

Can LLMs extract human-like fine-grained evidence for evidence-based fact-checking?

Antonín Jarolím*, Martin Fajčík and Lucia Makaiová

Brno University of Technology, Czech Republic

Abstract. Misinformation frequently spreads in user comments under online news articles, highlighting the need for effective methods to detect factually incorrect information. To strongly support or refute claims extracted from such comments, it is necessary to identify relevant documents and pinpoint the exact text spans that justify or contradict each claim. This paper focuses on the latter task — fine-grained evidence extraction for Czech and Slovak claims. We create new dataset, containing two-way annotated fine-grained evidence created by paid annotators. We evaluate large language models (LLMs) on this dataset to assess their alignment with human annotations. The results reveal that LLMs often fail to copy evidence verbatim from the source text, leading to invalid outputs. Error-rate analysis shows that the `llama3.1:8b` model achieves a high proportion of correct outputs despite its relatively small size, while the `gpt-oss-120b` model underperforms despite having many more parameters. Furthermore, the models `qwen3:14b`, `deepseek-r1:32b`, and `gpt-oss:20b` demonstrate an effective balance between model size and alignment with human annotations.

Keywords: Fact-checking · Fine-grained evidence · LLMs.

1 Introduction

On average, around three quarters of more than twelve thousand respondents across news sites expressed interest in having topic experts respond to comments in articles [15]. Unfortunately, this is highly impractical and therefore, an automatic discussion management is required. Major platforms have already recognized this need — for example, platform *X* (formerly Twitter) relies on human-generated “Community Notes”, while *Seznam.cz* supplements selected comments with fact-checked contextual information.

Another promising direction for managing discussions is the automatic provision of relevant documents drawn from collections of trustworthy sources such as Reuters. However, locating the precise evidence within a relevant document can be time-consuming, and providing the entire document as evidence is typically too coarse to be persuasive. Therefore, identifying fine-grained evidence within the document is crucial. Readers can evaluate a short, well-targeted text

* Correspondence to: ijarolim@fit.vut.cz

span more quickly without decreasing the accuracy of judgement [10]. Consequently, the fine-grained evidence is both *efficient* and *effective* for mitigating misinformation.

To automate the management of online discussions, large language models (LLMs) can be employed to extract claims from user comments. A retrieval model can then identify relevant documents, after which LLMs can locate fine-grained evidence within the retrieved content. This study focuses exclusively on the latter step — identifying fine-grained evidence in source documents. Specifically, the contributions of this paper include:

1. A *manually constructed* dataset comprising check-worthy claims written in Czech or Slovak languages is constructed, each paired with a highly relevant document identified by annotators using various search engines. Two annotators then highlight fine-grained evidence supporting each claim.
2. Fine-grained evidence is identified using several LLMs, including the 685B DeepSeek-R1 and 120B gpt-oss reasoning models, as well as a range of smaller open-weight models such as 27 billion (B) Gemma-3 and 14 B Phi4.
3. The performance of the LLMs is analysed in terms of error rates and alignment with human annotations in the context of fine-grained evidence generation.

In this work, we focus on the analysis of *supporting evidence*¹ only.

2 Related work

Automatic fact-checking. FactLens [13] decomposes complex claims into sub-claims and evaluates the veracity of each independently. Loki [11] extends this approach through an automated pipeline that identifies check-worthy claims, retrieves evidence, and verifies them. However, in both methods, the evidence used to determine a claim’s truthfulness is provided only at the passage level.

Furthermore, AmbiFC [5] introduces ambiguity into automated fact-checking by incorporating multiple sentence-level annotations with potentially divergent labels. It shows that models learning soft-label distributions for sentence-level evidence selection and veracity prediction achieve superior performance. This emphasizes the importance of fine-grained evidence for improving fact-checking accuracy.

Fact-checking datasets. While AmbiFC [5] claims to provide fine-grained evidence alongside claims, its annotations remain at the passage level. Other datasets, such as FEVER [17] and SciFact [18], offer finer granularity, as both include claims paired with sentence-level evidence. The former focuses on general claims derived from Wikipedia, whereas the latter provides rationale annotations within scientific paper abstracts.

¹ The same kind of analysis can be performed on refuting evidence. We leave this to future work.

However, no comparable dataset exists for Czech or Slovak data, and, to our knowledge, none has evidence annotation at the span level. To address this gap, we construct a new dataset comprising Czech and Slovak samples with manually annotated fine-grained evidence spans.

