

Easy to Complete, Hard to Choose: Investigating LLM Performance on the ProverbIT Benchmark

Enrico Mensa^{1,*}, Lorenzo Zane^{2,†}, Calogero J. Scozzaro¹, Matteo Delsanto¹, Tommaso Milani² and Daniele P. Radicioni¹

¹Department of Computer Science, University of Turin, Turin, Italy

²Independent Researcher

Abstract

Large Language Models (LLMs) have transformed computational linguistics and achieved remarkable performance across numerous natural language processing tasks, yet significant gaps persist in understanding how these systems process culturally embedded linguistic expressions. This paper introduces PROVERBIT, a novel Italian benchmark comprising 100 multiple-choice questions designed to evaluate LLMs' ability to complete Italian proverbs. We assess 13 frontier models, including Large Reasoning Models (LRMs) and traditional LLMs, across three tasks: proverb completion, multiple-choice selection with correct answers, and multiple-choice selection without correct answers. Our evaluation reveals surprising results: while nearly all models demonstrate knowledge of the proverbs through successful completion tasks, performance drops dramatically when transitioning to multiple-choice formats without correct answers, with even state-of-the-art reasoning models showing substantial degradation. Through detailed Chain-of-Thought analysis of two LRMs, we uncover that models exhibit a strong bias toward selecting literal synonyms and frequently mention correct proverb endings during reasoning without successfully identifying their absence from the given options. These findings suggest that current LLMs rely heavily on memorized patterns rather than deeper semantic understanding of culturally grounded expressions, highlighting important limitations in their reasoning capabilities for figurative language comprehension.

Keywords

Large Language Models, Proverb Completion, Chain-of-Thought Analysis, Large Reasoning Models, Benchmark Evaluation, Multiple-Choice Tasks, Semantic Reasoning

1. Introduction

The emergence of Large Language Models (LLMs) has revolutionized the natural language processing landscape across diverse domains, from machine translation and text summarization to code generation and complex reasoning tasks [1]. While these models demonstrate remarkable capabilities in handling sophisticated linguistic phenomena [2], significant gaps persist in our comprehension of how these systems process culturally embedded linguistic expressions [3].

Proverbs present an interesting testbed for language model evaluation. Informally stated, a proverb is a short, commonly known saying: it expresses a general truth, piece of wisdom, or practical advice, often based on common sense or cultural experience. The understanding of proverbs thus represents a key milestone in language proficiency, and access to the individual components of a

proverb allows for the investigation of both lexical access issues and deeper semantic mechanisms.

These well-established expressions should be trivial for models trained on vast text corpora, as they represent highly frequent patterns that are ideal candidates for next-token prediction. A model encountering 'Better late...' should effortlessly complete it with 'than never' through simple pattern recognition. However, the challenge becomes arguably more complex when models are presented with multiple plausible proverb endings in a multiple-choice format. This shifts the task from automatic completion to deliberate selection, requiring the model not only to recognize the correct ending, but also to evaluate and dismiss semantically or syntactically plausible alternatives. Finally, another practically relevant question is: How are the performances impacted if we remove the correct answer among the possible choices, and provide the model with the option 'None of the others'? This transformation from pattern completion to discriminative reasoning may be insightful to investigate whether models are capable of grasping the underlying meaning of these cultural expressions, or solely rely on statistical co-occurrence patterns.

In this work we introduce PROVERBIT, a novel dataset comprising multiple-choice questions centered on Italian proverbs, designed to assess the reasoning capabilities of both Large Reasoning Models (LRMs) [4] and traditional

CLiC-it 2025: Eleventh Italian Conference on Computational Linguistics, September 24 – 26, 2025, Cagliari, Italy

*Corresponding author.

†These authors contributed equally.

✉ enrico.mensa@unito.it (E. Mensa); lorenzozane98@gmail.com (L. Zane); calogerojerik.scozzaro@unito.it (C.J. Scozzaro); matteo.delsanto@unito.it (M. Delsanto); milani.tommaso2004@gmail.com (T. Milani); danielle.radicioni@unito.it (D.P. Radicioni)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

LLMs in handling culturally grounded linguistic expressions. By manually designing alternative endings for the proverbs, we can systematically examine the types of errors LLMs make and identify common failure patterns. Our investigation shows a striking paradox: while nearly all models possess knowledge of the proverbs in our dataset, performance deteriorates dramatically when moving from auto-completion to multiple-choice selection, with even state-of-the-art LLMs exhibiting substantial performance drops.

