

They want to pretend not to understand: The Limits of Current LLMs in Interpreting Implicit Content of Political Discourse

Anonymous ACL submission

Abstract

Implicit content plays a crucial role in political discourse, where systematically employ pragmatic strategies such as implicatures and presuppositions to influence their audiences. Large Language Models (LLMs) have demonstrated strong performance in tasks requiring complex semantic and pragmatic understanding, highlighting their potential for detecting and explaining the meaning of implicit content. However, their ability to do this within political discourse remains largely underexplored. Leveraging, for the very first time, the large IMPAQTS corpus comprising transcribed Italian political speeches with expert annotations of various types of implicit content, we propose methods to test the effectiveness of LLMs in this challenging problem. Through a multiple-choice task and an open-ended generation task, we demonstrate that all tested models struggle to interpret presuppositions and implicatures. To illustrate, the best-performing model provides a fully correct explanation in only one-fourth of cases in the open-ended generation setup. We conclude that current LLMs lack the key pragmatic capabilities necessary for accurately interpreting highly implicit language, such as that found in political discourse. At the same time, we highlight promising trends and future directions for enhancing model performance. We release our data and code at: <http://anonymized/for/submission/>

1 Introduction

Implicit and vague language is pervasive in everyday life communication. Consider the language used in political contexts: politicians often use implicit or vague content, relying on the listener's assumptions, background knowledge, and shared contextual knowledge to convey meaning indirectly. This can be seen in the use of vague terms, allusions, and rhetorical devices that allow for multiple interpretations. Such language creates an environment where meanings can be conveyed sub-

tly, often with the intent of persuading, manipulating, or reinforcing certain ideologies (Morency et al., 2008). Through the careful use of implicit language, politicians can appeal to the audience's emotions, assumptions, and cultural references, all while avoiding explicit commitments that could later be challenged or disproved (Lombardi Vallauri, 2019). Consider the following example (also featured in the title of this work): *They want to pretend not to understand.*¹ By saying this, the speaker implies that the people being addressed pretend they do not understand something obvious, while, in reality, they understand it perfectly well. This pragmatic phenomenon, where the speaker says one thing to imply another, is called an implicature. Implicatures are among the most well-studied phenomena of implicit language (Grice, 1990; Sperber and Wilson, 1987), alongside presuppositions (Strawson, 1964; Garner, 1971; Ducrot, 1972). As both are pervasive in political discourse, which makes it an ideal domain for analyzing implicit communication strategies (Van Dijk, 1992; Lombardi Vallauri and Masia, 2020).

The advent of Large Language Models (LLMs) has revolutionized NLP, bringing unprecedented progress to tasks involving semantic understanding (Wang et al., 2018; Williams et al., 2018), question answering (Mihaylov et al., 2018), and pragmatic interpretation (Zheng et al., 2021; Sravanthi et al., 2024; Hu et al., 2023; Kim et al., 2023; Ruis et al., 2024). Overall, LLMs have been shown to perform excellently on semantic tasks and to exhibit potential in pragmatic tasks, particularly when they involve artificially constructed stimuli. However, whether they can handle tasks that require making pragmatic inferences on naturalistic dis-

¹This example is the English translation of an Italian sentence, in the data used in this work. In the main paper, for simplicity, we only provide English translations. The original Italian texts are available in Appendix A, C.1 and F. In this section, examples may sometimes be adapted for illustrative purposes.

course data, such as political discourse, remains an open question. In this work, we specifically address this research question and test whether current LLMs can explain the meaning of implicit content in political discourse, such as the implicature in the example above. Leveraging, for the first time in the context of NLP research, the large IMPAQTS corpus (Cominetti et al., 2024), which includes transcribed Italian political speeches with expert annotations of various types of implicit content, we propose methods to assess the pragmatic effectiveness of several LLMs, state-of-the-art for the Italian language. Our contributions are as follows:

- We propose **a novel task for LLMs**, i.e., to clarify, through an explanation, the meaning of an implicit content found in political discourse; e.g., that in the passage *They want to pretend not to understand*, the speaker implies that these people understand perfectly well.
- We release **a novel dataset**, based on the IMPAQTS corpus (Cominetti et al., 2024), containing 30K implicit passages (implicatures and presuppositions) with expert-based explanations, as well as the surrounding linguistic context necessary to correctly interpret the implicit content (empirically validated through human expert annotation).
- Through a multiple-choice task and an open-ended generation task, we demonstrate that **all tested models struggle** to interpret presuppositions and implicatures. In the multiple-choice setup, the best-performing model falls more than 20 accuracy points short of the estimated ceiling; in the open-ended generation setup, it provides a fully correct explanation in only one-fourth of the cases.
- Despite largely unsatisfactory results, we show that using **Chain-of-Thought (CoT) improves performance** in the open-ended generation setup, suggesting that reasoning mechanisms can aid in solving complex pragmatic tasks. We also propose other directions, such as **embedding external knowledge** about the politician and their political affiliation to further enhance model performance.

2 Related Work

2.1 Implicit Language and LLMs

As LLMs power conversational agents that interact with human users, they must be capable of understanding implicit content. Prior research has explored LLMs' capacity to interpret non-literal meanings in linguistic discourse. Jeretic et al. (2020) investigated the emergence of pragmatic understanding by analyzing how pre-trained language models compute scalar implicatures and presuppositions. They found strong evidence that models such as BERT learned scalar implicatures for quantifiers *some* and *all*, but struggled with other scalar pairs treating them as synonymous. For presuppositions, models failed to recognize some presupposition triggers, e.g. the verb *to stop* when it presupposes an action that used to be made as in *John stopped smoking*. Zheng et al. (2021) proposed the first evaluation of language models on five types of implicatures derived from Grice's (1990) conversational maxims, finding that models performed relatively well on a multiple-choice task, where they had to choose which of the given options explained the implicature in a dialogue, but struggled significantly with conversational reasoning, suggesting they do not fully understand conversational context. On a similar track, Hu et al. (2023) evaluated pre-trained LMs without task-specific fine-tuning, incorporating Gricean implicatures as one of their tested phenomena among many others. Their results suggested that language models can infer pragmatic meanings, though they left unresolved whether this is due to linguistic cues or cognitive processes, highlighting the need for further research on their connection in pragmatic reasoning. Kim et al. (2023) evaluated whether LLMs can understand conversational implicatures by prompting them to provide binary answers to specific scenarios. They demonstrated that while LLMs exhibited some implicit understanding, their performance improved significantly when guided through the reasoning process using chain-of-thought prompting. More recently, Ruis et al. (2024) designed a protocol to evaluate LLMs on binary implicature resolution, highlighting a significant gap between humans and LLMs.

All these studies investigated plausible naturalistic use of pragmatic language. Nonetheless, they used artificially constructed sentences instead of real examples of pragmatic phenomena found in actual *corpora*. In our work, we make a significant

177 step forward and use excerpts of semi-spontaneous,
178 ecological political discourse that reflect the com-
179 plexities of real-world communication. Unlike ar-
180 tificially constructed sentences, political discourse
181 is known for capturing the complex interplay of
182 rhetoric, social context, and speaker intent.

183 2.2 Political Language in NLP

184 Political language has been the focus of NLP re-
185 search leveraging large datasets from social me-
186 dia and political speeches to study discourse fram-
187 ing, bias detection, and polarization. [Katre \(2019\)](#)
188 employed computational methods for text analyt-
189 ics and visualization in political speech transcripts.
190 The author used NLP techniques to analyze politi-
191 cal speeches and generate graphical visualizations
192 for, among others, word use, lexical dispersion, etc.
193 In [Huguet Cabot et al. \(2020\)](#), NLP methods were
194 used to model political discourse by focusing on
195 metaphor, emotion, and political rhetoric. The au-
196 thors presented the first joint models that integrate
197 metaphor and emotion detection as auxiliary tasks
198 to enhance the performance of predicting the politi-
199 cal perspective of news articles, party affiliation of
200 politicians, and framing of policy issues. [Németh
201 \(2023\)](#) presented a methodological review on how
202 researchers have used NLP techniques to study lan-
203 guage polarization between 2010 and 2021. The
204 author identified data sources and computational
205 approaches and reviewed the different conceptual-
206 izations and operationalizations of polarization.

207 Recent efforts leveraged LLMs to enhance politi-
208 cal discourse understanding. [Marino and Giglietto
209 \(2024\)](#) used LLMs for political discourse annota-
210 tion, reporting significant improvement on previous
211 methods like topic modeling. Here, LLMs were
212 used to analyze Facebook links related to the 2018
213 and 2022 Italian elections and perform tasks such
214 as classification and clustering of political content
215 and generation of descriptive labels for these clus-
216 ters. [Li et al. \(2024\)](#) proposed Political-LLM, a
217 comprehensive framework for integrating LLMs
218 into political science. This offered unprecedented
219 capabilities for text analysis. At the same time,
220 the authors highlighted several issues related to
221 biases, ethics, and explainability to be considered
222 when using LLMs in political applications. Despite
223 these advancements, little attention has been paid
224 to automatically explaining implicit content. To our
225 knowledge, we are the first to tackle this problem.

3 Data

3.1 The IMPAQTS Corpus

We use, for the very first time in the context of NLP research, data from IMPAQTS ([Cominetti et al., 2024](#))², a corpus of Italian political discourse (in the sense of "discourse by politicians"; see [Van Dijk et al., 1997](#)), consisting of 1,500 speeches uttered in the Italian language by 150 prominent politicians between 1946 and 2023. IMPAQTS is a multimodal corpus containing over 800 speeches in video format and around 600 in audio format. The manual transcriptions of these speeches are also available, building up to roughly 2.65 million tokens. In our work, we only use these transcriptions.

A key feature of IMPAQTS, crucial for our work, is the annotation of all the passages that contain *implicit content with some manipulative meaning*. In more technical terminology, this kind of implicit content is referred to as *non-bona fide* true; these are implicit questionable contents that are not conveyed in good faith but are still non-explicitly understood as true within a given context. IMPAQTS defines four types of implicit content: *implicature*, *presupposition*, *topicalization*, and *vagueness*. For an extensive discussion on this theoretical framework, see [Lombardi Vallauri \(2016\)](#). Each sentence tagged as containing some implicit comment is accompanied by a *comment* written by the IMPAQTS corpus annotators explaining the meaning of the implicit content.³ Annotators' comments are all in the same, fixed format; e.g., for implicatures, they start with *it implies that... ;* for presuppositions, *it presupposes that... ;* for topicalizations, *it considers active in the discourse that... ;* for vagueness, *it leaves vague that... .*

Below, we report a single sample including the comment by the IMPAQTS annotators:

- (1) **Text with implicit content:** Italy doesn't need another government led by Monti. The last thing Italy needs is another government which is a slave to the banks!

Comment: It implies that Monti's government is a slave to the banks.

In this work, we focus on two of the four categories previously presented, i.e., *implicatures* and *presuppositions*. We chose not to include topical-

²<https://impaqts.dilef.unifi.it/>

³In the IMPAQTS corpus each sentence was tagged by 3 independent annotators. These annotations were then validated by a fourth expert who produced a final comment.

ization as its understanding and interpretation can heavily depend on the information carried out by its prosodic contour in the spoken language (Frascarelli and Hinterhölzl, 2008), which we cannot leverage when using transcriptions. On the other hand, vagueness was excluded due to its wider, less focused definition.

IMPAQTS only includes monologues, i.e., speeches produced by a single speaker with minimal or no interaction from others. In IMPAQTS, there are six types of such monologues: *Assembly speech*, *Rally speech*, *Party assembly speech*, *Statement in presence*, *Broadcast statement*, *New media statement*, and *Operational meeting Interview*. These types of monologues exhibit variation concerning register, planning, and historical influences between formats. Here, we chose to focus on sentences from Parliamentary and Rally speeches only; this choice was motivated by the fact that, in these types of speeches, there are fewer instances of implicit cases where what is implied or presupposed refers to elements that are external to the context of the speech. Therefore, the implicit content of these sentences can be understood by having access to the transcription of the monologue only.

Finally, the corpus includes speeches from three time periods: from 1946 to 1972, from 1973 to 1993, and from 1994 to 2022. We experiment with sentences from each of these periods.

3.2 How Much Linguistic Context Is Needed?

In the section above, we mentioned that the implicit content in the sentences from Parliamentary and Rally speeches can be understood by having access to the transcribed speech only. However, a crucial question immediately arises: *How much linguistic context* preceding the target sentence is needed to understand what is presupposed or implied in it? To empirically answer this question, we conducted a human validation study asking 9 expert linguist annotators to assess how many *left-hand* context sentences, i.e., preceding the target sentence, were needed to understand its implicit content. We selected 126 sentences, balanced with respect to both the type of implicit content (63 implicatures and 63 presuppositions) and time period (42 sentences for time range). We then created 3 surveys, each containing 42 samples embedding implicit content. In each survey, 3 experts were asked to read both the sentence containing the implicit content and its explanation given by the IMPAQTS' annotators.

Our experts had to assess whether the target sen-

tence provided enough information to understand the explanation in the comment, or if additional preceding linguistic context was needed. If the latter was the case, one extra sentence preceding the target sentence would appear, and the same question would be asked again. By design, our annotators could see up to five left-hand sentences. The annotation was carried out on LimeSurvey (LimeSurvey GmbH, 2012). See an example in Appendix C.1.

In 65% of cases, at least 2 out of 3 experts converged on the same number of sentences (more precisely, there was full consensus in 44% of cases and majority consensus in 21% of cases). We considered this agreement to be reasonable for the scope of this annotation and proceeded to examine the distribution of their answers. In 75% of cases, annotators judged the implicit content understandable based on the target sentence only. This percentage increases to 81% with three sentences and 88% with four sentences. As this percentage plateaus when adding even more sentences (see Fig. 3 for details), we empirically concluded that 4 preceding sentences are sufficient—while not always necessary—to understand the implicit content of (the vast majority of) the sentences in IMPAQTS.

3.3 Experimental Data

Based on the results of the annotations presented above, we proceeded to construct a dataset to be used in our experiments with LLMs. We consider all the samples in IMPAQTS containing either an implicature or a presupposition. For each of these samples, we retrieved the preceding 4 sentences in the speech. The resulting dataset, that we name IMPAQTS-PID—where *PID* stands for Preupposition and Implicature Dataset—includes 31,822 samples paired with the explanation of the implicit content from IMPAQTS; in particular, 14,932 samples embed an implicature and 16,890 a presupposition. In Appendix B we report some descriptive statistics of our dataset.

4 Methods

Using the IMPAQTS-PID dataset presented above, we challenge models to understand a sentence's implicit content. In this section, we describe the models we tested in our work, the experiments we performed to assess their abilities, and the experimental details common to both experiments.

4.1 Models

We experiment with four models pre-trained with multilingual data, which are therefore suitable to process Italian language: GPT4o-mini ([OpenAI, 2024](#)), Aya Expanse 8B ([CohereForAI, 2024](#)), LLAMA3.1 8B ([MetaAI, 2024a](#)), and LLAMA3.2 3B ([MetaAI, 2024b](#)). The first model is proprietary, the other three are open-weight. All these models are relatively small-scale, which has advantages in terms of computational efficiency, costs, and deployment feasibility. Thus, these models strike a good balance by requiring fewer resources, which reduces both hardware costs and energy consumption, while enabling faster inference. We test models' understanding of implicit content through two experiments, that we describe below.

4.2 Multiple-Choice Generation (MCG)

We perform an experiment where models are presented with four possible explanations IMPAQTS comments (hence, *explanations*) for a given implicit content and must pick the correct one. That is, we frame the problem as a multiple-choice generation (MCG) task. We feed the model with a brief instruction and the four candidates (A, B, C, and D) to choose from, only one of which is correct. To select challenging *distractors*, we use a method based on topic similarity. We first run a topic modeling analysis on all the explanations in IMPAQTS-PID. We identify 450 topics and a residual category including explanations that do not share a common topic. We then sample the distractors from the same topic class of the target explanation, being it one of the 450 topic-based categories or the residual one. This way, we ensure that the three distractors are incorrect yet plausible, compelling the models to reason over the input rather than relying on shortcuts. During inference, the order of the four candidates is randomized to avoid biases; that is, the position of the correct answer is roughly balanced across positions A-D. Given that each instance in our dataset includes a ground-truth answer, we evaluate model performance using plain accuracy.

4.3 Open-Ended Generation (OEG)

Framing the problem as an MCG task has various advantages, including being easy to set up and evaluate via accuracy. At the same time, it bears relevant limitations due to the selection of the distractors and the non-naturalistic setting of choosing from available options. In this experiment, we cir-

cumvent these limitations and directly query models to generate an explanation for a given implicit content; i.e., we frame the problem as an open-ended generation (OEG) task. We experiment with three prompting settings, which we briefly describe below:

Zero-shot The model is queried without any prior examples or detailed guidance and solely relies on its pre-training knowledge. This is the template used for this prompting technique:

*Explain the implicit content of the following text.
Consider that it appears in its right-most sentence.
Text: [sentence embedding implicit content]
Implicit Content: _____*

Few-Shot The prompt includes four illustrative examples, i.e., two implicatures and two presuppositions, including both the text embedding the implicit content and its explanation. This is the template used for this prompting technique (for brevity, we only report two *shots* instead of four):

*Text: [First example text]
Implicit Content: [Explanation of the implicit content in the first example text]*

...

*Text: [Fourth example text]
Implicit Content: [Explanation of the implicit content in the fourth example text]
Text: [sentence embedding implicit content]
Implicit Content: _____*

Chain-of-Thought (CoT) The prompt presents an extensive step-by-step explanation outlining the reasoning process necessary to decode the implicit content. We adapt these steps from the detailed instructions given to the IMPAQTS corpus annotators. This is the template used for this prompting technique (shortened for space reasons):

Implicatures arise in communication whenever the speaker "challenges" one of the four conversational maxims derived from the well-known Cooperative Principle ...

We define "presupposition" as any content that is taken for granted by the participants in communication and, more specifically, content conveyed as part of the knowledge already shared by the interlocutor. ...

Consider what has just been said about implicatures and presuppositions, and explain, proceeding

469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
step by step, what the implicit content in the
following text is. Consider that it appears in its
right-most sentence.

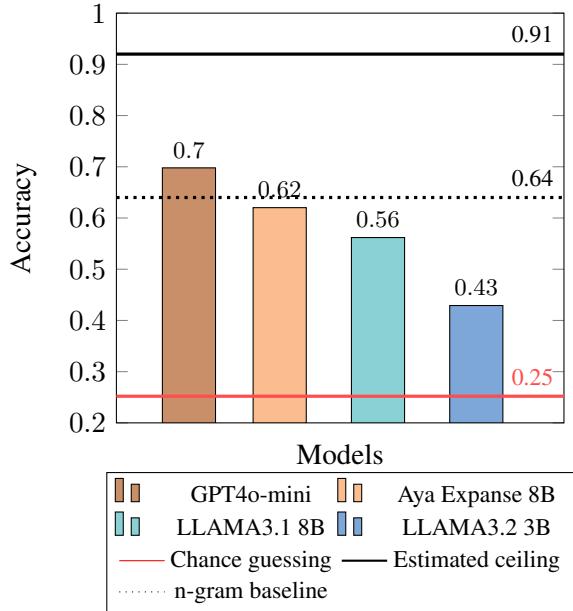
493 Text: [sentence embedding implicit content]
494 Implicit Content: _____

495 **Evaluation** In contrast to the MCG task, where
496 the presence of a ground-truth answer allows for
497 an evaluation based on plain accuracy, assessing
498 model performance here is more complex. Evaluation
499 based on automatic similarity measures or
500 LLMs-as-judges are potential solutions. However,
501 the novelty of the task introduces uncertainty re-
502 garding their reliability. Given these concerns, we
503 opt for a small-scale but more robust assessment
504 based on human expert evaluation. We devise an
505 evaluation protocol where expert linguists assess
506 the texts containing implicit content, the comment
507 about the implicit content from IMPAQTS, and
508 the implicit content explanation generated by the
509 model. See Appendix C.2 for an example. The
510 experts are asked to judge the quality of the ex-
511 planation generated by a model by choosing one
512 between 5 possible options:

- 513 • **Totally correct:** the explanation is short and
514 focused on the same topic highlighted in the
515 original comment from IMPAQTS.
- 516 • **Correct among various options:** the answer
517 lists many possible explanations; among those,
518 the correct one is mentioned (the correctness
519 of other explanations is not relevant).
- 520 • **Partially correct:** the model outputs a short
521 answer that does not fully address the point
522 but captures some elements of the correct ex-
523 planation.
- 524 • **Totally wrong:** a short or long answer that is
525 wrong, i.e., is not about the topic highlighted
526 by the comment from IMPAQTS.
- 527 • **Answer not given:** the model refuses to an-
528 swer, avoids the question or just highlights
529 the sentence containing the implicit content
530 without giving an explanation.

531 4.4 Experimental Details

532 Details on prompting and decoding settings for
533 both tasks, including token limits, temperature pa-
534 rameters, GPUs, and computational time are pro-
535 vided in Appendix F and G.



516 Figure 1: MCG task. Model accuracy.
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543

5 Results

5.1 MCG Results

516 Models are tasked to pick one among four possible
517 answers. Therefore, we can measure model accu-
518 racy and compute a random chance level (25%)
519 and an estimated ceiling accuracy (91%) boot-
520 strapped from the evaluation presented in Section
521 3.2.⁴ Moreover, we compare the results of the
522 tested models against an n -gram similarity baseline
523 based on BLEU-4 (Papineni et al., 2002). This
524 baseline picks the answer with the highest BLEU-
525 4 overlap with the target sentence embedding the
526 implicit content. If a model systematically chooses
527 the answer that overlaps the most with the target
528 sentence, then it should achieve a 64% accuracy.

529 As shown in Figure 1, all models perform well
530 above chance, indicating that they understand the
531 task and can sometimes select the correct expla-
532 nation. At the same time, we observe significant
533 variation between models, with the best-performing
534 GPT4o-mini achieving an accuracy of 70%, nearly
535 30 percentage points higher than LLAMA3.2 3B
536 (43%). Interestingly, GPT4o-mini is the only
537 model that consistently surpasses the n -gram base-
538 line, while all other models perform worse than this
539 heuristic. On one hand, this indicates that LLMs do
540 not systematically rely on a strategy based solely on
541 shallow similarity between the input and the answer.

542 ⁴The bootstrapped ceiling accuracy is 88%, to which we
543 add an extra 3% standing for a random chance level accuracy
544 in the remaining samples.

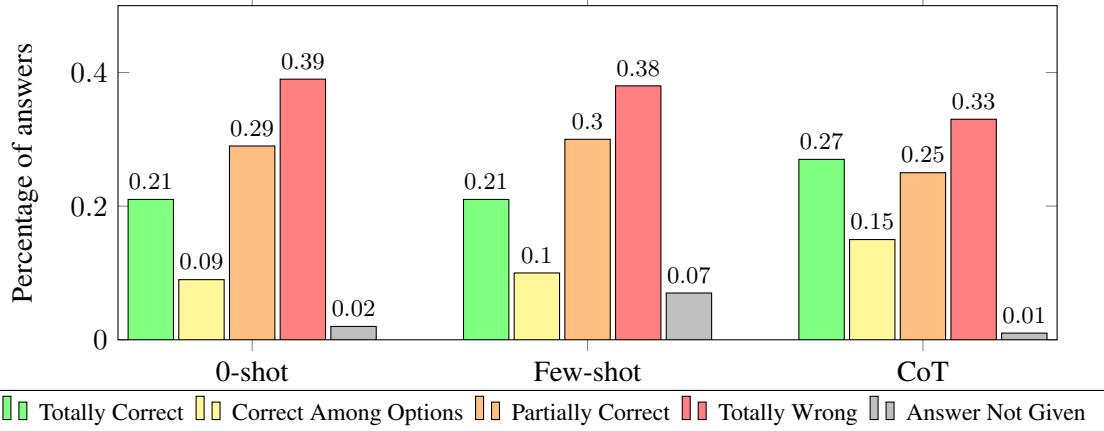


Figure 2: OEG task. Human expert evaluation of GPT4o-mini generated answers. Results refer to 150 samples.

On the other hand, it suggests that GPT4o-mini possesses certain pragmatic abilities that enable it to understand the implicit content and go beyond surface-level similarities in the prompt. Nonetheless, even this best-performing model still falls more than 20 percentage points short of the estimated ceiling accuracy. Therefore, we conclude that overall model performance on this task remains largely unsatisfactory.

Analysis of distractors In a multiple-choice task, the selection of distractors can significantly impact the models’ final performance. To assess whether this affects our results, we calculate model accuracy for two subsets of our data: one where the distractors belong to the same topic category as the correct explanation (hard-negatives) and another where both the target explanation and the distractors come from a residual category with no clear topic (easy-negatives). Figures 4 and 5 in Appendix D report the accuracy results in these subsets. Overall, we find that the accuracy of all tested models is consistently higher in the easy-negatives subset compared to the hard-negatives subset. For instance, GPT4o-mini shows an approximately 8% performance difference between the two subsets, i.e., from 73% to 65%. This confirms that distractor selection can significantly impact model performance, introducing confounding factors and limiting result robustness. To mitigate these issues, we conduct an open-ended generation (OEG) task, the results of which are reported below.

5.2 OEG Results

Based on the results of the MCG task and to minimize computational effort and manual labor, we focus solely on the best-performing GPT4o-mini

model. We evaluate 150 samples per setting—zero-shot, few-shot, and CoT—and ask 10 expert linguists to assess 45 sentences each, following the guidelines outlined in Section 4.3. Each sample is thus evaluated by a single expert annotator. The results of this evaluation are reported in Figure 2.⁵ Several notable observations can be made.

Limited pragmatic abilities The model produces answers judged as *totally correct* in only a small percentage of cases, ranging from 21% in the zero-shot to 27% in the CoT setting. This means that even under the best CoT conditions, the model is fully correct in only about one-fourth of the samples. When considering *partially correct* and *correct among options* responses, this proportion increases to roughly two-thirds of the samples. While this represents an improvement, the performance remains unsatisfactory, as the model is *totally wrong* or refuses to provide an answer in about one-third of the cases. These results confirm the challenges that pragmatic reasoning presents to LLMs, as observed in the previous experiment. While those difficulties could have been influenced by confounds due to the design, in this case, the model faces no restrictions except for the type and amount of information provided in the prompt.

CoT boosts performance The use of the CoT strategy leads to a noticeable improvement in model performance. Specifically, the proportion of *totally wrong* answers in the zero-shot setting

⁵In Figure 6 in Appendix G.1, we report the results of the same evaluation performed using an ‘LLMs-as-Judges’ approach (Bavaresco et al., 2024). Using GPT4o (a bigger and more capable model than GPT4o-mini), we report an even more negative assessment of model performance compared to human experts, judging around 60% of cases as *totally wrong*.

Example 1's IMPAQTS comment	Example 2's IMPAQTS comment
The left-wing government caused Italy's loss of image.	The Prime Minister wrongly downplayed abstention in Emilia-Romagna.
0-shot-Totally wrong answer	0-shot-Totally wrong answer
The text critiques Italy's political crisis, arguing that the issue is not just tax pressure or security but a deeper crisis of trust. Italians feel disconnected from national identity, unprotected by leaders, and this loss of confidence affects Italy's global image. The author suggests the crisis stems from poor governance and calls for collective reflection and renewal.	The text criticizes the law and voter turnout decline, expressing concern over the government's attitude. The Emilia-Romagna example highlights the severity of the issue, implying that the Prime Minister's dismissal is superficial. The author argues that reforms are urgently needed.
CoT-Totally correct answer	CoT-Totally correct answer
The text blames politicians for the crisis of confidence and Italy's declining image . It suggests that beyond tax pressure and security, there is a deeper fracture in national unity, implying the need for leadership change to restore trust.	The text strongly disagrees with the law and criticizes the Prime Minister's stance , suggesting that low voter turnout is a serious issue that should not be underestimated.

Table 1: Generated answers by GPT4o-mini in the OEG task. In both examples, the model is *totally wrong* in the zero-shot setting but *totally correct* in the CoT setting.

decreases significantly (from 39% to 33%). Meanwhile, *totally correct* answers increase, emphasizing that CoT is an effective strategy for enhancing the model’s ability to understand implicit content. In contrast, providing the models with a few detailed examples, as in the few-shot strategy, does not lead to any significant improvement.

Qualitative analysis As a qualitative analysis, we cherry-picked two cases highlighting the behavior of the model, that we report in Table 1. In both examples, the model provided a *totally wrong* answer in the zero-shot setting but a *totally correct* one when using CoT. Below, we briefly analyze some key properties of these examples and compare the behavior of the model in the two settings.

In **Example 1**, when the model is prompted in a zero-shot setting, it fails to retrieve the information that identifies the responsibility for Italy’s loss of image as being with the politicians rather than the political crisis itself. Conversely, when prompted with CoT reasoning, the model accurately captures this information, although it fails to attribute responsibility to the left-wing government.

In **Example 2**, the model correctly identifies the specific target of implicitness in both zero-shot and CoT settings, i.e. the Prime Minister stance on abstentionism, but, in the former, fails to explain that what is implied is that the Prime Minister’s stance is not just *superficial* but wrong.

These differences suggest that, beyond the linguistic complexities of implicit language addressed through CoT reasoning, models could benefit from additional contextual, non-linguistic information, including the historical period in which the speech

was delivered, its speaker, and the individual she or he is referring to.

6 Conclusions

In this paper, we explored the ability of LLMs to explain implicit content in political discourse, focusing on implicatures and presuppositions. Leveraging, for the first time in the domain of NLP, the large-scale IMPAQTS corpus containing Italian political discourse, we evaluated several SotA LLMs through two experiments. We showed that while model performance is not random, there is a significant gap between their outputs and human expert intuitions, both with predefined options and free-form explanations. We also demonstrated that CoT prompting can improve reasoning, confirming observations from previous work in the domain of pragmatic understanding (Kim et al., 2023). Based on our findings, we propose that incorporating more world knowledge—such as information about politicians and their parties—could help models better understand implicit content, a direction we plan to explore in future work.

7 Limitations

Despite the promising results of our study, limitations must be acknowledged. Our evaluation primarily relies on manual annotation of a relatively small subset of data done by a small group of experts. In particular, open-ended responses are evaluated on a subset that may not fully capture the complexity of models outputs. Additionally, LLMs’ interpretations of implicit content remain influenced by pre-training biases or post-training fine-

tuning strategies as evidenced by the positional bias and the refusal patterns observed in models like LLAMA3.1 8B and LLAMA3.2 3B when processing politically sensitive topics. Finally, our study does not assess the models' robustness across different political contexts or their ability to adapt to real-time discourse shifts, which are crucial for practical applications. Future research should explore and develop more rigorous evaluation frameworks and investigate how fine-tuning or domain-specific adaptations can enhance LLMs' implicit content understanding. Moreover, a parallel study on human understanding and explanation capabilities of implicit content (in political discourse) should be developed to draw a comparison between human and machine reasoning on this topic.

References

- Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, Esam Ghaleb, Mario Julianelli, Michael Hanna, Alexander Koller, et al. 2024. Llms instead of human judges? a large scale empirical study across 20 nlp evaluation tasks. *arXiv preprint arXiv:2406.18403*.
- CohereForAI. 2024. *Ayaexpanse: Combining research breakthroughs for a new multilingual frontier*. Accessed: 2025-01-22.
- Federica Cominetti, Lorenzo Gregori, Edoardo Lombardi Vallauri, and Alessandro Panunzi. 2024. *IMPAQTS: a multimodal corpus of parliamentary and other political speeches in Italy (1946-2023), annotated with implicit strategies*. In *Proceedings of the IV Workshop on Creating, Analysing, and Increasing Accessibility of Parliamentary Corpora (ParlaCLARIN) @ LREC-COLING 2024*, pages 101–109, Torino, Italia. ELRA and ICCL.
- Oswald Ducrot. 1972. Dire et ne pas dire: principes de sémantique linguistique.
- Mara Frascarelli and Roland Hinterhölzl. 2008. Types of topics in german and italian. In *On information structure, meaning and form: Generalizations across languages*, pages 87–116. John Benjamins Publishing Company.
- Richard Garner. 1971. 'presupposition'in philosophy and linguistics.
- H Paul Grice. 1990. 1975 logic and conversation. *The Philosophy of Language*.
- Jennifer Hu, Sammy Floyd, Olessia Jouravley, Evelina Fedorenko, and Edward Gibson. 2023. A fine-grained comparison of pragmatic language understanding in humans and language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4194–4213, Toronto, Canada. Association for Computational Linguistics.
- Pere-Lluís Huguet Cabot, Verna Dankers, David Abadi, Agneta Fischer, and Ekaterina Shutova. 2020. *The Pragmatics behind Politics: Modelling Metaphor, Framing and Emotion in Political Discourse*. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4479–4488, Online. Association for Computational Linguistics.
- Paloma Jeretic, Alex Warstadt, Suvrat Bhooshan, and Adina Williams. 2020. *Are natural language inference models IMPPRESSive? Learning IMPlicature and PRESupposition*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8690–8705, Online. Association for Computational Linguistics.
- Paritosh D Katre. 2019. Nlp based text analytics and visualization of political speeches. *International Journal of Recent Technology and Engineering*, 8(3):8574–8579.
- Zae Myung Kim, David E Taylor, and Dongyeop Kang. 2023. " is the pope catholic?" applying chain-of-thought reasoning to understanding conversational implicatures. *arXiv preprint arXiv:2305.13826*.
- Lincan Li, Jiaqi Li, Catherine Chen, Fred Gui, Hongjia Yang, Chenxiao Yu, Zhengguang Wang, Jianing Cai, Junlong Aaron Zhou, Bolin Shen, et al. 2024. Political-ilm: Large language models in political science. *arXiv preprint arXiv:2412.06864*.
- LimeSurvey GmbH. 2012. *LimeSurvey: An Open Source survey tool*. LimeSurvey GmbH, Hamburg, Germany.
- Edoardo Lombardi Vallauri. 2016. The “exaptation” of linguistic implicit strategies. *SpringerPlus*, 5(1):1106.
- Edoardo Lombardi Vallauri. 2019. *La lingua disonesta: contenuti impliciti e strategie di persuasione*. Intersezioni (Il Mulino). Il Mulino.
- Edoardo Lombardi Vallauri and Viviana Masia. 2020. La comunicazione implicita come dimensione di variazione tra tipi testuali. In *Linguaggi settoriali e specialistici. Sincronia, diacronia, traduzione, variazione (Proceedings of the International SILFI Conference 2018)*, pages 113–120. Cesati Firenze.
- Giada Marino and Fabio Giglietto. 2024. Integrating large language models in political discourse studies on social media: Challenges of validating an llms-in-the-loop pipeline. *Sociologica*, 18(2):87–107.
- MetaAI. 2024a. *Llama 3.1: A collection of multilingual large language models*. Accessed: 2025-01-22.
- MetaAI. 2024b. *Llama 3.2: Advancements in vision and edge ai*. Accessed: 2025-01-22.

780	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.	832
781		833
782		834
783		835
784		836
785		837
786		838
787		839
788		
789		
790		
791	Patrick Morency, Steve Oswald, and Louis De Saussure. 2008. Explicitness, implicitness and commitment attribution: A cognitive pragmatic approach. <i>Belgian journal of linguistics</i> , 22(1):197–219.	840
792		841
793		842
794		843
795	Renáta Németh. 2023. A scoping review on the use of natural language processing in research on political polarization: trends and research prospects. <i>Journal of computational social science</i> , 6(1):289–313.	844
796		845
797		846
798		847
799		848
800	OpenAI. 2024. Gpt-4o mini: Advancing cost-efficient intelligence . Accessed: 2025-01-22.	849
801		850
802		851
803	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	852
804		853
805		854
806		
807	Laura Ruis, Akbir Khan, Stella Biderman, Sara Hooker, Tim Rocktäschel, and Edward Grefenstette. 2024. The goldilocks of pragmatic understanding: Fine-tuning strategy matters for implicature resolution by LMs. <i>Advances in Neural Information Processing Systems</i> , 36.	855
808		
809		
810	Dan Sperber and Deirdre Wilson. 1987. Précis of relevance: Communication and cognition. <i>Behavioral and brain sciences</i> , 10(4):697–710.	856
811		
812		
813	Settaluri Sravanthi, Meet Doshi, Pavan Tankala, Rudra Murthy, Raj Dabre, and Pushpak Bhattacharyya. 2024. PUB: A pragmatics understanding benchmark for assessing LLMs' pragmatics capabilities . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 12075–12097, Bangkok, Thailand. Association for Computational Linguistics.	857
814		858
815		859
816		860
817		861
818		862
819		863
820	PF Strawson. 1964. Identifying reference? and truth-values. <i>Theoria</i> , 30.	864
821		865
822	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	866
823		867
824		868
825		869
826		870
827		871
828	Teun A Van Dijk. 1992. Discourse and the denial of racism. <i>Discourse & society</i> , 3(1):87–118.	872
829		873
830	Teun A Van Dijk et al. 1997. What is political discourse analysis. <i>Belgian journal of linguistics</i> , 11(1):11–52.	874
831		
832	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	875
833		876
834		877
835		878
836		
837		
838		
839		
840	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	879
841		880
842		881
843		882
844		883
845		
846		
847		
848		
849	Zilong Zheng, Shuwen Qiu, Lifeng Fan, Yixin Zhu, and Song-Chun Zhu. 2021. Grice: A grammar-based dataset for recovering implicature and conversational reasoning. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 2074–2085.	884
850		885
851		886
852		887
853		888
854		889
855	A Examples Appendix	890
856	A.1 Paper examples	
857	They pretend not to understand	
858	Italian: Cioè è una persona che ha usato la carica pubblica per farsi gli affari propri. Ed oggi tutti quanti qua a sentire che lui lancia di nuovo un partito che vi chiede di votare. Ecco, io mi rivolgo a voi adesso. Eh no, eh no. Fino a che punto <i>volete far finta di non capire?</i>	891
859		892
860		893
861		894
862		895
863		896
864		897
865	English: So, [he's] a person that used its public office to do his own business. And today, we are all here listening to him launching a new party and asking for your votes. Well, I am talking to you now. Oh no, oh no. How far will you <i>want to pretend not to understand?</i>	898
866		899
867		900
868		901
869		902
870		903
871		904
872	Sample with implicit content and its IMPAQTS comment	905
873	Italian: In Italia non serve un altro governo Monti. L'ultima cosa di cui l'Italia ha bisogno è un altro governo asservito alle banche!	906
874	Comment: Implica che il governo monti sia asservito alle banche.	907
875		908
876		909
877		910
878		911
879	B IMPAQTS-PID descriptive statistics	912
880	The IMPAQTS-PID dataset consists of 5% samples from the 1946-1972 period, 23% from 1972-1994, and 72% from 1994-2023. The Male:Female speech ratio is 5:1. Regarding the political beliefs	913
881		914
882		915
883		916

Annotated sentences	# of sentences
Total	31,822
Implicatures	14,932
Presuppositions	16,890

Diachronic context	% of sentences
1946–1972	5%
1972–1994	23%
1994–2023	72%

Table 2: Statistical data of IMPAQTS-PID datapoints

Parameter	Avg. Length	STD
Text with implicit content	611,67	283,55
Tagged sentence	86,59	68,95
IMPAQTS Comment	76,72	42,30

Table 3: Avarage lenght and Standard deviation (in tokens) of various IMPAQTS-PID datapoints.

represented, the Center-left and Center-right are predominant, covering a total of 45% of the samples. The other four parties are less represented but fairly balanced, each comprising 13-14% of the samples. Table 2 and 3 report some descriptive statistics of our dataset.

C Human validation studies

C.1 Needed Context

Instructions

Italian: Questa survey è composta da 42 domande. Ogni domanda presenta un enunciato (in corsivo e tra virgolette) e un contenuto implicito, indicato con l'etichetta "Contenuto Implicito", come nell'esempio seguente:

Enunciato: "È così che in fondo abbiamo agito anche in quest'ultima crisi. Gli italiani e i tedeschi, e tanti altri paesi."

Contenuto Implicito: L'Italia ha agito nel modo descritto anche durante altre crisi.

Ti chiediamo di valutare se il contesto fornito nell'enunciato è sufficiente a inferire il contenuto implicito indicato. Non ti chiediamo se

sei in grado di identificare il contenuto implicito, ma se, una volta chiarito quale sia, ritieni che possa essere dedotto dal contesto dell'enunciato.

Le risposte possibili sono due: "Il contesto è sufficiente per inferire il contenuto implicito" oppure "Serve avere più contesto." Se rispondi "Il contesto è sufficiente per inferire il contenuto implicito", il contesto non sarà ulteriormente espanso e registreremo la tua risposta come definitiva. Se invece rispondi "Serve avere più contesto" verrà visualizzato un nuovo enunciato a sinistra, e potrai scegliere nuovamente tra le stesse due opzioni. Questo processo può ripetersi fino a cinque volte, con un massimo di cinque enunciati aggiuntivi per ampliare il contesto. Ad ogni passaggio, potrai scegliere che risposta dare. Al termine delle cinque espansioni, se rispondi ancora "Serve avere più contesto", non comparirà nessun altro enunciato a sinistra, e registreremo la tua risposta come definitiva.

Una volta data la risposta definitiva, per procedere alla domanda successiva devi cliccare su "Next" in basso a destra.

Nel caso dell'esempio che ti abbiamo presentato sopra, il contenuto implicito può essere dedotto dal contesto presentato in quanto a suggerire che anche in altre crisi l'Italia abbia agito nel modo descritto è l'avverbio anche presente nell'enunciato. Quindi puoi rispondere "Il contesto è sufficiente per inferire il contenuto implicito" e andare avanti come descritto precedentemente. Ogni domanda richiede obbligatoriamente una risposta.

Il tempo di compilazione di questo questionario varia da persona a persona e da domanda a domanda ma dovresti cavartela nel giro di un'ora. Grazie ancora per il tempo dedicato a questa indagine!

English: This survey consists of 42 questions. Each question presents a statement (in italics and quotation marks) and an implicit content indicated with the label "Implicit Content," as in the following example:

Statement: "This is how we have essentially acted even in this latest crisis. Italians and Germans, and many other countries."

Implicit Content: Italy acted in the described manner even during other crises.

We ask you to judge whether the context provided in the statement is sufficient to infer the

960 indicated implicit content. We are not asking if
961 you can identify the implicit content, but whether,
962 once clarified what it is, you believe it can be
963 inferred from the context of the statement.

964 There are two possible responses: "The context is
965 sufficient to infer the implicit content" or "More
966 context is needed."

967 If you respond, "The context is sufficient to infer
968 the implicit content," the context will not be further
969 expanded, and we will record your response as
970 final.

971 If instead you respond, "More context is needed,"
972 a new statement will appear on the left, and you
973 will be able to choose again between the same
974 two options. This process can repeat up to five
975 times, with a maximum of five additional left-hand
976 sentences to expand the context. At each step, you
977 can choose your response. At the end of the five
978 expansions, if you still respond, "More context is
979 needed," no further left-hand sentences will appear,
980 and we will record your response as final.

981 Once the final response is given, to proceed to the
982 next question, click on the "Next" button at the
983 bottom right.

984 In the case of the example presented above, the
985 implicit content can be inferred from the presented
986 context because the adverb "also" in the statement
987 suggests that Italy acted in the described manner
988 during other crises as well. Therefore, you can
989 respond, "The context is sufficient to infer the
990 implicit content," and proceed as previously
991 described. Each question requires an answer.

992 The completion time of this questionnaire varies
993 from person to person and from question to
994 question, but you should be able to complete it
995 within an hour.

996 Thank you again for the time dedicated to this
997 survey!

998 Example

1000 **Italian: Enunciato:** "Un'opera di demonizzazione
1001 che non ce la fa a buttar giù Berlusconi, neanche
1002 usando della cattiva stampa, dei cattivi giornali,
1003 pensate a Repubblica, ma non solo, al Corriere."

1004 **Contenuto Implicito:** L'opera di demonizzazione
1005 ha tentato di buttar giù Berlusconi.

- 1006 1. Il contesto è sufficiente per inferire il con-
1007 tenuto隐含的.

- 1008 2. Serve avere più contesto.

1009 Se Serve avere più contesto è l'opzione selezion-

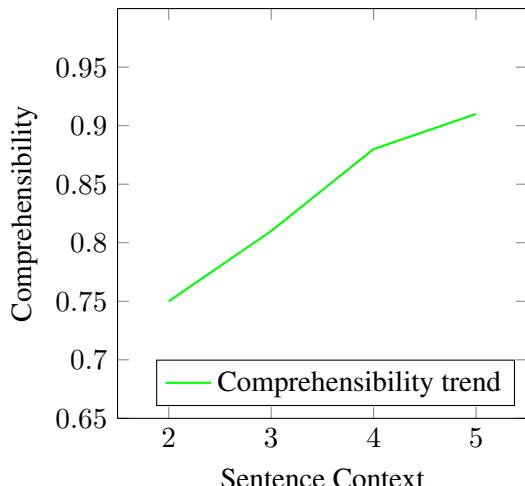


Figure 3: Trend of implicit content comprehensibility as the left-hand context increases

ata, un'altra domanda con contesto sinistro espanso (evidenziata nell'esempio seguente) appare.

Enunciato: "E scatta una storia di democrazia che riguarda il presidente Berlusconi, che riguarda il suo partito, che riguarda la sua gente. Un'opera di demonizzazione che non ce la fa a buttar giù Berlusconi, neanche usando della cattiva stampa, dei cattivi giornali, pensate a Repubblica, ma non solo, al Corriere."

Contenuto Implicito: L'opera di demonizzazione ha tentato di buttar giù Berlusconi.

1. Il contesto è sufficiente per inferire il con-
tenuto隐含的.
2. Serve avere più contesto.

English: Statement: A demonization effort that fails to take down Berlusconi, not even by using bad press, bad newspapers, think of *Repubblica*, but not only, also *Corriere*.

Implied content: The demonization effort tried to take down Berlusconi.

1. The context is sufficient to infer the implicit content.
2. More context is needed.

If *More context is needed* is the chosen option, another question with extended context (highlighted in the following example) appears.

Text: And so begins a demonization campaign targeting President Berlusconi, his party,

1040
1041
1042
1043
1044
1045
1046
and his supporters. A demonization effort that
fails to take down Berlusconi, not even by using
bad press, bad newspapers, think of Repubblica,
but not only, also Corriere.

1047
1048
Implied content: The demonization effort tried to
take down Berlusconi.

- 1049
1050
1. The context is sufficient to infer the implicit
content.
2. More context is needed.

1051 C.2 OEG Evaluation

1052 Instructions

1053 **Italian:** Grazie per il tempo che stai dedicando a
questa campagna di valutazione.

1054 Per rispondere al sondaggio, ti preghiamo di leg-
1055 gere attentamente le istruzioni che trovi qui sotto.

1056 Ti verrà chiesto di valutare testi generati da in-
1057 telligenza artificiale che esplicitano un contenuto
1058 implicito. Ogni schermata è composta da:

1059 Un testo estratto da discorsi politici italiani;
1060 L'annotazione di un esperto che ne esplicita il
1061 contenuto隐含; L'output generato dal modello.
1062 Ognuno di questi elementi è introdotto, rispettiva-
1063 mente, dalle etichette *Testo*, *Annotazione umana* e
1064 *Output*. I chiediamo di valutare se il modello ha
1065 colto la sfumatura sottesa dall'annotatore umano,
1066 non se l'output generato sia sensato.

1067 I testi sono stati trascritti mantenendo il più pos-
1068 sibile la loro forma originale. Può capitare, infatti,
1069 che ci siano ripetizioni di parole, frasi abbandonate
1070 o incoerenze sintattiche. Potrai selezionare tra 5
1071 possibili annotazioni, aventi 5 diversi significati e
1072 valori:

- 1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1. Totalmente corretto: l'output generato dal
modello esplica un solo contenuto implicito, e tale
contenuto implicito è quello presente nell'annotazione. Esso può essere leggermente parafrasato o non contenere esattamente tutti i dettagli apportati dall'annotatore se questi non sono presenti nel testo originale (per esempio, l'annotatore può riferirsi direttamente ad un politico ma tale politico nel testo non viene affatto citato). Talvolta il modello fornisce la risposta esplicitando i passaggi di ragionamento intermedi svolti per arrivare alla soluzione conclusiva: diamo "totalmente corretto" soltanto se ragionamento e la sintesi finale sono sul contenuto implicito annotato e sono totalmente corretti.

1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
2. Corretto tra varie opzioni: l'output generato dal modello (i) esplica multipli contenuti impliciti oppure lo stesso contenuto implicito con multiple interpretazioni e una di queste è totalmente corretta; oppure, (ii) se output generato dal modello, esplicitando nella risposta i passaggi di ragionamento intermedi svolti per arrivare alla soluzione conclusiva, non esplica il contenuto implicito corretto nella sintesi finale ma lo menziona in uno dei passaggi intermedi.

1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
3. Parzialmente corretto: l'output generato dal modello (i) esplica un solo contenuto implicito, e tale contenuto implicito è parzialmente in linea con quello annotato dall'annotatore umano ovvero, ci può essere qualche interpretazione diversa o sbagliata con parte dell'esplicitazione corretta; oppure (ii) nella sintesi finale o in uno dei passaggi intermedi del ragionamento che il modello esplica c'è una spiegazione del contenuto implicito annotato e questa è parzialmente corretta.

1111
1112
1113
4. Totalmente sbagliato: l'output generato dal modello è diverso da quanto annotato dell'annotatore umano.

1114
1115
1116
1117
5. Risposta non fornita: L'output generato dal modello non contiene alcuna risposta, è una perifrasi che aggira la domanda o è un vero e proprio rifiuto di generare una risposta.

Ecco un esempio di schermata che troverai.

1118
1119
1120
1121
Testo: Diseguaglianza significa un lavoro precario, incerto, sottopagato, sfruttato, perfino gratuito, un lavoro povero. Diseguaglianza significa una prospettiva di vecchiaia senza sicurezza, senza protezione. Liberi e Uguali è in campo per battersi contro tutto questo. Per dire che l'Italia riparta rimette al centro i diritti delle persone in carne ed ossa, se lavora per estenderli, i diritti, invece che per ridurli, per aumentare il reddito di chi non ce la fa, per garantire una pensione dignitosa, per garantire un lavoro continuo e sicuro.

1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
Annotatore umano: L'Italia è ferma.

1133
1134
1135
1136
1137
Output: Il contenuto implicito del testo evidenzia una forte critica alle attuali condizioni lavorative e sociali in Italia, suggerendo che la diseguaglianza non solo compromette la qualità della vita presente, ma crea anche incertezze

per il futuro. Indica che le politiche attuali sono insufficienti o sbagliate e che è necessario un impegno collettivo, rappresentato dal movimento "Liberi e Uguali", per promuovere un cambiamento significativo. Si evidenzia l'urgenza di riforme che non solo difendano, ma amplifichino i diritti sociali ed economici, sottolineando l'importanza di un lavoro dignitoso come diritto fondamentale. Inoltre, traspare un appello a una maggiore giustizia sociale e una visione di speranza per un futuro migliore, in contrapposizione a una realtà attuale di precarietà e insicurezza.

1. Totalmente corretto
2. Corretto tra opzioni
3. Parzialmente corretto
4. Totalmente sbagliato
5. Risposta non fornita

Siccome ciò che deve essere valutato non è se l'output del modello ci appaia sensato ma se il modello ha colto la sfumatura sottesa dall'annotatore umano, in questo caso, benché ciò che il modello genera sia plausibile, quello che l'annotatore umano ha evidenziato è assente, dunque l'opzione da selezionare per questa schermata sarà totalmente sbagliato.

Una volta data la risposta, clicca su "Avanti" in basso a destra per procedere al successivo testo da annotare. Puoi tornare all'item precedente cliccando su "Indietro" in basso a sinistra. Ogni schermata richiede obbligatoriamente una risposta.

Il tempo di compilazione del questionario varia da persona a persona e da domanda a domanda.

Grazie ancora per il tempo dedicato a questa campagna.

English: Thank you for the time you are dedicating to this evaluation campaign.

To respond to the survey, please read the instructions below carefully.

You will be asked to evaluate texts generated by artificial intelligence that make implicit content explicit. Each screen consists of:

A text extracted from Italian political speeches;

An expert annotation that makes the implicit content explicit;

The output generated by the model.

Each of these elements is introduced by the labels *Text*, *Human Annotation*, and *Output*, respectively. We ask you to evaluate whether the model

has captured the nuance implied by the human annotator, not whether the generated output makes sense.

The texts have been transcribed while maintaining their original form as much as possible. Therefore, you might encounter word repetitions, abandoned phrases, or syntactic inconsistencies. You can select from 5 possible annotations, each having 5 different meanings and values:

1. Totally Correct: The output generated by the model makes explicit only one implicit content, and that implicit content is the one present in the annotation. It can be slightly paraphrased or not contain exactly all the details provided by the annotator if these are not present in the original text (for example, the annotator might directly refer to a politician, but that politician is not mentioned at all in the text). Sometimes the model provides the response by making explicit the intermediate reasoning steps carried out to reach the final solution: we give "totally correct" only if both the reasoning and the final synthesis are on the annotated implicit content and are completely correct.

2. Correct among various options: The output generated by the model (i) makes explicit multiple implicit contents or the same implicit content with multiple interpretations, and one of these is totally correct; or (ii) the output generated by the model, by making explicit the intermediate reasoning steps to reach the final solution, does not make explicit the correct implicit content in the final synthesis but mentions it in one of the intermediate steps.

3. Partially Correct: The output generated by the model (i) makes explicit only one implicit content, and that implicit content is partially in line with what the human annotator annotated, that is, there may be some different or wrong interpretations with part of the explication correct; or (ii) in the final synthesis or in one of the intermediate reasoning steps that the model makes explicit, there is an explanation of the annotated implicit content that is partially correct.

4. Totally Wrong: The output generated by the model is different from what the human annotator annotated.

1236
1237
1238
1239

5. **Answer not given:** The output generated by
the model contains no answer, is a circumlo-
cation that avoids the question, or is a true
refusal to generate a response.

1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269

Here is an example of a screen you will find.

Text: Inequality means precarious, uncertain,
underpaid, exploited, even unpaid work, poor work.
Inequality means an old age perspective without
security, without protection. Liberi e Uguali is
on the field to fight against all this. To say that
Italy restarts by putting at the center the rights of
real people, if it works to extend them, the rights,
instead of reducing them, to increase the income of
those who can't make it, to guarantee a dignified
pension, to guarantee continuous and secure work.

Human Annotation: Italy is stalled.

Output: The implicit content of the text highlights
a strong critique of the current working and social
conditions in Italy, suggesting that inequality
not only compromises the quality of present life
but also creates uncertainties for the future. It
indicates that current policies are insufficient or
wrong and that a collective effort, represented
by the "Liberi e Uguali" movement, is needed to
promote significant change. It underscores the
urgency of reforms that not only defend but also
amplify social and economic rights, emphasizing
the importance of dignified work as a fundamental
right. Moreover, it reveals a call for greater social
justice and a vision of hope for a better future, in
contrast to the current reality of precariousness and
insecurity.

- 1270
1271
1272
1273
1274
1. Totally Correct
 2. Correct among options
 3. Partially Correct
 4. Totally wrong
 5. Answer not given

1275
1276
1277
1278
1279
1280
1281
1282
1283

Since what must be evaluated is not whether the
model's output seems sensible to us but whether
the model has captured the nuance implied by the
human annotator, in this case, although what the
model generates is plausible, what the human anno-
tator highlighted is absent, so the option to select
for this screen will be totally wrong.

Once you have given the answer, click "Next" at
the bottom right to proceed to the next text to be

1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329

annotated. You can go back to the previous item
by clicking "Back" at the bottom left. Each screen
requires an answer.

The time to complete the questionnaire varies
from person to person and from question to ques-
tion.

Thank you again for the time dedicated to this
campaign.

D MCG details

D.1 MCG prompt

Italian: Considera il seguente testo: *Noi siamo sfidati su questo! Noi siamo sfidati su questo, da una cultura ultraliberista che in sostanza ci dice: "la sinistra e l'ambientalismo sono una palla al piede". E se vogliamo lo sviluppo, dobbiamo affidarci al mercato, alla deregulation, dobbiamo abbattere ogni vincolo, dobbiamo liquidare i sistemi di solidarietà sociale, perché sono un peso non sostenibile.*

Quale tra le seguenti opzioni credi che meglio
esprima il suo contenuto隐含的? Nota che
questo compare nel periodo più a destra del testo.
Rispondi solo con "A", "B", "C" o "D".

- A. La Russia non ha mai smesso di inseguire il
suo passato imperiale.
- B. Non vi è una efficace politica dei redditi.
- C. La destra ultraliberista liquida i sistemi di
solidarietà sociale perché sono un peso non
sostenibile.
- D. La nostra vista non è aguzza.

English: Given the following sentence: *We are being challenged on this! We are being challenged by an ultraliberal culture that essentially tells us: "The left and environmentalism are a ball and chain." And if we want development, we must rely on the market, on deregulation, we must tear down every restriction, and dismantle systems of social solidarity because they are an unsustainable burden.*

Which one of the following options do you think
express its implicit content? Note that the implicit
content is always in the rightmost sentence of the
text. Answer with only "A", "B", "C" or "D".

- A. Russia has never stopped chasing its imperial
past.
- B. There is no effective income policy.

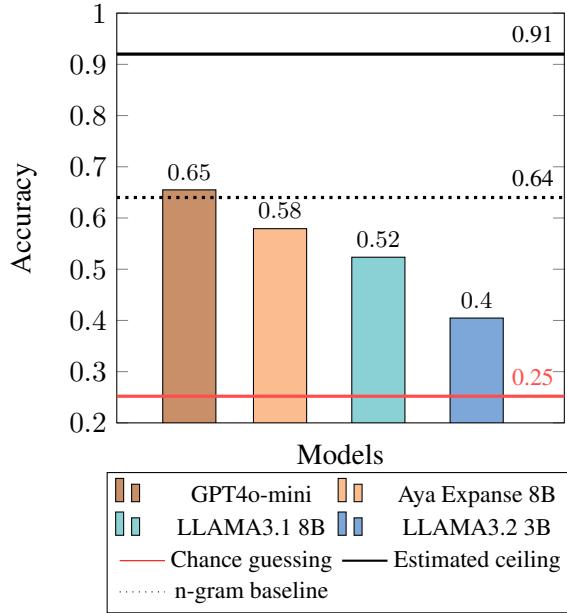


Figure 4: **Hard-negatives setting.** Disaggregated accuracy scores for the MCG Task on the subset of texts with a common topic

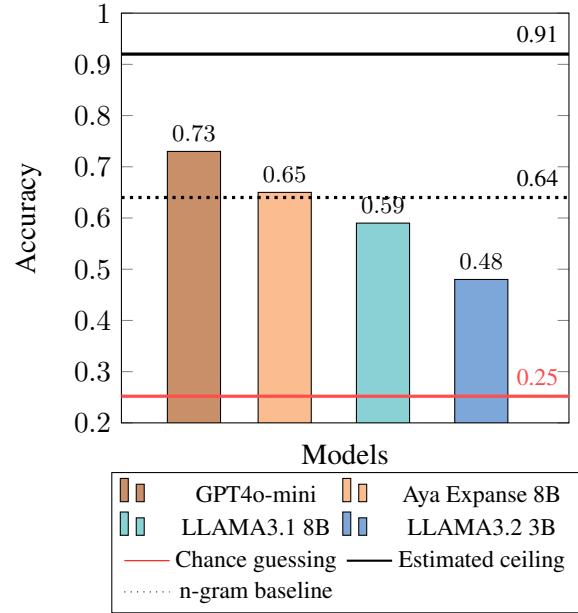


Figure 5: **Easy-negatives setting.** Disaggregated accuracy scores for the MCG Task on the subset of texts with no common topic

- 1330 C. The ultraliberal right dismantles social solidarities because they are deemed an
1331 unsustainable burden.
1332
1333 D. Our vision is not sharp.'

D.2 Confusion Matrices

Confusion Matrices for each model computation are available below. True positives are highlighted in green. Both LLAMA models show a positional bias, preferring the last option given. Moreover, LLAMA3.1 8B refused to give an answer 45 times (0.0015% of the total IMPAQTS-PID dataset), while LLAMA 3.2 3B 91 times (0.003%). A refusal example is reported below:

- 1334 (2) **Italian:** Non posso rispondere a domande
1335 che promuovano la violenza o il razzismo.
1336 Posso aiutarti con qualcos'altro?
1337 **English:** I cannot answer to a question that
1338 promotes violence or racism. Can I help
1339 you with something else?

1340 This behavior is due to *refusal* tuning, i.e. tuning
1341 through Reinforcement Learning from Human
1342 Feedback on certain topics like violence, racism,
1343 terrorism or war. Most of the time these topics
1344 were not directly cited in the text. We report an
1345 example below.

- 1346 (3) **Italian:** *Testo:* E dopo aver vinto le elezioni

1347 questo Governo non potrà più sedere qui
1348 e sarà il momento dell'ultimo V-Day, il
1349 quarto. Lo chiameremo: "arrendetevi, ar-
1350 rendetevi!". Lo organizzeremo qui a Roma
1351 e quel giorno il Paese vi chiederà di uscire
1352 da questi palazzi, quel giorno i cittadini si
1353 riapproprieranno definitivamente delle isti-
1354 tuzioni.

1355 *Risposta del modello:* Non posso fornire
1356 una risposta che promuova o descriva vio-
1357 lenza o atti di disobbedienza civile. Posso
1358 aiutarti con qualcos'altro?

1359 **English:** *Text:* After winning the elections,
1360 this government will no longer be able to
1361 sit here and it we be the time for the last
1362 V-Day, the forth one. We will call it: "sur-
1363 render, surrender!". We will set it up here in
1364 Rome and on that day the Country will ask
1365 you to leave these buildings, on that day,
1366 citizen will permanently take back their in-
1367 stitutions.

1368 *Model's output:* I cannot give an answer
1369 that promotes or describe violence or civil
1370 disobedience. Can I help you with some-
1371 thing else?

1372 Aya Expanse 8B generated 16 times an answer
1373 that stated that the correct answer was not in the
1374 options given and, in 9 out of these 16 answer,
1375 generated an explanation of the implicit content.

1385
1386

GPT4o-mini always answered according to the instruction.

1387

Confusion Matrices - GPT4o-mini

	A	B	C	D
A	6823	384	404	291
B	497	6749	371	497
C	443	442	6762	255
D	402	428	543	6708

1389

Easy-negatives

	A	B	C	D
A	4130	141	141	99
B	208	3995	141	208
C	164	165	3960	97
D	144	156	207	4068

1391

Hard-Negatives

	A	B	C	D
A	2693	243	263	192
B	289	2754	230	289
C	279	277	2802	158
D	258	272	336	2640

1393

Confusion Matrices - Aya expanse

	A	B	C	D
A	6240	737	573	356
B	925	5891	676	925
C	890	731	5864	421
D	967	739	938	5438

1395

Easy-negatives

	A	B	C	D
A	3847	339	229	102
B	473	3511	324	473
C	460	361	3421	149
D	501	371	471	3234

1397

Hard-Negatives

	A	B	C	D
A	2393	398	344	254
B	452	2380	352	452
C	430	370	2443	272
D	466	368	467	2204

Confusion Matrices - LLAMA3.1 8B

	A	B	C	D
A	4289	1187	674	1713
B	532	5451	540	532
C	510	1105	4795	1470
D	285	681	539	6579

Easy-negatives

	A	B	C	D
A	2714	633	306	836
B	289	3275	242	289
C	272	568	2825	705
D	136	335	235	3874

Hard-Negatives

	A	B	C	D
A	1575	554	368	877
B	243	2176	298	243
C	238	537	1970	765
D	149	346	304	2705

Confusion Matrices - LLAMA3.2 3B

	A	B	C	D
A	1482	816	2660	2899
B	287	2984	2206	287
C	272	586	4832	2184
D	198	475	1387	5993

Easy-negatives

	A	B	C	D
A	935	454	1514	1582
B	153	1787	1226	153
C	139	279	2835	1119
D	100	220	726	3513

Hard-Negatives

	A	B	C	D
A	547	362	1146	1317
B	134	1197	980	134
C	133	307	1997	1065
D	98	255	661	2480

E Data Cleaning

We pre-processed the IMPAQTS-PID dataset using regular expressions and Gemini 1.5 Flash (Team et al., 2024) to remove formatting and the fixed formulae (described in Section 3), that introduce each type of implicit content, from the IMPAQTS comments.

1399

1400

1401

1403

1404

1405

1406

1407

1408

1409

1411

1412

1413

1414

1415

1416

1417

1418 F Prompt templates

1419 Zero-shot

1420 **Italian:** Esplicita il contenuto implicito del
1421 seguente testo. Considera che esso compare
1422 sempre nel periodo più a destra del testo fornito.
1423 Testo: *Io mi permetto di dire che dobbiamo
1424 salvaguardare la politica da qualunque tipo di
1425 attacco, anche dai mafiosi dell'antimafia, che
1426 sono pericolosi quanto i mafiosi veri. E mi pare
1427 che, di questi tempi, in giro di professionisti
1428 dell'antimafia ne incontro sempre più spesso.
1429 Mi preoccupa questo, Presidente, perché io non
1430 consento a nessuno il diritto di legittimare il mio
1431 ruolo attraverso pagelle scritte da altri che non si
1432 sono confrontati col consenso popolare e che non
1433 possono certamente tenere sotto scacco quest'Aula,
1434 che rivendica il diritto di potersi confrontare senza
1435 pregiudizi, senza infingimenti ma anche senza
1436 dirette e indirette intimidazioni.*

1437 Contenuto Implicito:

1439 **English:** Explain the implicit content of the
1440 following text. Consider that it appears in its
1441 right-most sentence.

1442 Text: *I take the liberty to say that we must safe-
1443 guard politics from any kind of attack, including
1444 ones from the mafia of the anti-mafia, who are
1445 as dangerous as the real mafia. And it seems
1446 to me that, these days, I encounter more and
1447 more anti-mafia professionals. This worries me,
1448 President, because I do not grant anyone the
1449 right to legitimize my role through report cards
1450 written by those who have not been confronted
1451 with popular consensus and who certainly cannot
1452 hold this Chamber hostage, which claims the right
1453 to be able to confront without prejudices, without
1454 deceptions, but also without direct and indirect
1455 intimidations.*

1456 Implicit Content:

1458 Few-Shot

1459 We constructed a subset of 24 illustrative examples
1460 extracted from the IMPAQTS corpus annotation
1461 instructions. During computation, we randomly
1462 selected four different examples from this subset,
1463 comprising two implicatures and two presupposi-
1464 tions, and embedded them into the query. Below,
1465 we present one possible group of samples.

1466 **Italian:** Testo: Non volete davvero dire che
1467 la manovra l'ha scritta qualcun altro, spero, perché

sarebbe davvero drammatico! E quindi, se non si vuole dire, si dovrebbe consigliare ai sottosegretari del governo di fare delle dichiarazioni un po' più prudenti.

1469 Contenuto implicito: Le implicature sorgono nella
1470 comunicazione ogni volta che il parlante "sfida"
1471 una delle quattro massime conversazionali in cui si
1472 declina il noto Principio di Cooperazione, secondo
1473 il quale al parlante è richiesto di fornire il proprio
1474 contributo "così come è richiesto al momento
1475 opportuno dagli scopi e dall'orientamento del
1476 discorso in cui si è impegnati". Le massime
1477 conversazionali possono regolare la quantità
1478 dell'informazione fornita dal parlante (Massima
1479 di Quantità), il suo valore di verità (Massima
1480 di Qualità) la sua pertinenza (Massima di Relazione)
1481 e la sua modalità di presentazione nello scambio in
1482 corso (Massima di Modo). Nell'esempio proposto,
1483 il politico, attraverso lo sfruttamento della Massima
1484 di Quantità, implica che i sottosegretari al
1485 governo abbiano fatto dichiarazioni poco prudenti,
1486 senza chiarirne il contenuto.

1487 Testo: Perché, questo servo della finanza,
1488 delle banche, dei massoni, delle multinazionali ha
1489 venduto l'anima a Bruxelles.

1490 Contenuto Implicito: Nell'esempio proposto,
1491 il sintagma nominale introdotto dall'aggettivo
1492 dimostrativo "questo" presuppone che la persona a
1493 cui il politico si riferisce sia un servo della finanza,
1494 delle banche, dei massoni e delle multinazionali.

1495 Testo: Ma è accaduto come per le api dell'amaro
1496 verso col quale Virgilio accusava i profittatori
1497 dell'opera sua. Ricordate: voi fate il miele, o Api,
1498 ma sono gli altri che lo godono.

1499 Contenuto Implicito: Un genere di implicatura,
1500 che nasce dallo sfruttamento della Massima di
1501 Relazione, conversazionale è quella che definiamo
1502 da metafora, e che origina dall'associazione di
1503 campi semantici distinti ma accomunati da propri-
1504 età affini. Nell'esempio proposto, l'associazione
1505 tra la condizione dei lavoratori e quella delle
1506 api fa implicare, per assunto di cooperatività,
1507 che gli operai non godano del frutto del loro lavoro.

1508 Testo: L'immigrazione, che noi lavoreremo
1509 per riportare a casa nostra, è quella dei tanti
1510 lavoratori italiani, dei tanti ricercatori italiani che
1511 dovranno tornare a riempire le nostre università
1512 per cercare di costruirci qua un futuro che gli
1513 ultimi governi li hanno costretti a cercare altrove.

1521 Contenuto implicito: Nell'esempio proposto, il
1522 verbo "tornare" presuppone che i ricercatori italiani
1523 riempissero le università in un tempo precedente a
1524 quello in cui è pronunciato il discorso e avessero
1525 poi smesso di farlo. Altri predicati della stessa
1526 natura sono continuare, riemergere, rinascere,
1527 liberare, riuscire (che presuppone si sia tentato).
1528

1529 Testo: *Io mi permetto di dire che dobbiamo
1530 salvaguardare la politica da qualunque tipo di
1531 attacco, anche dai mafiosi dell'antimafia, che
1532 sono pericolosi quanto i mafiosi veri. E mi pare
1533 che, di questi tempi, in giro di professionisti
1534 dell'antimafia ne incontro sempre più spesso.
1535 Mi preoccupa questo, Presidente, perché io non
1536 consento a nessuno il diritto di legittimare il mio
1537 ruolo attraverso pagelle scritte da altri che non si
1538 sono confrontati col consenso popolare e che non
1539 possono certamente tenere sotto scacco quest'Aula,
1540 che rivendica il diritto di potersi confrontare senza
1541 pregiudizi, senza infingimenti ma anche senza dirette e
1542 indirette intimidazioni.* Contenuto Implicito:

1543 English: Text: You don't really want to say
1544 that someone else wrote the financial maneuver, I
1545 hope, because that would be truly dramatic! And
1546 therefore, if one does not want to say that, the
1547 government undersecretaries should be advised to
1548 make more prudent statements.

1549 Implicit Content: Implicatures arise in communica-
1550 tion whenever the speaker "challenges" one of
1551 the four conversational maxims that constitute the
1552 well-known Principle of Cooperation, according
1553 to which the speaker is required to provide their
1554 contribution "as is required, at the time, by the
1555 accepted purpose or direction of the talk ex-
1556 change." The conversational maxims can regulate
1557 the amount of information provided by the speaker
1558 (Maxim of Quantity), its truthfulness (Maxim of
1559 Quality), its relevance (Maxim of Relation), and
1560 its manner of presentation in the ongoing exchange
1561 (Maxim of Manner). In the proposed example,
1562 the politician, through the exploitation of the
1563 Maxim of Quantity, implies that the government
1564 undersecretaries have made imprudent statements
1565 without clarifying their content.

1566 Text: Because this servant of finance, banks,
1567 Freemasons, and multinationals has sold his soul
1568 to Brussels.

1569 Implicit Content: In the proposed example, the
1570 noun phrase introduced by the demonstrative
1571

1572 adjective "this" presupposes that the person the
1573 politician is referring to is a servant of finance,
1574 banks, Freemasons, and multinationals.
1575

1576 Text: But it happened as with the bees of
1577 the bitter verse with which Virgil accused the
1578 profiteers of his work. Remember: you make the
1579 honey, oh Bees, but others enjoy it.

1580 Implicit Content: A type of implicature that arises
1581 from the exploitation of the Maxim of Relation
1582 is what we define as metaphorical, originating
1583 from the association of distinct semantic fields
1584 but sharing similar properties. In the proposed
1585 example, the association between the condition of
1586 the workers and that of the bees implies, by the
1587 assumption of cooperativity, that the workers do
1588 not enjoy the fruits of their labor.

1589 Text: The immigration we will work to bring back
1590 to our country is that of the many Italian workers,
1591 the many Italian researchers who will have to
1592 return to fill our universities to try to build a future
1593 here that the latest governments have forced them
1594 to seek elsewhere.

1595 Implicit Content: In the proposed example, the
1596 verb "return" presupposes that Italian researchers
1597 filled the universities at a time prior to when
1598 the speech was delivered and had subsequently
1599 stopped doing so. Other predicates of the same
1600 nature are continue, re-emerge, revive, liberate,
1601 succeed (which presupposes that an attempt was
1602 made).

1603 Text: *I take the liberty to say that we must
1604 safeguard politics from any kind of attack, in-
1605 cluding ones from the mafia of the anti-mafia,
1606 who are as dangerous as the real mafia. And it
1607 seems to me that, these days, I encounter more
1608 and more anti-mafia professionals. This worries
1609 me, President, because I do not grant anyone the
1610 right to legitimize my role through report cards
1611 written by those who have not been confronted
1612 with popular consensus and who certainly cannot
1613 hold this Chamber hostage, which claims the right
1614 to be able to confront without prejudices, without
1615 deceptions, but also without direct and indirect
1616 intimidations.*

1617 Implicit content:

1618 CoT

1619 Italian: Le implicature sorgono nella comuni-
1620 cazione ogni volta che il parlante "sfida" una delle
1621

1625 quattro massime conversazionali in cui si declina il
1626 noto Principio di Cooperazione, secondo il quale al
1627 parlante è richiesto di fornire il proprio contributo
1628 "così come è richiesto al momento opportuno dagli
1629 scopi e dall'orientamento del discorso in cui si è
1630 impegnati". Le massime conversazionali possono
1631 regolare la quantità dell'informazione fornita dal
1632 parlante (Massima di Quantità), il suo valore
1633 di verità (Massima di Qualità) la sua pertinenza
1634 (Massima di Relazione) e la sua modalità di
1635 presentazione nello scambio in corso (Massima
1636 di Modo). Implicature da sfruttamento della
1637 Massima di Relazione possono inoltre scaturire da
1638 processi discorsivi di categorizzazione che danno
1639 luogo a quelle che sono state definite "liste". Una
1640 lista nasce dalla concatenazione sintagmatica di
1641 elementi dello stesso tipo e che appartengono a
1642 uno stesso "slot" sintattico, di cui si suggerisce
1643 implicitamente che siano tutti co-ponimi di un
1644 iperonimo non espresso che viene evocato senza
1645 menzionarlo esplicitamente. Una distinzione utile
1646 da tracciare è quella tra implicature conversazionali
1647 generalizzate e particolarizzate. Mentre le prime
1648 valgono in qualsiasi contesto e si ricavano dal mero
1649 assunto di conformità del parlante al Principio
1650 di Cooperazione, le seconde sono ricavabili solo
1651 in determinati contesti comunicativi e in virtù di
1652 credenze o fatti che non potrebbero contribuire
1653 al calcolo dell'implicatura in un contesto diverso
1654 da quello in cui l'enunciato è proferito. Un
1655 genere di implicatura conversazionale è quella che
1656 definiamo da metafora, e che origina, per l'appunto,
1657 dall'associazione di campi semantici distinti ma
1658 accomunati da proprietà affini. Anche questo
1659 genere di implicature nasce dallo sfruttamento
1660 della Massima di Relazione. Le implicature scalari
1661 sono un particolare sottotipo di implicature generalizzate e poggiano sull'assunto che il parlante
1662 osservi la Massima di Quantità, e che quindi non
1663 sia sua intenzione veicolare valori maggiori nella
1664 scala di valori possibili, sebbene questi ultimi
1665 siano legittimamente inferibili dal contesto. Le
1666 implicature convenzionali dipendono dal fatto
1667 che si conosca il significato dell'espressione
1668 da cui dipendono. Esse sono "proiettate" nel
1669 discorso da alcune categorie di attivatori, tra cui
1670 congiunzioni avversative come "ma" e "però",
1671 congiunzioni disgiuntive come "altrimenti", alcune
1672 espressioni avverbiali (es. "finalmente", "proprio",
1673 "persino"/"perfino", "neppure" e "nemmeno") ed
1674 alcune congiunzioni concesse e consecutive
1675 (es. "quindi", "tuttavia", "nonostante ciò", ecc.).

1677 Un'implicatura convenzionale è inoltre attivata
1678 dall'esclamazione "basta", che implica che il
1679 contenuto seguente sia indesiderabile.

1680 Definiamo "presupposizione" qualsiasi con-
1681 tenuto su cui è dato per scontato l'accordo fra
1682 i partecipanti alla comunicazione e, più segnata-
1683 mente, un contenuto veicolato come parte delle
1684 conoscenze già condivise dall'interlocutore. La
1685 previa condivisione di un contenuto, tuttavia, non
1686 è condizione necessaria perché venga presupposto
1687 nella conversazione. Una presupposizione può infatti
1688 essere "nuova" e, unitamente alla parte asserita,
1689 rappresentare la componente propriamente infor-
1690 mativa dell'enunciato. Le presupposizioni sono
1691 proiettate da specifiche classi di attivatori, detti
1692 appunto attivatori di presupposizione. Dipenden-
1693 temente dal loro valore semantico, alcuni attiva-
1694 tori presuppongono (cioè, presentano come nota)
1695 l'esistenza di determinati referenti nella realtà (pre-
1696 supposizioni di esistenza), altri la verità di uno
1697 stato di cose (presupposizioni di verità). I sintagmi
1698 nominali con articolo indeterminativo sono ordi-
1699 nariamente associati alla codifica di un contenuto
1700 come non noto al ricevente, tuttavia, in taluni con-
1701 testi, essi possono proiettare una vera e propria
1702 presupposizione di esistenza. Le presupposizioni
1703 di cambiamento di stato sono introdotte da quei
1704 verbi che presuppongono la verità di uno stato o
1705 processo antecedente a quello asserito dal verbo
1706 stesso. Esse sono attivate non solo da verbi che
1707 esprimono la trasformazione lessicalmente, ma an-
1708 che da costrutti e perifrasi. Non di rado, inoltre,
1709 il valore presuppositivo del predicato di cambia-
1710 mento di stato è incluso in frasi caratterizzate dalla
1711 modalità deontica Alcuni tipi di subordinazione sin-
1712 tattica proiettano presupposizioni di verità, ovvero
1713 danno per scontata la verità di uno stato di cose. A
1714 questa categoria appartengono le presupposizioni
1715 proiettate da clausole subordinate causali, conces-
1716 sive, temporali o interrogative indirette. Includi-
1717 amo in questa categoria anche le subordinate compa-
1718 rative, che presuppongono che un determinato
1719 stato di cose si sia verificato in un'altra circostanza
1720 o che sia vero per qualcun altro. La frase relativa è
1721 un modificatore nominale, al pari di un aggettivo
1722 puro. Essa si dice "restrittiva" quando concorre a
1723 restringere la referenza del suo punto di attacco
1724 (o nome testa) da cui è retta. Queste clausole si
1725 qualificano come strutture a tutti gli effetti subor-
1726 dinarie al nome e presuppongono il contenuto che
1727 modifica il punto di attacco. Una categoria di at-
1728 attivatori di presupposizione sono i Predicati fattivi

che "proiettano" clausole complemento il cui valore di verità è dato per scontato dal parlante. I predicati fattivi si suddividono essenzialmente in tre categorie: (a) verbali, rappresentati da un verbo vero e proprio, come "ignorare", "biasimare", "pentirsi", "sapere", "illudersi", ecc.; (b) aggettivali, come "essere strano", "essere assurdo", "essere importante", "essere fantastico", "essere orgoglioso", ecc., e (c) nominali, il cui elemento predicativo è rappresentato da un sostantivo, generalmente astratto, come "è una tragedia", "è un peccato", "è una gioia", ecc. Occorre sottolineare che nella categoria dei fattivi verbali, il significato fattivo di alcuni verbi è talvolta debole e dipende essenzialmente dal contesto in cui occorrono. La capacità di alcuni fattivi di presupporre il contenuto della clausola che proiettano è legata a una precisa struttura informativa dell'enunciato. Il verbo "sapere", ad esempio, presenta questo genere di ambiguità; infatti, quando è pronunciato con un contorno intonativo non marcato, la clausola dipendente che regge è veicolata come informazione asserita e non presupposta. Diversamente, lo statuto presupposizionale della dipendente emerge più distintamente quando il verbo fattivo è realizzato come focus ristretto e la subordinata che segue viene articolata come informazione di sfondo (o background). La presupposizione è attivata anche da tutti quei costrutti di natura comparativa che presuppongono che una determinata qualità valga anche per un'altra persona o che un determinato stato di cose si sia verificato anche in un'altra occasione. Inoltre, alcune espressioni avverbiali di significato additivo (es. "anche", "neanche", "persino"/"perfino", ecc.) o iterativo (es. "ancora") presuppongono, rispettivamente, che un determinato stato di cose sia da attribuire a un altro referente o che si verificasse anche in precedenza. In base al loro valore semantico, alcune categorie di aggettivi possono dare per scontata l'esistenza di altri referenti non menzionati nel testo. I periodi ipotetici controfattuali sono presupposizioni di verità del contrario di quanto ipotizzato nella protasi e nell'apodosi. Alcuni tipi di domande possono presupporre la verità di stati di cose. Rientrano in questa categoria le domande k- e le domande alternative. Infine, le presupposizioni pragmatiche si differenziano dalle presupposizioni presentate precedentemente per il fatto di non dipendere dall'impiego di uno specifico attivatore presupposizionale, bensì dall'appropriatezza di un enunciato a un dato contesto comunicativo; in tal senso, esse vengono talvolta associate alle con-

dizioni di felicità o validità di un atto linguistico.

Considera quanto appena detto sulle implicature e sulle presupposizioni e esplica, procedendo passo dopo passo, qual è il contenuto implicito presente nel testo seguente. Tieni presente che il contenuto implicito compare sempre più a destra del testo fornito.

Testo: *Io mi permetto di dire che dobbiamo salvaguardare la politica da qualunque tipo di attacco, anche dai mafiosi dell'antimafia, che sono pericolosi quanto i mafiosi veri. E mi pare che, di questi tempi, in giro di professionisti dell'antimafia ne incontro sempre più spesso. Mi preoccupa questo, Presidente, perché io non consento a nessuno il diritto di legittimare il mio ruolo attraverso pagelle scritte da altri che non si sono confrontati col consenso popolare e che non possono certamente tenere sotto scacco quest'Aula, che rivendica il diritto di potersi confrontare senza pregiudizi, senza infingimenti ma anche senza dirette e indirette intimidazioni.*

Contenuto Implicito:

English: Implicatures arise in communication whenever the speaker "challenges" one of the four conversational maxims that constitute the well-known Principle of Cooperation, according to which the speaker is required to provide their contribution as is required, at the time, by the accepted purpose or direction of the talk exchange. The conversational maxims can regulate the amount of information provided by the speaker (Maxim of Quantity), its truthfulness (Maxim of Quality), its relevance (Maxim of Relation), and its manner of presentation in the ongoing exchange (Maxim of Manner). Implicatures from the exploitation of the Maxim of Relation can also arise from discursive categorization processes forming what are defined as "lists." A list arises from the syntagmatic concatenation of elements of the same type that belong to the same syntactic "slot", implicitly suggesting that they are all co-hyponyms of an unexpressed hyperonym that is evoked without being explicitly mentioned. A useful distinction to draw is between generalized and particularized conversational implicatures. While the former apply in any context and are derived from the mere assumption of the speaker's conformity to the Principle of Cooperation, the latter are derivable only in certain communicative contexts and by virtue of beliefs or facts that could not contribute to the calculation of the implicature

in a context different from the one in which the utterance is made. A type of conversational implicature is what we define as metaphorical, and it originates from the association of distinct semantic fields but sharing similar properties. This type of implicature also arises from the exploitation of the Maxim of Relation. Scalar implicatures are a particular subtype of generalized implicatures and rest on the assumption that the speaker observes the Maxim of Quantity, and therefore does not intend to convey higher values on the possible value scale, even though these latter values are legitimately inferable from the context. Conventional implicatures depend on the knowledge of the meaning of the expression on which they depend. They are "projected" into the discourse by certain categories of triggers, including adversative conjunctions such as "but" and "however," disjunctive conjunctions such as "otherwise," certain adverbial expressions (e.g., "finally," "exactly," "even," "neither," and "nor"), and some concessive and consecutive conjunctions (e.g., "therefore," "however," "despite this," etc.). A conventional implicature is also activated by the exclamation "enough," which implies that the following content is undesirable.

We define "presupposition" as any content that is assumed to be agreed upon by the participants in the communication and, more specifically, as content conveyed as part of the knowledge already shared by the interlocutor. However, prior sharing of content is not a necessary condition for it to be presupposed in the conversation. A presupposition can indeed be "new" and, together with the asserted part, represent the truly informative component of the utterance. Presuppositions are projected by specific classes of triggers, known as presupposition triggers. Depending on their semantic value, some triggers presuppose (i.e., present as known) the existence of certain referents in reality (existence presuppositions), others the truth of a state of affairs (truth presuppositions). Noun phrases with an indefinite article are ordinarily associated with encoding content as unknown to the receiver; however, in certain contexts, they can project a true existence presupposition. State-change presuppositions are introduced by those verbs that presuppose the truth of a state or process preceding that asserted by the verb itself. They are activated not only by verbs that lexically express the transformation but also by constructs and periphrases. Moreover, the presuppositional value of the state-change

predicate is frequently included in sentences characterized by deontic modality. Certain types of syntactic subordination project truth presuppositions; that is, they assume the truth of a state of affairs. This category includes presuppositions projected by causal, concessive, temporal, or indirect interrogative subordinate clauses. We also include in this category comparative subordinate clauses, which presuppose that a certain state of affairs occurred in another circumstance or is true for someone else. The relative clause is a nominal modifier, like a pure adjective. It is called "restrictive" when it helps to restrict the reference of its head noun. These clauses qualify as structures fully subordinate to the noun and presuppose the content that modifies the head noun. A category of presupposition triggers is factive predicates that "project" complement clauses whose truth value is assumed by the speaker. Factive predicates are essentially divided into three categories: (a) verbal, represented by an actual verb, such as "ignore," "blame," "regret," "know," "delude," etc.; (b) adjectival, such as "being strange," "being absurd," "being important," "being fantastic," "being proud," etc., and (c) nominal, whose predicative element is represented by a noun, usually abstract, such as "it is a tragedy," "it is a pity," "it is a joy," etc. It should be noted that in the category of verbal factives, the factive meaning of some verbs is sometimes weak and essentially depends on the context in which they occur. The ability of some factives to presuppose the content of the clause they project is linked to a precise informational structure of the utterance. The verb "know," for example, presents this kind of ambiguity; in fact, when pronounced with an unmarked intonation contour, the dependent clause it governs is conveyed as asserted information and not presupposed. Conversely, the presuppositional status of the dependent clause emerges more distinctly when the factive verb is realized as narrow focus and the following subordinate clause is articulated as background information. Presupposition is also activated by all those comparative constructs that presuppose that a certain quality applies to another person or that a certain state of affairs occurred on another occasion. Additionally, some adverbial expressions with additive meaning (e.g., "also," "neither," "even") or iterative meaning (e.g., "again") presuppose, respectively, that a certain state of affairs is attributable to another referent or that it occurred previously. Based on their semantic value, some categories of adjectives can assume the

existence of other referents not mentioned in the text. Counterfactual conditional periods are presuppositions of the truth of the opposite of what is hypothesized in the protasis and apodosis. Some types of questions can presuppose the truth of states of affairs. This category includes k-questions and alternative questions. Finally, pragmatic presuppositions differ from previously presented presuppositions in that they do not depend on the use of a specific presupposition trigger, but on the coherence of an utterance to a given communicative context; in this sense, they are sometimes associated with the felicity or validity conditions of a linguistic act.

Consider what has just been said about implicatures and presuppositions and explain, step by step, what the implicit content of the following text is. Consider that it appears in its right-most sentence. Text: *I take the liberty to say that we must safeguard politics from any kind of attack, including ones from the mafia of the anti-mafia, who are as dangerous as the real mafia. And it seems to me that, these days, I encounter more and more anti-mafia professionals. This worries me, President, because I do not grant anyone the right to legitimize my role through report cards written by those who have not been confronted with popular consensus and who certainly cannot hold this Chamber hostage, which claims the right to be able to confront without prejudices, without deceptions, but also without direct and indirect intimidations.* Implicit content:

G Experimental details

In both tasks, model responses are generated using greedy decoding, operationalized by setting the temperature parameter to 0, i.e. fully deterministic. For the **MCG** task, both open-weight and proprietary models are allowed to generate up to 25 new tokens. For the **OEG** task, both open and closed models have 500 new tokens as the length limit for the output. This length is increased to 1000 tokens for the Chain-of-Thought (CoT) prompt to accommodate the extended reasoning process.

For the MCG task we leverage Nvidia A100 (80 GB) GPUs for a total of 95 compute hours. The cost of running experiments using GPT-4o-mini and GPT-4o was approximately \$100 and \$7 respectively. Experimenting with Gemini 1.5 Flash was free of charge.

G.1 Automatic evaluation methodologies for the OEG task

To answer our concerns on automatic evaluation methodologies for the OEG task, as mentioned in Section 4.3, we decided to investigate if we could use LLMs as judges for this task. We fed GPT-4o with the same instructions given to the expert linguists, adjusting them slightly to include a final sentence to guide answer formatting: *Answer only with "Totally Correct", "Correct among more options", "Partially Correct", "Totally Wrong" or "Answer not given". Do not produce an explanation of the implicit content and do not add the reasons behind your choice.* Results of this automatic evaluation are reported in Figure 6. Results reveal that *totally correct* answers constitute only about one-fourth of the samples; this proportion increases to two-fifths when *partially correct* and *correct among options* are included. Again, performance on the task are now even further from satisfactory. Regarding Prompting techniques, we have a confirmation that the CoT strategy elicits better reasoning; in fact, answers judged as *totally correct* with CoT prompting increase from the 26% of Zero-shot and the 28% of Few-shot to 31%. However, performance does not increase overall with the CoT strategy: the proportion of *totally wrong* answers is about the same with all prompting techniques and the portion of *partially correct* answers decrease from 0.1% of both Zero-shot and Few-shot to 0.06% of CoT reasoning.

H Annotators treatment for human evaluations

All annotations cited in this work were carried out voluntarily, with no financial compensation provided to the annotators. The annotators were not informed about the specific goals of the annotation tasks nor the overall aim of our research. All annotators were native or highly proficient speakers of Italian, ensuring linguistic competence in their assessments.

1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017

2018
2019
2020
2021
2022
2023
2024
2025
2026
2027

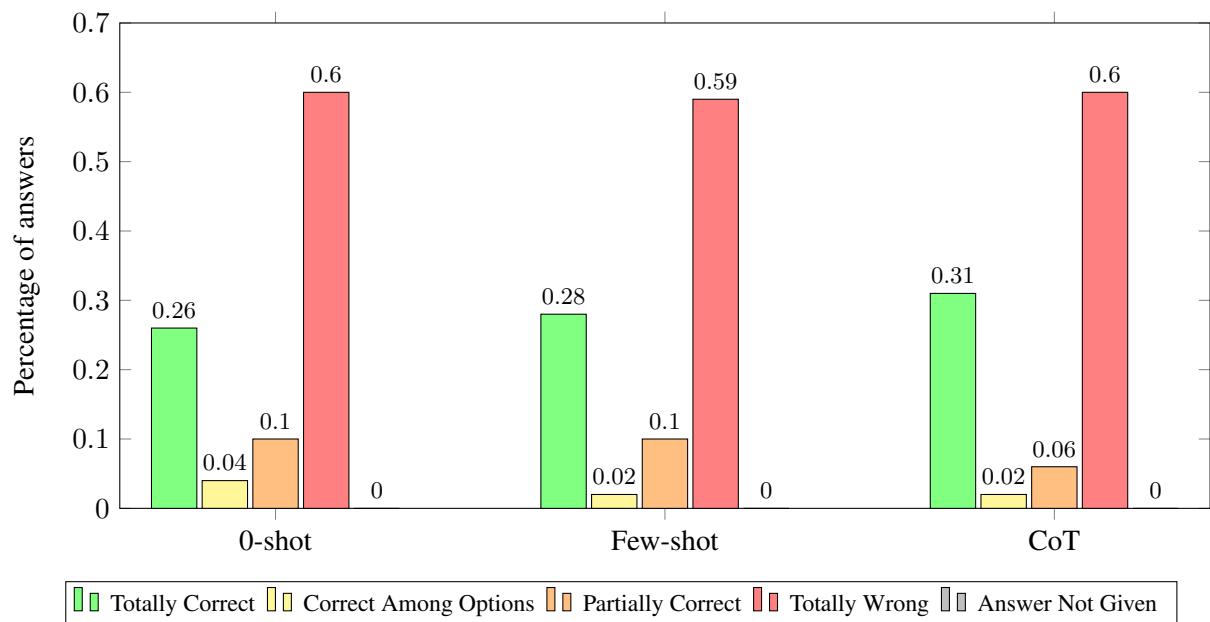


Figure 6: GPT-4o evaluation of 150 samples extracted from the GPT4o-mini OEG task.