Reasoning capabilities of LLMs. The performance of LLMs continues to improve with increasing model size, architectural advancements [12], and enhanced reasoning capabilities [7]. Recent work further demonstrates that LLMs exhibit strong reasoning performance on Czech and Slovak datasets [2].

These developments open the possibility of employing LLMs for automatic fine-grained evidence extraction. As a first step toward automated fact-checking with fine-grained evidence, we evaluate the alignment between human annotators and LLMs on this task.

3 Method

This section outlines the methods for fine-grained evidence extraction. Given a claim and its associated text, the task is formally defined as follows:

Problem statement. Given a *claim* and a tokenized *text* $t = (t_1, t_2, \dots, t_N)$, select a set of spans $S = \{s_1, s_2, \dots, s_M\}$, where $N, M \in \mathbb{N}^+$ and each span $s_m = (t_i, \dots, t_j)$, with $i, j \in \mathbb{N}^+$ and $i \leq j$, denotes a contiguous subsequence of t that supports the claim.

The following section outlines how this task is approached by human annotators, baselines, and LLMs.

Two-Way Annotation of Fine-Grained Evidence Dataset. We collected 186 claim-text pairs and obtained two independent fine-grained evidence annotations for each sample, created by different annotators from a pool of eight non-expert annotators. The first annotation was produced using a custom annotation tool, while the second was created in Label Studio². Annotation guidelines stated:

- Highlight the *minimal* portion of text that supports or refutes the claim.
- Highlight the part that most convinces you that the given statement is *true*.

Fortunately, both annotation interfaces allowed annotators to directly highlight evidence spans. Thus, each selected subsequence corresponds to a contiguous span appearing in the text. In contrast, LLMs were required to regenerate the selected span, allowing generation of subsequences not appearing in the text, as explained below.

Evidence Extraction using LLMs. We employ a diverse set of LLMs of varying sizes to perform fine-grained evidence extraction. Specifically, we use

² <https://labelstud.io/>

qwen2.5 (72B, 32B) [20], llama3.3 (70B) [6], gemma2 (27B) [16], phi4 (14B) [1], llama3.1 (8B) [16], gemma3 (27B, 12B, 4B) [3] and mixtral (8×7B) [8]. Additionally, we include Chain-of-Thought (CoT) LLMs that produce intermediate reasoning steps before providing the final answer gpt-oss-120b (20B, 120B) [14], deepseek-r1 (685B, 32B) [7], qwen3 (32B, 14B) [19].

Similarly to the human annotators, the LLMs were presented with claim-text pairs, along with the comment from which the claim was originally extracted to provide additional context. The models were instructed to identify the smallest possible segments of the text that directly support the claim and to output a JSON-formatted list of the selected spans (see Appendix B for the complete prompt). Although the LLMs were explicitly instructed to generate only spans appearing verbatim in the text, this constraint was not enforced technically. Consequently, as discussed in the following section, the models occasionally produced spans that did not occur in the source text.

Baseline Approaches. Additionally, we include non-neural approaches. Firstly, the *claim baseline* tokenizes the claim c into word tokens $c = (c_1, c_2, \dots, c_O)$ and overlaps them with the corresponding text t , constructing the set of selected spans $S_C = \{(t_i, \dots, t_j) | \exists k, l \in \mathbb{N} : (t_i, \dots, t_j) = (c_l, \dots, c_k) \wedge l \leq k\}$. The queries used to search for claim evidences during the fine-grained evidence annotations were also stored. These queries form the basis of the *query baseline*, which overlaps the text with the query in the same manner as the claim baseline. Finally, the *random baseline* uniformly samples contiguous spans whose number and length match those of a randomly chosen annotator for each sample.

4 Results

All experiments are conducted on the manually annotated dataset described above, comprising 186 samples with two independent sets of fine-grained evidence annotations.

Empty Annotations & Model Parameters																
invalid %	4.3	7.0	13.4	16.7	17.7	24.2	27.4	28.5	30.1	30.6	32.8	35.5	36.0	40.3	57.5	61.8
Params (B)	72	685	8	32	27	27	70	32	20	12	120	32	14	14	4	7
qwen2.5:72b	deepseek-r1	llama3.1:8b	qwen2.5:32b	gemma2:27b	gemma3:27b	llama3.3:70b	deepseek-r1:32b	gpt-oss:20b	gemma3:12b	gpt-oss:120b	qwen3:32b	phi4:14b	qwen3:14b	gemma3:4b	mixtral:8x7b	

Fig. 1. Various LLMs used to generate fine-grained extraction, their number of parameters in billions, and the percentage of incorrectly annotated samples (invalid %) they generated.

Error Analysis of LLM-Generated Evidence. In general, model performance tends to improve with the number of parameters [6,12]. This trend is also evident in the present task, where the inability to generate valid output is interpreted as a failure to follow instructions.

The error rates of all evaluated LLMs are presented in Figure 1. The highest number of incorrect annotations 61.8% was produced by the `mixtral:8x7b` model. This result is consistent with the fact that `mixtral:8x7b` is among the smallest models in the comparison. Conversely, the `qwen2.5:72b` model achieved the lowest error rate, which aligns with its considerably larger parameter count.

A few models deviate from this overall pattern. Notably, `llama3.3:72b` exhibits a relatively high number of invalid outputs despite its large size, and `gpt-oss:120b`, which would be expected to perform among the best, also shows unexpectedly poor reliability. Interestingly, `llama3.1:8b`, one of the smallest models, has a number of incorrectly generated samples similar compared with huge `deepseek-r1` model.

Given that the dataset contains only 186 samples and that some models produce up to 116 incorrect outputs, the overall error rates remain substantial. Consequently, future work should explore methods that enforce the generation of semantically and structurally valid evidence. Approaches such as constrained decoding or structured output generation could substantially reduce the occurrence of erroneous outputs [4].

LLMs Extraction Performance. Before presenting the results, it should be emphasized, that certain LLM pairs matches are evaluated using a substantially reduced subset of data. This is because, as previously discussed, LLMs can fail to generate valid annotations. Excluding these missing data points may introduce bias — for instance, by disproportionately removing more challenging samples and thereby artificially inflating the performance metrics of certain models. Nevertheless, this exclusion is necessary to ensure that model pairs are evaluated on the same set of available samples.

To reduce noise and improve comparability in the evaluation, we remove stop words (see Appendix A for entire list of stop words) from all evidence sets before computing F1-scores. The degree of overlap between two annotation sets, is then computed using token-level F1 score. Since neither human annotators nor LLMs are instructed to produce exhaustive span selections, the number of annotated spans may differ between the two sets. To not penalize different degree of exhaustiveness, we use the Hungarian matching algorithm [9] to find optimal assignment between two annotation sets. First, we compute the F1-score for all possible span pairs across the two sets. Second, we apply the Hungarian algorithm to solve the assignment problem, ensuring that each span is matched to at most one span in the other set while maximizing the total F1. Finally, the average F1 of the resulting optimal matching represents the token-level F1 for a single data point.

The average F1-score between two annotators is 48, as we can see in the Figure 2. LLMs `deepseek-r1:32b` and `qwen3:14b` achieve the strongest performance among the evaluated models, with alignment scores of 55 and 56, respec-

		Best scores Hungarian matching - F1																					
Prediction	ann 1 (LS)	100	48	38	42	43	34	55	39	47	41	48	45	47	56	42	35	52	29	10	17	12	
	ann 2	48	100	34	35	27	28	33	32	34	25	34	39	30	35	30	20	29	21	7	20	16	
	deepseek-r1	38	34	100	68	32	45	46	42	52	33	44	61	30	58	42	21	34	26	5	21	17	
	gpt-oss-120b	42	35	68	100	34	47	48	44	52	33	53	67	34	59	43	23	36	30	6	21	17	
	qwen2.5:72b	43	27	32	34	100	46	62	51	53	49	49	49	50	58	49	40	49	36	11	19	16	
	llama3.3:70b	34	28	45	47	46	100	50	52	60	43	50	59	41	66	52	32	47	37	6	20	16	
	deepseek-r1:32b	55	33	46	48	62	50	100	60	64	49	58	55	66	68	55	46	57	35	13	17	14	
	qwen2.5:32b	39	32	42	44	51	52	60	100	60	45	55	58	51	63	52	40	54	34	10	18	15	
	qwen3:32b	47	34	52	52	53	60	64	60	100	48	58	60	52	67	56	38	50	34	10	21	17	
	gemma2:27b	41	25	33	33	49	43	49	45	48	100	52	48	48	56	51	39	49	34	13	16	13	
	gemma3:27b	48	34	44	53	49	50	58	55	58	52	100	66	54	64	58	41	53	42	10	22	16	
	gpt-oss:20b	45	39	61	67	49	59	55	58	60	48	66	100	47	65	55	35	49	39	7	21	17	
	phi4:14b	47	30	30	34	50	41	66	51	52	48	54	47	100	62	48	48	66	44	12	20	17	
	qwen3:14b	56	35	58	59	58	66	68	63	67	56	64	65	62	100	59	40	51	42	14	23	17	
	gemma3:12b	42	30	42	43	49	52	55	52	56	51	58	55	48	59	100	43	55	40	11	25	17	
	llama3.1:8b	35	20	21	23	40	32	46	40	38	39	41	35	48	40	43	100	39	38	11	15	13	
	mixtral:8x7b	52	29	34	36	49	47	57	54	50	49	53	49	66	51	55	39	100	30	15	21	20	
	gemma3:4b	29	21	26	30	36	37	35	34	34	34	42	39	44	42	40	38	30	100	9	19	14	
random	10	7	5	6	11	6	13	10	10	13	10	7	12	14	11	11	15	9	100	7	4		
claim	17	20	21	21	19	20	17	18	21	16	22	21	20	23	25	15	21	19	7	99	62		
query	12	16	17	17	16	16	14	15	17	13	16	17	17	17	17	13	20	14	4	62	88		
		Ground Truth																					

Fig. 2. Token-level F1 scores (on a 1–100 scale) between human annotations, large language models (LLMs), and non-neural baseline models. Baselines **random**, **claim**, and **query** are included for comparison. **ann 1 (LS)** refers to annotations created in Label Studio, while **ann 2** denotes annotations collected through the custom annotation interface.

tively, against the annotations created in Label Studio — surpassing even the human–human agreement. However, it should be emphasized that these scores are calculated only on samples where LLM annotations were generated correctly and human annotations were available.

All the non-neural baseline methods proves to be weak, as all of them have F1-score less then 18. Further analysis shows, that the precision for the **claim** and **query** baselines is high — around 30, showcasing that words in the claim and query are often also used as evidence by annotators. However, the recall is very low, therefore decreasing the overall F1 performance significantly.

Lastly, a comparison between the two sets of human annotations reveals systematic differences. All evaluated LLMs achieve higher alignment with the annotations created in Label Studio, than with those produced in the custom an-

notation environment, suggesting that the two annotation settings yield slightly different annotation styles or levels of granularity. Further analysis is required to reveal the differences between annotations.

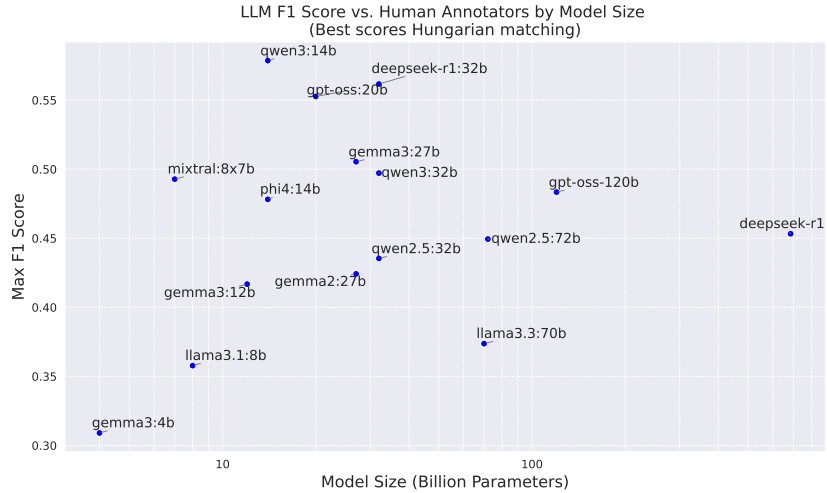


Fig. 3. Token-level F1 scores between large language models (LLMs) and human annotations against the model size in billions (log scale). For each data point, the maximum F1 score across the two annotators is used, reflecting alignment with at least one human annotator.

Alignment with Human Annotators. To directly evaluate alignment with human annotators, each LLM-generated evidence span is compared separately with both annotators’ labels. The higher of the two token-level F1 scores is then selected, reflecting agreement with at least one annotator.

The results of this analysis are presented in Figure 3, where the maximum token-level F1 score is plotted against model size. While performance generally improves from small to medium-sized models, adding substantially more parameters (e.g., 685B `deepseek-r1` or 120B `gpt-oss`) does not yield further gains. This suggests that beyond a certain threshold, model size alone does not correlate with better extraction performance. Nevertheless, the models 14B `qwen3`, 32B `deepseek-r1`, and 20B `gpt-oss` exhibit a favourable trade-off between parameter count and alignment with human annotations.

5 Conclusion

We introduce a new dataset of Czech and Slovak texts with fine-grained evidence annotations, produced by two independent annotators for each sample. This

dataset enables the computation of inter-annotator agreement, which we measure using the Hungarian matching algorithm and the Token-F1 metric, resulting in a score of 47. Additionally, it provides a foundation for evaluating the alignment of LLM-generated evidence with human judgments — filling an existing gap in available resources.

Using this dataset, we analysed the ability of various LLMs to generate valid fine-grained evidence. We observed a clear relationship between model size and the proportion of valid outputs: smaller models such as 4 B **gemma3** and 8 B **mixtral** exhibited error rates exceeding 50%. This highlights the requirement to employ constrained decoding mechanisms in the further work.

Despite these generation errors, our results on LLMs extraction performance show diminishing returns with increasing model size. Performance improves from small to medium-sized models, but beyond a certain threshold additional parameters do not yield better extraction accuracy. Notably, 14 B **qwen3**, 32 B **deepseek-r1** and 20 B **gpt-oss** offer the best trade-off between model size and alignment with human annotations.

Acknowledgements. This work was supported by the Technology Agency of the Czech Republic (TAČR) under the SIGMA Programme, 8th Public Competition, Sub-objective 4: Bilateral Cooperation, project TQ16000028.

References

1. Abdin, M., Aneja, J., Behl, H., Bubeck, S., Eldan, R., Gunasekar, S., Harrison, M., Hewett, R.J., Javaheripi, M., Kauffmann, P., et al.: Phi-4 technical report. arXiv preprint arXiv:2412.08905 (2024)
2. Fajcik, M., Docekal, M., Dolezal, J., Ondrej, K., Beneš, K., Kapsa, J., Smrz, P., Polok, A., Hradis, M., Neverilova, Z., et al.: Benczechmark: A czech-centric multitask and multimetric benchmark for large language models with duel scoring mechanism. *Transactions of the Association for Computational Linguistics* **13**, 1068–1095 (2025)
3. Gemma Team et al.: Gemma 3 technical report (2025), <https://arxiv.org/abs/2503.19786>
4. Geng, S., Josifoski, M., Peyrard, M., West, R.: Grammar-constrained decoding for structured nlp tasks without finetuning. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. pp. 10932–10952 (2023)
5. Glockner, M., Staliūnaitė, I., Thorne, J., Vallejo, G., Vlachos, A., Gurevych, I.: Ambifc: Fact-checking ambiguous claims with evidence. *Transactions of the Association for Computational Linguistics* **12**, 1–18 (2024)
6. Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al.: The llama 3 herd of models. ArXiv preprint **abs/2407.21783** (2024), <https://arxiv.org/abs/2407.21783>
7. Guo, D., Yang, D., Zhang, H., Song, J., Wang, P., Zhu, Q., Xu, R., Zhang, R., Ma, S., Bi, X., et al.: Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature* **645**(8081), 633–638 (2025)
8. Jiang, A.Q., Sablayrolles, A., et al., A.R.: Mixtral of experts (2024), <https://arxiv.org/abs/2401.04088>

9. Kuhn, H.W.: The hungarian method for the assignment problem. *Naval research logistics quarterly* **2**(1-2), 83–97 (1955)
10. Leonhardt, J., Rudra, K., Anand, A.: Extractive explanations for interpretable text ranking. *ACM Trans. Inf. Syst.* **41**(4) (2023). <https://doi.org/10.1145/3576924>, <https://doi.org/10.1145/3576924>
11. Li, H., Han, X., Wang, H., Wang, Y., Wang, M., Xing, R., Geng, Y., Zhai, Z., Nakov, P., Baldwin, T.: Loki: An open-source tool for fact verification. In: *Proceedings of the 31st International Conference on Computational Linguistics: System Demonstrations*. pp. 28–36 (2025)
12. Liu, E., Bertsch, A., Sutawika, L., Tjuatja, L., Fernandes, P., Marinov, L., Chen, M., Singhal, S., Lawrence, C., Raghunathan, A., et al.: Not-just-scaling laws: Towards a better understanding of the downstream impact of language model design decisions. *arXiv preprint arXiv:2503.03862* (2025)
13. Mitra, K., Zhang, D., Rahman, S., Hruschka, E.: Factlens: Benchmarking fine-grained fact verification. In: *Findings of the Association for Computational Linguistics: ACL 2025*. pp. 18085–18096 (2025)
14. OpenAI team et al.: gpt-oss-120b & gpt-oss-20b model card (2025), <https://arxiv.org/abs/2508.10925>
15. Stroud, N.J., Van Duyn, E., Alizor, A., Alibhai, A., Lang, C.: Comment section survey across 20 news sites. Report published by the Center for Media Engagement, University of Texas at Austin (2017)
16. Team, G., Riviere, M., Pathak, S., Sessa, P.G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al.: Gemma 2: Improving open language models at a practical size. *ArXiv preprint abs/2408.00118* (2024), <https://arxiv.org/abs/2408.00118>
17. Thorne, J., Vlachos, A., Christodoulopoulos, C., Mittal, A.: Fever: a large-scale dataset for fact extraction and verification. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. pp. 809–819 (2018)
18. Wadden, D., Lin, S., Lo, K., Wang, L.L., van Zuylen, M., Cohan, A., Hajishirzi, H.: Fact or fiction: Verifying scientific claims. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pp. 7534–7550 (2020)
19. Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al.: Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025)
20. Yang, A., Yang, B., gpt-oss-120bet al., B.Z.: Qwen2.5 technical report (2025), <https://arxiv.org/abs/2412.15115>

A List of Stop-Words

aby	dalsi	kam	ne	prvni	tohle	a	když	ní	takže
aj	design	kde	nebo	pta	toho	aniž	která	nové	te
ale	dnes	kdo	nejsou	re	tohoto	ano	které	nový	tě
ani	do	kdyz	neni	si	tom	až	který	o	těma
asi	email	ke	nez	strana	tomto	budeš	kterí	ode	této
az	ho	ktera	nic	sve	tomuto	být	ku	on	tím
bez	jak	ktere	nove	svych	tu	což	máte	práve	tímto
bude	jako	kteri	novy	svym	tuto	či	me	proč	toto
budem	je	kterou	od	svymi	ty	článek	mě	protože	tvůj
budes	jeho	ktery	pak	ta	tyto	článku	mít	první	u
by	jej	ma	po	tak	uz	články	mně	před	už
byl	jeji	mate	pod	take	vam	další	mnou	přede	v
byla	jejich	mezi	podle	takze	vas	i	můj	přes	vám
byli	jen	mi	pokud	tato	vase	já	může	při	váš
bylo	jeste	mit	pouze	tedy	ve	její	my	s	vaše
byt	ji	muj	prave	tema	vice	jenž	ná	se	více
ci	jine	muze	pred	ten	vsak	ještě	nám	sice	však
clanek	jiz	na	pres	tento	za	jiné	napište	své	všechn
clanku	jsem	nad	pri	teto	zda	již	náš	svůj	vy
clanky	jes	nam	pro	tim	zde	jseš	naši	svých	z
co	jsme	napiste	proc	timto	ze	jšte	necht'	svým	zpět
coz	jsou	nas	proto	tipy	zpet	k	není	svými	zprávy
cz	jste	nasi	protoze	to	zpravy	každý	než	také	že

B Prompt Used to Extract Fine-Grained Evidence

Fine-Grained Evidence Extraction Instructions Prompt

Comment (for context): {{source_comment}}
Claim (extracted from the comment): {{claim}}
Text: {{text}}

Your task is to identify the smallest possible parts of the Text that directly support the claim.

Focus on the phrases that most clearly justify or confirm the truth of the claim. Avoid selecting entire sentences unless absolutely necessary—choose the shortest meaningful spans that stand on their own. You may select multiple spans if more than one part of the text provides evidence.

Do not modify, correct, or rewrite the text. Preserve all grammatical and syntactic errors exactly as they appear in the original.

Return only a JSON object of type `Dict[str, List[str]]` with the key `'spans'` and a list of selected spans as its value.

Important: Every selected span must be an exact substring of the given Text — character-for-character identical. Do not paraphrase, retype, or alter any character (including punctuation, spacing, or capitalization).

Fig. 4. Prompt used for large language models to extract fine-grained evidence supporting the claim. During the generation, {{source_comment}}, {{claim}} and {{text}} placeholders were substituted with the source comment containing claim, text of the claim and the text to extract fine-grained relevance from, respectively.