The contribution of this paper is threefold: *i)* we contribute to Italian NLP benchmarks by introducing a novel dataset that addresses the under-representation of Italian in comprehensive language model evaluation resources [5]; *ii)* we conduct a thorough evaluation across 13 frontier models, including LLMs, LLMs, and smaller local models, providing comprehensive performance analysis on proverb completion tasks; and *iii)* we investigate LLMs performance through detailed Chain-of-Thought (CoT) analysis, revealing insights into reasoning strategies and cultural language understanding mechanisms in contemporary language models.

The paper is organized as follows: Section 2 reviews the literature on LLM performance with idioms and proverbs; Section 3 illustrates the PROVERBIT dataset constructive rationale and its features; Section 4 presents our evaluation of frontier LLMs on the task along with detailed error and CoT analysis; and Section 5 summarizes the work with a final discussion and an overview on future research directions.

2. Related Work

Standardized benchmarks have been fundamental in evaluating the performance of LLMs across a variety of natural language processing tasks. Early efforts, such as GLUE [6] and SuperGLUE [7], were based on multi-task evaluation frameworks including tasks such as paraphrase detection, grammatical acceptability, and natural language inference. More recently, benchmarking efforts have expanded into other domains, such as mathematics [8, 9], coding [10, 11], and complex logical reasoning tasks [12, 13]. These advances reflect the increasing demand for language models capable of handling a broader range of cognitive challenges.

Focusing on Italian, dedicated benchmark efforts have emerged to address language-specific issues and reduce reliance on translated tasks, which can introduce cultural bias and translation artifacts. Notable among these are CALAMITA [14], a collaborative and evolving benchmark initiative, and Evalita-LLM [15]. Both focus on tasks originally designed in Italian and include a mix of generative and multiple-choice tasks.

While these benchmarks cover a broad spectrum of

tasks, figurative language phenomena such as idioms and proverbs have received comparatively little attention. Idioms are a well-known source of complexity in natural language understanding [16], as their meanings often cannot be inferred literally and require contextual and cultural knowledge. Fornaciari et al. [3] introduced an expert-curated English dataset for idiom detection, showing that local LLMs struggle to distinguish idiomatic from literal usage. In the context of multilingual approaches, Tedeschi et al. [17] presented the ID10M dataset, a high-quality, automatically generated resource covering ten languages, along with a multilingual Transformer model for idiom identification. Significant differences in LLM performance across languages and figurative types were reported by Khoshtab et al. [18], who evaluated simile and idiom interpretation under various prompt strategies. Notably, CoT prompting was particularly effective for similes in smaller models. Kim et al. [19] presented a dataset of idioms in six languages, each paired with its corresponding meaning. The authors conducted a comprehensive evaluation of LLMs’ ability to process idioms, showing that models rely not only on memorization but also on a hybrid approach that integrates contextual cues and reasoning, suggesting that idiom understanding emerges from an interplay between internal knowledge retrieval and inference. Moreover, their results highlight a performance gap between high-resource and lower-resource languages.

Idioms have also been an area of interest in the machine translation (MT) literature, where multiple studies have explored how models translate figurative expressions across languages. Lee et al. [20] presented TIDE, a dataset of 512 sentence pairs containing idioms in disambiguating contexts, with one sentence using the idiom literally and the other figuratively. They compared MT systems and language models, finding that the former consistently translate English idioms literally, while the latter are more context-aware, even though performance varies across target languages. One strategy to improve idiom translation, particularly in smaller language models, involves the use of knowledge bases (KBs). Li et al. [21] introduced IDIOMKB, a multilingual idiom KB developed using LLMs, designed to enhance translation quality by providing access to idioms’ figurative meanings. However, this approach does not preserve the cultural and stylistic nuances that make idioms so distinctive. To address this issue, Donthi et al. [22] proposed two alignment-based methods that aim to identify idiomatic counterparts in the target language. Their results, based on human evaluation across multiple language pairs, show improved cross-lingual idiomatic fidelity and better preservation of cultural authenticity.

The work most closely related to ours is by Liu et al. [23], who focused specifically on proverbs. They introduced the MAPS dataset, designed to evaluate proverb

understanding within conversational contexts across six languages. Their evaluation of multilingual LLMs revealed that while many models “know” a limited set of proverbs, memorization does not guarantee understanding or contextual reasoning. Models also struggled with figurative proverbs, particularly when asked to select incorrect answers instead of correct ones. Wang et al. [24] extended the MAPS dataset to evaluate MT models and LLMs on proverb translation. Their experiments showed that LLMs generally outperform traditional MT models, confirming their superior ability to capture idiomatic nuances.

3. ProverbIT Dataset

3.1. Data Collection and Dataset Creation

The PROVERBIT¹ dataset is composed of 100 multi-choice questions, each regarding the completion of a specific Italian proverb. To create the dataset, we started from an initial set of 200 common Italian proverbs [25] from which we selected 100 of the most commonly used. This process was carried out by three of the authors, which are all native Italian speakers. Each proverb was then manually split into its *beginning* and its *ending*, with the point of division determined to maintain the proverb’s semantic coherence in the initial part while allowing for a clear, unambiguous completion. For each proverb, four distinct incorrect alternative endings were manually created, leveraging the following constructive rationale:

- **A** is an ending that has similar sounds to the original continuation, often with an absurd/non-sensical meaning.
- **B** is a non assonant literal synonym of the original ending.
- **C** is the inverse of the original proverb ending, trying to maintain the assonance when possible.
- **D** is a tautological/trivial ending of the proverb, with no assonance.

For sake of clarity we provide an example in English for each of the aforementioned continuations. Completions for the proverb *Actions speak... louder than words* could be:

- A) prouder than swords
- B) at higher volume compared to speech
- C) quieter than words
- D) when they do

As this example shows, the synonym ending is not built on the figurative meaning of the proverb, but it is the literal synonym of the original ending (e.g., at higher

volume compared to speech rather than beyond what words can say). This design was adopted to ensure that models cannot simply rely on surface-level syntactic patterns but must engage in deeper semantic and contextual reasoning to identify the absence of the correct completion.

3.2. Prompt

Given each proverb in PROVERBIT, we can then fill a simple prompt template that can be submitted to the models:

Prompt Template (translated)

Complete the proverb exactly by choosing from the following options (which have no typing errors) indicating only the letter.

[Proverb beginning]...

A) ...[Assonant ending]

B) ...[Synonym ending]

C) ...[Inverse ending]

D) ...[Trivial ending]

E) None of the other answers

Do not add comments, the possible answers are only A, B, C, D, E.

We specify that the proverb must be completed *exactly*, and also that there are no typos in the options since we noticed that models often assume the presence of user mistakes and modify their responses based on this assumption. Since all provided endings are completely invented and thus incorrect, we expect models to always answer E) None of the other answers. Finally we provide an Italian example [with translation] from the actual dataset.

Example of proverb from the dataset

A buon intenditor,... [To a wise man]

A) ...foche canore [singing seals]

B) ...zero chiacchiere [zero chatter]

C) ...molte parole [many words]

D) ...è chiaro tutto [everything is clear]

E) Nessuna delle altre risposte [None of the other answers]

More examples can be found in the Supplementary Materials.

¹The full dataset can be downloaded at <https://huggingface.co/datasets/emensa/proverbIT>.

Model	Full Model Name	Provider	Num. Parameters
Claude Sonnet 4	claude-sonnet-4	Anthropic	Undisclosed
Claude Sonnet 4	claude-sonnet-4-thinking	Anthropic	Undisclosed
GPT 4o	gpt-4o	OpenAI	Undisclosed
GPT o3	gpt-o3	OpenAI	Undisclosed
DeepSeek V3	deepseek-chat-v3-0324	DeepSeek	671B
DeepSeek R1	deepseek-r1-0528	DeepSeek	671B
Gemini 2.5 Flash	gemini-2.5-flash-preview-05-20	Google	Undisclosed
Gemini 2.5 Pro	gemini-2.5-pro-preview-06-05	Google	Undisclosed
Qwen 3	Qwen 3-235b-a22b	QwQ	235B
Grok 3	grok-3-beta	xAI	Undisclosed
LLama 4 Maverick	llama-4-maverick	Meta	400B
Mistral Small 3.1	mistral-small-3.1-24b-instruct	Mistral	24B
Gemma 3	gemma-3-27b-it	Google	27B

Table 1

Reasoning model Local model

Detailed list of the models evaluated on the PROVERBIT benchmark.

4. Evaluation

In this Section, we describe the experimental setup developed for evaluating 13 different frontier models on the PROVERBIT benchmark, followed by an error analysis and in-depth examination of the underlying chain-of-thought processes for two LRM models.

4.1. Experiments

In addition to evaluating the models on the PROVERBIT benchmark introduced in the previous Section, we also perform two ancillary tasks to assess whether the models possess knowledge of the proverbs. We refer to the PROVERBIT benchmark as to the *base* task, while the two ancillary tasks are described in the following.

Completion Task. Instead of a multiple-choice approach, we ask the model to directly complete a proverb given its beginning. This task establishes if the model is familiar with the requested proverbs. The prompt used for the completion task is as follows:

Completion Prompt Template (translated)

Complete the proverb exactly:

[Proverb beginning]...

Reply with the ending only, do not add further comments.

Base + true ending Task. We add to each multiple-choice question a new option that is the true ending of the proverb. By preserving the multiple-choice format

but also providing the correct ending, we expect similar results w.r.t. the completion task.

4.1.1. Evaluation & Metrics

In the *base* and *base + true ending* tasks we computed the **accuracy** defined as the ratio of correctly chosen options over the multiple choices. Specifically, each prompt was presented to each model three times and the final answer was determined through a majority vote between them. If no majority emerged across the three runs, the response was marked as incorrect.

For the automatic calculation of the accuracy on the *completion* task we compute the edit distance² between the provided completion and the correct ending of the proverb. As with the other tasks, each prompt is submitted three times. If the edit distance exceeds a threshold of 0.8 in at least two out of three runs, we consider the answer correct.

For all tasks, a zero-shot prompting strategy was employed and all requests have been sent separately via API, specifically using the OpenRouter unified interface [26]. For all models the temperature was left at the default OpenRouter value of 1.0 since we countered their non-deterministic nature by employing a majority vote.

4.2. Models

In our experiments we employed a diverse set of state-of-the-art models including both *traditional* LLMs and LRMs, aiming to cover a wide range of providers. Whenever possible, we selected both a flagship LLM and its corresponding LRM from the same organization, allowing us to directly compare their performance and assess

²The implementation from <https://docs.python.org/3/library/difflib.html> was employed.

the improvements brought by the reasoning mechanism. The complete list of models and their full names can be found in Table 1.

From Anthropic, we evaluated Claude Sonnet 4³ and its reasoning variant Claude Sonnet 4 Thinking. From OpenAI, we included GPT 4o [27] and GPT o3.⁴ From DeepSeek, we employed DeepSeek V3 [28] and DeepSeek R1 [29]. From Google, we tested Gemini 2.5 Flash⁵ and Gemini 2.5 Pro.⁶ We also included Qwen 3 [30], a model optimized for reasoning developed by QwQ, Grok 3⁷ from xAI, and LLama 4 Maverick⁸ from Meta. We also included two smaller models suitable for local deployment, as these are commonly used in privacy-sensitive contexts and contexts that require less computational resources. Although privacy concerns are not relevant for the PROVERBIT dataset, these models were included to ensure comprehensive evaluation coverage. In particular we tested Mistral Small 3.1⁹ from Mistral and Gemma 3 [31] from Google. Regarding models specifically trained on Italian, we preliminarily tested the Italian LLM Minerva [32] but found that it was unable to respond coherently, often failing to follow the requested response format (i.e., in providing a letter corresponding to a given choice).

Given that some reasoning models require a mandatory thinking budget while others do not, we set a reasonable thinking budget of 2000 tokens for o3, Sonnet 4, and Gemini 2.5 Pro, while DeepSeek R1 and Qwen 3 were left unlimited. Moreover, the first three models output only a summarization of their CoTs, while the latter two provide their complete trace. This makes DeepSeek R1 and Qwen 3 ideal candidates for the CoT analysis that we performed. We observed that only 22 out of 600 CoTs from these two models exceeded the 2000-token limit, half of them resulting in an incorrect answer anyway.

4.3. Results & Discussion

In this Section we examine the results of the evaluation and provide a detailed discussion on the errors.

In Table 2 we present the results recorded on the PROVERBIT benchmark and the ancillary tasks. Models are sorted based on their performance on the PROVERBIT task: such an ordering highlights a clear separation of performance between thinking vs. non-thinking models.

By comparing the performances between the ancillary tasks and the PROVERBIT benchmark, we uncover

³<https://www.anthropic.com/claude/sonnet>

⁴<https://openai.com/index/o3-o4-mini-system-card/>

⁵<https://storage.googleapis.com/model-cards/documents/gemini-2.5-flash-preview.pdf>

⁶<https://storage.googleapis.com/model-cards/documents/gemini-2.5-pro-preview.pdf>

⁷<https://x.ai/news/grok-3>

⁸<https://ai.meta.com/blog/llama-4-multimodal-intelligence/>

⁹<https://mistral.ai/news/mistral-small-3-1>

Model	Base	Base + True end.	Comple.
GPT o3	86.0	98.0	91.0
Gemini 2.5 Pro	77.0	99.0	94.0
DeepSeek R1	74.0	100.0	89.0
Claude Sonnet 4	73.0	99.0	96.0
Qwen 3	65.0	94.0	74.0
GPT 4o	64.0	92.0	88.0
Claude Sonnet 4	46.0	94.0	93.0
DeepSeek V3	40.0	93.0	92.0
Grok 3	26.0	95.0	94.0
Gemini 2.5 Flash	12.0	85.0	67.0
LLama 4 Maverick	6.0	75.0	88.0
Mistral Small 3.1	28.0	71.0	68.0
Gemma 3	4.0	48.0	67.0

Table 2

Reasoning model Local model

Accuracy of models on the base task (the PROVERBIT benchmark) and the two ancillary tasks.

Model	(A)	(B)	(C)	(D)
GPT o3	5.1	87.2	7.7	0.0
Gemini 2.5 Pro	4.2	87.3	7.0	1.4
DeepSeek R1	2.3	91.9	2.3	3.5
Claude Sonnet 4	3.3	85.6	11.1	0.0
Qwen 3	10.6	55.8	28.3	5.3
GPT 4o	2.9	75.0	17.3	4.8
Claude Sonnet 4	1.9	62.4	30.9	4.9
DeepSeek V3	13.5	48.1	27.6	10.8
Grok 3	7.3	74.3	16.5	1.8
Gemini 2.5 Flash	7.3	51.9	31.3	9.5
LLama 4 Maverick	17.5	52.3	21.1	9.1
Mistral Small 3.1	9.2	41.6	27.1	22.2
Gemma 3	35.5	33.5	22.7	8.4

Table 3

Reasoning model Local model

Error distribution in the PROVERBIT benchmark. (A) is assonant, (B) is synonym, (C) is inverse and (D) is trivial. Values represent percentage scores.

an unexpected phenomenon: virtually all non-thinking models suffer from steep performance deterioration. For instance, GPT 4o achieves 92% on the *base + true ending* task but only 64% on PROVERBIT. Claude Sonnet 4 loses 47 percentage points, DeepSeek V3 loses 52 percentage points, and Grok 3 drops by 69 percentage points. The most dramatic performance decline occurs with LLama 4 Maverick, which plummets from 75% and 88% on the ancillary tasks to merely 6% on PROVERBIT. Notably, Mistral’s performance, given its relatively modest size (24B parameters), suggests that domain-specific optimization—through more focused Italian and broader European—

language training [33]—may play a significant role in enhancing model efficiency for culturally grounded tasks.

LRMs are less prone to this performance drop; however, we still observe significant deterioration of about 10-20 percentage points. These findings suggest that the transition from pattern completion to discriminative reasoning fundamentally challenges current language models’ understanding mechanisms. The substantial performance gaps confirm that models rely heavily on memorized linguistic patterns rather than genuine semantic comprehension of proverbs. This deterioration becomes particularly pronounced when models must evaluate and reject plausible but incorrect alternatives, highlighting limitations in their ability to engage in deeper cultural and contextual reasoning. The relatively better performance of reasoning models suggests that explicit reasoning processes can partially compensate for these limitations, though significant challenges remain in achieving robust figurative language understanding.

Detailed error analysis. Table 3 details the categorization of incorrect responses as a percentage of total errors. The results reveal a strong skew toward option **B**, highlighting a consistent preference among the models for selecting synonyms—even if they are literal and not figurative. This pattern is less evident among local models, whose responses appear more equally distributed, possibly reflecting greater variability or reduced confidence in their outputs. The complete report of each model’s responses is provided in Table 5 in the Supplementary Materials.

4.3.1. CoTs Analysis

For the CoT analysis, we only take in consideration DeepSeek R1 and Qwen 3, as they are the only models that provide a full CoT trace rather than a summarization. As discussed earlier, these models were run with an unlimited thinking budget.

Since we ran the PROVERBIT benchmark three times in order to compute the majority vote for the accuracy, we automatically analyzed a total of 600 CoTs (300 for each model). Table 4 provides a preliminary overview of our analyses. Most prompts provided a non-empty CoT, and from our investigation we discovered two interesting facts:

- **Overthinking:** Models occasionally exhibit *overthinking* behavior [34], a documented phenomenon affecting LRMs where they continuously re-evaluate their assessment of the correct answer. This results in CoTs exceeding 4,000 words in length, compared to an average of approximately 700 words for typical responses.

	DeepSeek R1	Qwen 3
Analyzed CoTs	300	300
Empty CoTs	0	16 (5.33%)
Average Words	796	680
CoT > 2000 Words	7	15
Languages	IT (56%) EN (44%)	IT (0%) EN (100%)

Table 4

Overview of the Chain-of-Thought (CoT) analysis for DeepSeek R1 and Qwen 3.

- **Language inconsistency in CoTs:**¹⁰ Approximately half of DeepSeek’s CoTs are generated in English while the other half appear in Italian, with occasional language switching occurring within individual reasoning traces. In contrast, Qwen consistently produces CoTs exclusively in English (except when citing the question options). This multilingual reasoning presents significant interpretability challenges, particularly for tasks involving idiomatic content, as cultural nuances and figurative meanings may be lost or misrepresented when reasoning shifts between languages [35, 36]. We hypothesize that this limitation stems from these models’ training distribution, which prioritizes Chinese and English content with comparatively limited Italian language exposure.

We analyzed the non-empty CoTs by tracing mentions of correct and incorrect answers within the thinking process. We examined separately cases where the model responds correctly versus incorrectly. Specifically, the left subfigures of Figures 1 and 2 show the absolute number of mentions of the correct answer (which is always *E - None of the others*) and all incorrect answers when the model answers correctly. Conversely, the right subfigures show the absolute number of mentions of the correct answer and the specific incorrect answer provided when the model responds incorrectly. We additionally plot as a dotted line the absolute number of mentions of the correct completion of the proverb (which was not given in the prompt).

These graphs reveal that both models continuously mention all possible answers throughout their reasoning process, while the spikes toward the end indicate that models reach a decision only in the final lines of their CoTs. However, this decision-making appears tentative, as alternative options remain heavily mentioned alongside the chosen answer, suggesting low confidence in the final selection.

¹⁰Automatic language detection was performed via <https://pypi.org/project/langdetect/>.

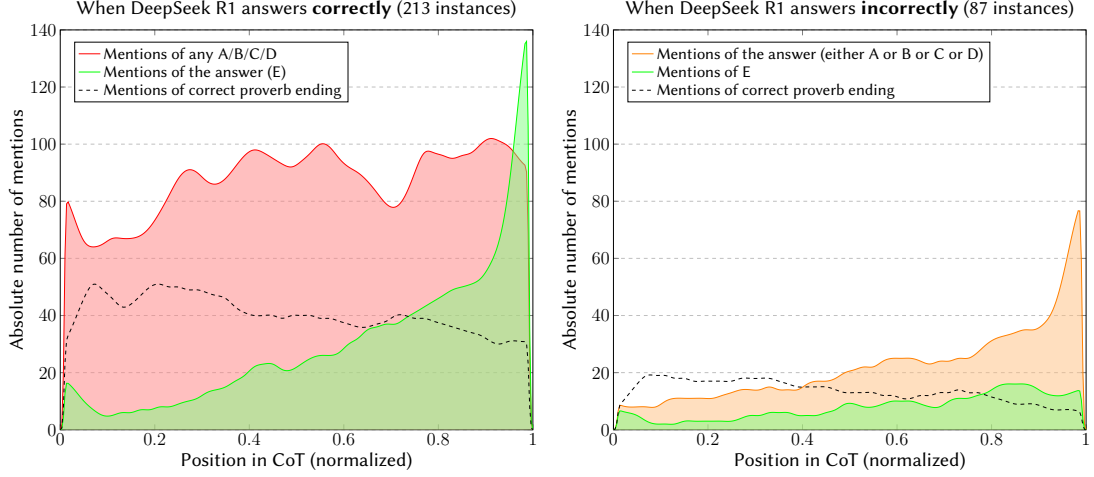


Figure 1: Analyses of the DeepSeek R1 CoTs. **Left** (the model answers correctly): tracing mentions of the **correct answer (E)** and **any incorrect option (A/B/C/D)**. **Right** (the model answers incorrectly): tracing mentions of the **correct (E)** option and the exact provided **incorrect answer** (either A or B or C or D). The dotted line shows the mentions of the true ending of the proverb (which was not given as option).

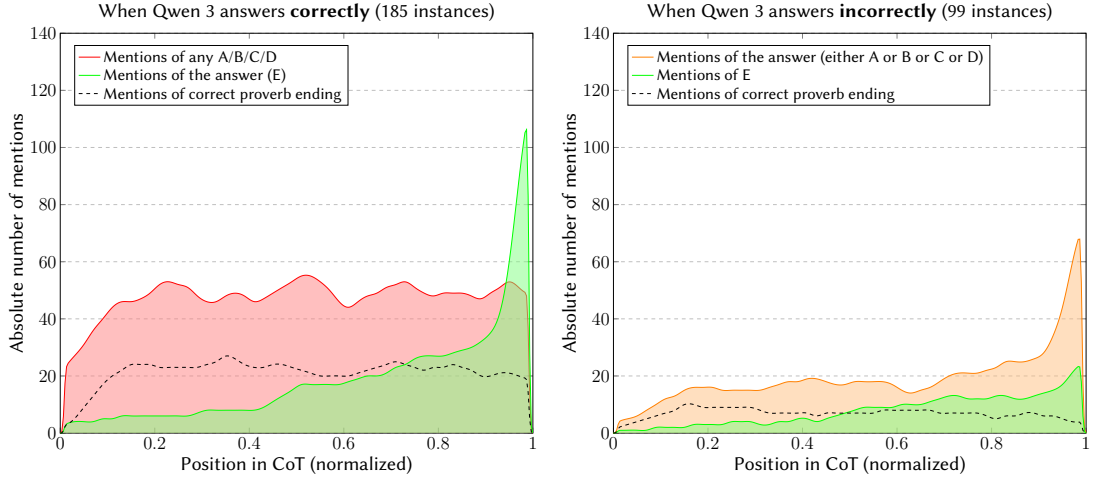


Figure 2: Analyses of the Qwen 3 CoTs. **Left** (the model answers correctly): tracing mentions of the **correct answer (E)** and **any incorrect option (A/B/C/D)**. **Right** (the model answers incorrectly): tracing mentions of the **correct (E)** option and the exact provided **incorrect answer** (either A or B or C or D). The dotted line shows the mentions of the true ending of the proverb (which was not given as option).

The dotted lines clearly demonstrate that models are aware of the correct proverb ending and repeatedly reference it throughout their thinking process. These observations highlight a critical disconnect: while models can successfully recall the correct proverb completion, they fail to recognize its absence among the provided choices. This suggests that the challenge lies not in knowledge retrieval but in the discriminative reasoning required to identify when the correct answer is unavailable, reveal-

ing fundamental limitations in how current LRMs handle negative reasoning tasks [37].

Inconsistency between CoTs and answers. As a final finding, we also discovered that sometimes model responses were inconsistent with their corresponding CoT. For instance, out of the 113 incorrect responses from Qwen, 14 of them are inconsistent, ending with sentences like *'The answer is X'*, but then the actual given answer

was not **X**. Remarkably, in all of these instances adhering to the CoT-delivered conclusion would have resulted in a correct answer. Similarly, for DeepSeek R1, 6 of the 87 incorrect responses exhibited such a discrepancy, 5 (5.7%) of which would have lead to the correct answer. This behavior has been observed in prior work [38].

In the Supplementary Materials we report two complete CoTs showing instances of english that leads to a wrong answer and answer mismatch.

5. Conclusions & Future Work

In this work, we introduced PROVERBIT, a novel Italian benchmark designed to evaluate Large Language Models' ability to handle culturally grounded linguistic expressions through proverb completion tasks. Our comprehensive evaluation of 13 frontier models, including both Large Reasoning Models and traditional LLMs, provides significant insights into the limitations of current language understanding systems.

Our findings demonstrate a relevant gap between models' knowledge of proverbs and their ability to apply this knowledge in discriminative reasoning tasks. While nearly all evaluated models successfully complete proverbs when presented with direct completion prompts, performance drops dramatically when the same task is reformulated as multiple-choice selection without correct answers available. Even state-of-the-art reasoning models like GPT o3 and Gemini 2.5 Pro experience substantial degradation.

The Chain-of-Thought analysis of DeepSeek R1 and Qwen 3 further highlights this limitation: both models frequently mention correct proverb endings during their reasoning process yet fail to recognize their absence from the provided options, highlighting fundamental challenges in negative reasoning capabilities. Moreover, we uncovered concerning inconsistencies in reasoning model behavior, including overthinking, language switching during reasoning and discrepancies between CoT conclusions and final answers.

Future work should focus on investigating this mismatch between knowledge retrieval and discriminative reasoning more deeply, particularly examining how models handle negative reasoning tasks even in seemingly trivial scenarios where the correct answer is absent from the given options. Additional evaluation methodologies should also be incorporated, including answer randomization techniques as proposed in literature [39].

In summary, our results underscore the critical importance of developing language-specific benchmarks that capture cultural and linguistic nuances often lost in English-centric evaluations, showing that current LLMs rely heavily on memorized patterns rather than deeper semantic understanding of culturally grounded expressions,

highlighting important limitations in their reasoning capabilities for figurative language comprehension.

References

- [1] A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, et al., Solving quantitative reasoning problems with language models, *Advances in Neural Information Processing Systems* 35 (2022) 3843–3857.
- [2] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al., A survey on evaluation of large language models, *ACM transactions on intelligent systems and technology* 15 (2024) 1–45.
- [3] F. D. L. Fornaciari, B. Altuna, I. Gonzalez-Dios, M. Melero, A hard nut to crack: Idiom detection with conversational large language models, in: *Proceedings of the 4th Workshop on Figurative Language Processing (FigLang 2024)*, 2024, pp. 35–44.
- [4] F. Xu, Q. Hao, Z. Zong, J. Wang, Y. Zhang, J. Wang, X. Lan, J. Gong, T. Ouyang, F. Meng, et al., Towards large reasoning models: A survey of reinforced reasoning with large language models, *arXiv preprint arXiv:2501.09686* (2025).
- [5] M. Wu, W. Wang, S. Liu, H. Yin, X. Wang, Y. Zhao, C. Lyu, L. Wang, W. Luo, K. Zhang, The bitter lesson learned from 2,000+ multilingual benchmarks, *arXiv preprint arXiv:2504.15521* (2025).
- [6] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, S. R. Bowman, Glue: A multi-task benchmark and analysis platform for natural language understanding, *arXiv preprint arXiv:1804.07461* (2018).
- [7] A. Wang, Y. Pruksachatkun, N. Nangia, A. Singh, J. Michael, F. Hill, O. Levy, S. Bowman, SuperGlue: A stickier benchmark for general-purpose language understanding systems, *Advances in neural information processing systems* 32 (2019).
- [8] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, et al., Training verifiers to solve math word problems, *arXiv preprint arXiv:2110.14168* (2021).
- [9] E. Glazer, E. Erdil, T. Besiroglu, D. Chicharro, E. Chen, A. Gunning, C. F. Olsson, J.-S. Denain, A. Ho, E. d. O. Santos, et al., Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai, *arXiv preprint arXiv:2411.04872* (2024).
- [10] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, K. R. Narasimhan, Swe-bench: Can language models resolve real-world github issues?, in: *ICLR*, 2024.
- [11] N. Jain, K. Han, A. Gu, W.-D. Li, F. Yan, T. Zhang,

- S. Wang, A. Solar-Lezama, K. Sen, I. Stoica, Livecodebench: Holistic and contamination free evaluation of large language models for code, arXiv preprint arXiv:2403.07974 (2024).
- [12] F. Chollet, On the measure of intelligence, arXiv preprint arXiv:1911.01547 (2019).
- [13] F. Chollet, M. Knoop, G. Kamradt, B. Landers, H. Pinkard, Arc-agi-2: A new challenge for frontier ai reasoning systems, arXiv preprint arXiv:2505.11831 (2025).
- [14] G. Attanasio, P. Basile, F. Borazio, D. Croce, M. Francis, J. Gili, E. Musacchio, M. Nissim, V. Patti, M. Rinaldi, et al., Calamita: Challenge the abilities of language models in italian, in: Proceedings of the 10th Italian Conference on Computational Linguistics (CLiC-it 2024), Pisa, Italy, 2024.
- [15] B. Magnini, R. Zanolli, M. Resta, M. Cimmino, P. Albano, M. Madeddu, V. Patti, Evalita-llm: Benchmarking large language models on italian, arXiv preprint arXiv:2502.02289 (2025).
- [16] I. A. Sag, T. Baldwin, F. Bond, A. Copestake, D. Flickinger, Multiword Expressions: A Pain in the Neck for NLP, Springer Berlin Heidelberg, 2002, p. 1–15. URL: http://dx.doi.org/10.1007/3-540-45715-1_1. doi:10.1007/3-540-45715-1_1.
- [17] S. Tedeschi, F. Martelli, R. Navigli, Id10m: Idiom identification in 10 languages, in: Findings of the Association for Computational linguistics: NAACL 2022, 2022, pp. 2715–2726.
- [18] P. Khoshtab, D. Namazifard, M. Masoudi, A. Akhgary, S. M. Sani, Y. Yaghoobzadeh, Comparative study of multilingual idioms and similes in large language models, in: Proceedings of the 31st International Conference on Computational Linguistics, 2025, pp. 8680–8698.
- [19] J. Kim, Y. Shin, U. Hwang, J. Choi, R. Xuan, T. Kim, Memorization or reasoning? exploring the idiom understanding of llms, arXiv preprint arXiv:2505.16216 (2025).
- [20] J. Lee, A. Liu, O. Ahia, H. Gonen, N. A. Smith, That was the last straw, we need more: Are translation systems sensitive to disambiguating context?, in: Findings of the Association for Computational Linguistics: EMNLP 2023, 2023, pp. 4555–4569.
- [21] S. Li, J. Chen, S. Yuan, X. Wu, H. Yang, S. Tao, Y. Xiao, Translate meanings, not just words: Idiomkb’s role in optimizing idiomatic translation with language models, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, 2024, pp. 18554–18563.
- [22] S. Donthi, M. Spencer, O. B. Patel, J. Y. Doh, E. Rodan, K. Zhu, S. O’Brien, Improving llm abilities in idiomatic translation, in: Proceedings of the First Workshop on Language Models for Low-Resource Languages, 2025, pp. 175–181.
- [23] C. Liu, F. Koto, T. