decoding 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 also available at: <https://github.com/<anonymized>>.

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.

```
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: {\n\"Age\": [], \"
\"Symptom\": [\nপ্রসাবের রাস্তায়
জ্বালাপোড়া\", \"ব্যাথা\", \"প্রসাবে মাঝে মাঝে গন্ধ\", \"
\"Medicine\": [], \"
\"Health_Condition\": [], \"
\"Specialist\": [], \"
\"Medical_Procedure\": []]\n\n")
```

Figure 4: Illustration of 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.

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. We designed the examples in a way that each entity type appears at least once, and the set includes both short and long examples.

3.4 Prompt Variations

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 4.

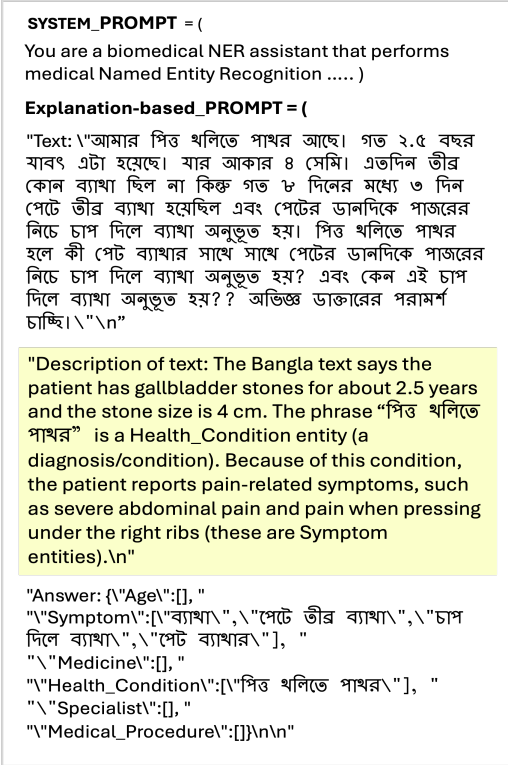


Figure 5: Illustration of an explanation-based prompt for Bangla biomedical NER, where the yellow-highlighted text provides a natural-language explanation of the clinical context and entity semantics.

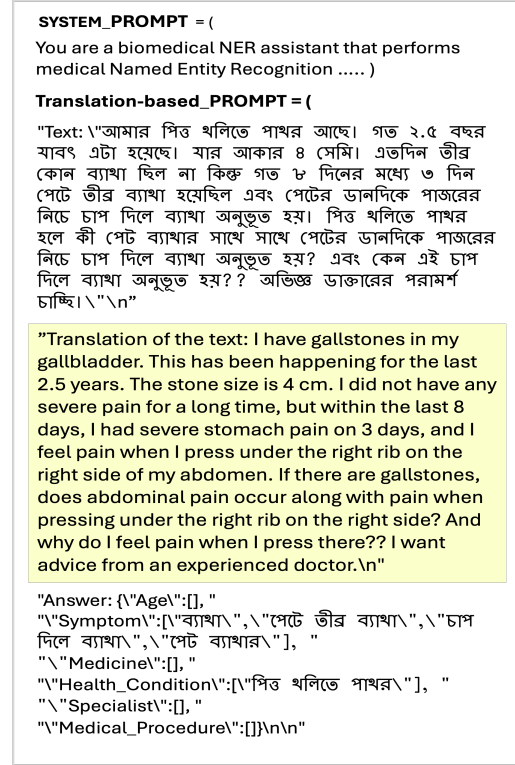


Figure 6: Illustration of a translation-based prompt for Bangla biomedical NER, where the yellow-highlighted text shows the English translation of the original Bangla clinical text used for entity extraction.

Explanation-based prompting (expl) is a prompting strategy (Figure 5) 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 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 6) 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** (Team, 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

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.

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

4.1 Output Schema Comparison: BIO Tagging vs Span-Based

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* yields consistent absolute F1 gains: +0.291 (Bangla; 0.054 \rightarrow 0.345), +0.208 (Basque; 0.104 \rightarrow 0.312), +0.426 (Spanish; 0.019 \rightarrow 0.445), and +0.407 (English; 0.120 \rightarrow 0.527). BIO tagging is likely harder for LLMs because the model must label every token in the exact order, with even one missing/extra token (or a tokenization mismatch in Bangla/Basque) 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 Prompting Strategy Comparison Across Languages

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 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

Prompt	Bangla	Basque	Spanish	English
<i>json_zero</i>	0.291	0.241	0.403	0.422
<i>json_few</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

low-resource languages than on high-resource languages. For example, 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 Biomedical Entity-type-wise Comparison

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 Model Comparison 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), while Qwen3-8B achieves the highest score on English (F1 = 0.541) and Basque (F1 = 0.550); Llama-3-8B performs best on Spanish (F1 = 0.529); Aya performs worse on all languages. This suggests that Qwen may have stronger English and Basque knowledge from pretraining, 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 im-

	Age	Symptom	Medicine	Health_Condition	Specialist	Medical_Procedure
<i>json_zero</i>	0.4211	0.0674	0.4000	0.4000	0.5143	0.1333
<i>json_few</i>	0.2000	0.1978	0.4898	0.4000	0.6207	0.1667
<i>json_few_qa</i>	0.4706	0.2151	0.6000	0.5333	0.6000	0.5455
<i>json_few_trans</i>	0.0099	0.1957	0.5600	0.3333	0.6250	0.2222
<i>json_few_expl</i>	0.2353	0.2526	0.6000	0.5185	0.5000	0.3636

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>	0.2153	0.3678	0.1401
<i>json_few</i>	0.3137	0.4413	0.1833
<i>json_few_qa</i>	0.4268	0.5151	0.5576
<i>json_few_trans</i>	0.3795	0.4100	0.3195
<i>json_few_expl</i>	0.4108	0.4093	0.3278

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*).

portance of LLM selection when working with low-resource languages.

5 Error Analysis

5.1 Error Categories

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.

5.2 Error Analysis Explanation

Figure 7 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), across all prompts, HALLUCINATION dominates the errors (about 38–52%), indicating that over-generation is the primary failure mode in Basque. However, the prompts shift the secondary errors: simple statement-based prompting (*json_few*) noticeably 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

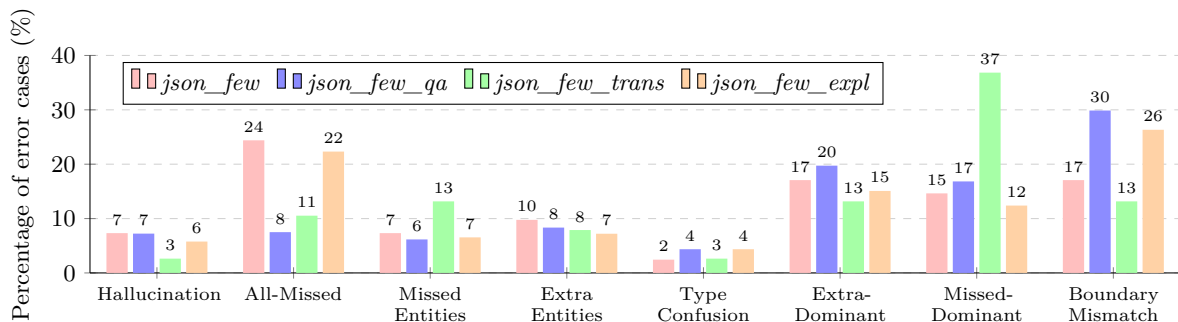


Figure 7: Error type distribution for LLM-extracted entities for Bangla.

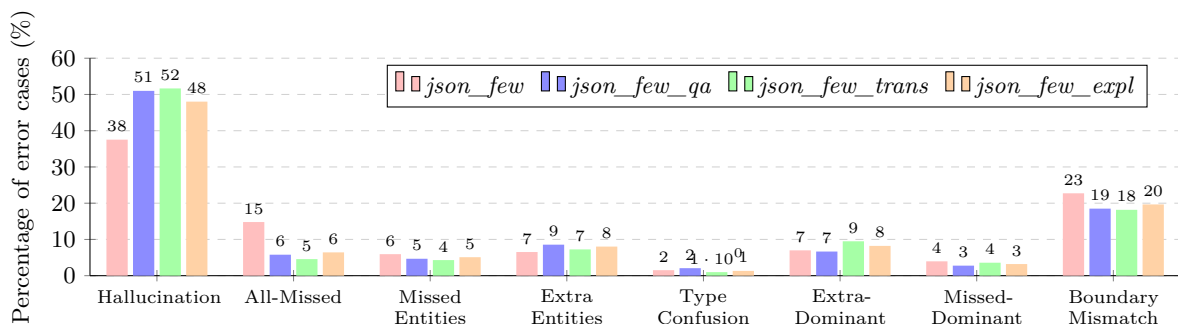


Figure 8: Error type distribution for LLM-extracted entities for Basque.

with higher EXTRA ENTITIES (8.6%) and substantial BOUNDARY MISMATCH (18.5%).

6 Conclusion and Future Work

We systematically evaluated prompt design choices and output schemas for multilingual biomedical NER, focusing on the low-resource, morphologically rich languages Bangla and Basque, while comparing against the high-resource languages Spanish and English. Overall, question-style prompting is the most consistent way to improve exact-match entity-level F1 across languages. Our error analysis also shows clear trade-offs: question-style prompts usually help recall (fewer empty outputs), while translation-based prompting can make the model more conservative and reduce hallucinations and boundary errors in some settings. Together, these results provide practical guidance for building more robust biomedical NER pipelines in low-resource languages without fine-tuning.

We also find that BIO-formatted output is less suitable for LLMs than text-span based outputs. 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 tagging requires an LLM to assign a correct label to every token and maintain strict consistency. In practice, this 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 simply 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.

515 Limitations

516 A key limitation of this study is that we eval-
517 uate only three open LLMs around 8B param-
518 eters and two datasets, Bangla HealthNER
519 and Basque E3C, so further work is needed to
520 determine whether the findings generalize to
521 other model sizes, closed models, or biomed-
522 ical text styles. We also do not cover addi-
523 tional low-resource languages because validat-
524 ing prompts, examples, and error cases reliably
525 requires language expert support. In addition,
526 we did not evaluate paid or proprietary GPT
527 models, and our strict exact match scoring can
528 penalize near-correct span boundaries in mor-
529 phologically rich languages.

530 References

531 Sayantan Adak, Pauras Mangesh Meher, Paramita
532 Das, and Animesh Mukherjee. 2025. [REVer-](#)
533 [Sum: A multi-staged retrieval-augmented gen-](#)
534 [eration method to enhance Wikipedia tail bi-](#)
535 [ographies through personal narratives](#). In *Pro-*
536 *ceedings of the 31st International Conference*
537 *on Computational Linguistics: Industry Track*,
538 pages 732–750, Abu Dhabi, UAE. Association
539 for Computational Linguistics.

540 Monica Agrawal, Stefan Hegselmann, Hunter Lang,
541 Yoon Kim, and David Sontag. 2022. [Large lan-](#)
542 [guage models are few-shot clinical information](#)
543 [extractors](#). In *Proceedings of the 2022 Con-*
544 *ference on Empirical Methods in Natural Lan-*
545 *guage Processing*, pages 1998–2022, Abu Dhabi,
546 United Arab Emirates. Association for Compu-
547 tational Linguistics.

548 Viraat Aryabumi, John Dang, Dwarak Talupuru,
549 Saurabh Dash, David Cairuz, Hangyu Lin,
550 Bharat Venkitesh, Madeline Smith, Yi Chern
551 Tan, Kelly Marchisio, Max Bartolo, Sebastian
552 Ruder, Acyr Locatelli, Julia Kreutzer, Nick
553 Frosst, Aidan Gomez, Phil Blunsom, Marzieh
554 Fadaee, Ahmet Üstün, and Sara Hooker. 2024.
555 [Aya 23: Open weight releases to further multi-](#)
556 [lingual progress](#). *Preprint*, arXiv:2405.15032.

557 Dhananjay Ashok and Zachary C. Lipton. 2023.
558 [PromptNER: Prompting For Named Entity](#)
559 [Recognition](#). *arXiv preprint*. ArXiv:2305.15444
560 [cs].

561 Reza Averly and Xia Ning. 2025. [Entity decom-](#)
562 [position with filtering: A zero-shot clinical named](#)
563 [entity recognition framework](#). In *Proceedings of*
564 *the 2025 Conference of the Nations of the Amer-*
565 *icas Chapter of the Association for Computa-*
566 *tional Linguistics: Human Language Technolo-*
567 *gies (Volume 1: Long Papers)*, pages 2935–2951,

Albuquerque, New Mexico. Association for Com-
putational Linguistics.

Israel Abebe Azime, Atnafu Lambebo Tonja,
Tadesse Destaw Belay, Yonas Chanie,
Bontu Fufa Balcha, Negasi Haile Abadi,
Henok Biadgign Ademtew, 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 under-](#)
[standing](#). 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 Com-
putational Linguistics.

Jiawei Chen, Qing Liu, Hongyu Lin, Xianpei Han,
and Le Sun. 2022. [Few-shot named entity recog-](#)
[nition with self-describing networks](#). In *Proceed-*
ings of the 60th Annual Meeting of the Associ-
ation for Computational Linguistics (Volume 1:
Long Papers), pages 5711–5722, Dublin, Ireland.
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 Ro-
nen, Abinew Ali Ayele, Naqee Rizwan, Florian
Schneider, Xintong Wang, Seid Muhie Yimam,
Daniil Moskovskiy, Elisei Stakovskii, Eran Kauf-
man, Ashraf Elnagar, Animesh Mukherjee, and
Alexander Panchenko. 2025. [Multilingual and](#)
[explainable text detoxification with parallel cor-](#)
[pora](#). In *Proceedings of the 31st International*
Conference on Computational Linguistics, pages
7998–8025, Abu Dhabi, UAE. Association for
Computational Linguistics.

Aaron Grattafiori and 1 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](#)

568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624

625	few-shot prompting with LLMs: A comparative study with fine-tuned models for Bangla sentiment analysis. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 17808–17818, Torino, Italia. ELRA and ICCL.	683
626		684
627		685
628		686
629		687
630		
631		
632	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, <i>Proceedings of 2nd International Conference on Artificial Intelligence: Advances and Applications</i> , Algorithms for Intelligent Systems, pages 131–140. Springer Nature Singapore.	688
633		689
634		690
635		691
636		692
637		693
638		694
639		695
640		696
641	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 <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5768–5774, Singapore. Association for Computational Linguistics.	697
642		698
643		699
644		700
645		701
646		702
647		703
648		704
649	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 <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 9209–9223, Abu Dhabi, UAE. Association for Computational Linguistics.	705
650		706
651		707
652		708
653		709
654		710
655		711
656		712
657		713
658	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 . <i>arXiv preprint</i> . ArXiv:2305.18624 [cs].	714
659		715
660		716
661		717
662		718
663		719
664	Mingchen Li and Rui Zhang. 2024. How far is language model from 100% few-shot named entity recognition in medical domain . <i>arXiv preprint arXiv:2307.00186</i> .	720
665		721
666		722
667		723
668	Xiaoqian Li, Ercong Nie, and Sheng Liang. 2023b. Crosslingual retrieval augmented in-context learning for Bangla . In <i>Proceedings of the First Workshop on Bangla Language Processing (BLP-2023)</i> , pages 136–151, Singapore. Association for Computational Linguistics.	724
669		725
670		726
671		727
672		728
673		729
674	Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020. A unified MRC framework for named entity recognition . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5849–5859, Online. Association for Computational Linguistics.	730
675		731
676		732
677		733
678		734
679		735
680		736
681	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 <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1743–1761, Toronto, Canada. Association for Computational Linguistics.	737
682		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

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 Sazed. 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: 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 Finnan, 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.

Qwen Team. 2025. *Qwen3 technical report*. *Preprint*, arXiv:2505.09388.

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,

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 ব্যথা 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", "সাত্বে", "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.

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.

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

Roberto Zanolli, Alberto Lavelli, Daniel Verdi do Amarante, and Daniele Toti. 2024b. *Assessment of the E3C corpus for the recognition of disorders in clinical texts*. *Natural Language Engineering*, 30:851–869.

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.

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

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। ফাটিয়ে দিবেন না। ইনফেকশন হয়ে যাবে। আপাতত এলাজি উল্লেখকারী খাবার, ধূলাবালি এড়িয়ে চলুন। এলাট্রিল খান। না দেখে সমাধান দেয়া যাচ্ছে না। সরাসরি চর্মরোগ বিশেষজ্ঞ দেখিয়ে পরামর্শ নিতে হবে ধন্যবাদ</p>	Medicine: এলাট্রিল Specialist: চর্মরোগ বিশেষজ্ঞ	Medicine: এলাট্রিল Health_Condition: ইনফেকশন Specialist: চর্মরোগ বিশেষজ্ঞ
2	json_ 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 findout 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 বাজাকে আদা, মধু, গরম পানি খাওয়ান। সিরাপ এমব্রোক্স খাওয়াতে পারেন আধা চামচ করে তিন বার একজন শিশু বিশেষজ্ঞ কে দেখিয়ে নিন।</p>	Medicine: আদা গরম পানি মধু সিরাপ এমব্রোক্স ; Specialist: শিশু বিশেষজ্ঞ	Medicine: সিরাপ এমব্রোক্স ; Specialist: শিশু বিশেষজ্ঞ
3	json_ 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 পচত পায়ের হার এর বেথা। আমি কি ঔষধ খাব বলেন স্যার। Calcium BP 500mg কিনছি খাব কি স্যার।</p>	Age: ১৩ বছর ; Symptom: পচত পায়ের হার এর বেথা ; Medicine: Calcium BP 500mg ; Specialist: স্যার	Age: ১৩ ; Symptom: পচত 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 Bihotzhodietako Kirurgiako Zerbitzuarekin harremanetan jarri gara.</i></p>	Disorder: bihotz- gutxiegitasuna estenosi aortiko larria sinkopea ; H- Professional: Bihotz- hodietako Kirur- giako Zer- bitzuarekin	Disorder: Estenosi aor- tiko bihotz- gutxiegitasuna sinkopea]
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 primaria.</i></p>	Disorder: Sifilis pri- maria Txankro sifil- itikoa	Disorder: Sifilis Txankro sifil- itikoa
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 daqo bena sakonetako tronbosi-zeinurik.</i></p>	Disorder: tronbosi- zeinurik	Disorder: Edemarik bena sakonetako tronbosi- zeinurik

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.

misses the HEALTH_CONDITION (ইনফেকশন), 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 recall by encouraging semantic inference, yet it also leads to over-generation, with the model extracting implicit or negated find-

Error Type	Input Example	Gold Output	Pred Output
Hallucinations	ধনাবাদ আপনাকে প্রশ্নের জন্য। জ্বর ৯৯ হলে ওষুধ খাওয়াবেন। নয়ত দরকার নেই। ধনাবাদ	<i>Empty</i>	Symptom: জ্বর
All-Missed	আমি। আমার বয়স 27 বছর। সাম্প্রতিক সময়ে আমি আমার শরীরে অস্বস্তি...গত হয়ে গেছি, আগের থেকে শরীরের ওজন বৃদ্ধি পেয়েছে, পেটে প্র...	Age: 27 বছর ; Symptom: একাধারে মাটিতে বসে থাকলে মাজার ... শরীরের ওজন বৃদ্ধি হাটলেই মনে হয় ক্লান্ত হয়ে	<i>Empty</i>
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 ধনাবাদ প্রশ্ন করার জন্য। আপনি গরম পানির ধোয়া নাকে নিয়...রাতে একটা করে খান। একজন মেডিসিন বিশেষজ্ঞ পরামর্শ নিন। ধনাবাদ	Medicine: এলাট্রিল ; Specialist: মেডিসিন বিশেষজ্ঞ	Medicine: এলাট্রিল গরম পানির ধোয়া ; 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 বের হয় 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: <i>Burua eta lepoa</i>
All-Missed	<i>Ez dago aurkikuntza patologikorik.</i>	Disorder: <i>aurkikuntza</i>	<i>Empty</i>
Missed Entities	<i>Estenosi aortiko larria eta bihotz-gutxiegitasuna.</i>	Disorder: <i>Estenosi aortiko larria bihotz-gutxiegitasuna</i>	Disorder: <i>bihotz-gutxiegitasuna</i>
Extra Entities	<i>Sabela.</i>	<i>Empty</i>	Patient: <i>Sabela</i>
Type Confusion	<i>Biguna eta zanpagarria.</i>	<i>Empty</i>	Disorder: <i>Biguna zanpagarria</i>
Mixed Errors	<i>Biriketako murmurio normala.</i>	<i>Empty</i>	Disorder: <i>Biriketako murmurio</i>
Boundary Mismatch	<i>Ez dago aurkikuntza patologikorik.</i>	Disorder: <i>aurkikuntza patologikorik</i>	Disorder: <i>aurkikuntza</i>

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.

ings (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

900 text: short phrases and anatomical references
901 are frequently over-labeled as medical enti-
902 ties even when the gold annotation contains
903 none. We also see **Boundary Mismatch** in
904 cases where the model selects a shorter head
905 noun rather than the full descriptive mention,
906 which is penalized under exact-match evalu-
907 ation. Compared with Bangla, Basque ex-
908 hibits more frequent over-extraction and fewer
909 purely recall-only failures, consistent with the
910 overall error distribution where FP-heavy er-
911 rors occupy a larger fraction of error-only
912 cases.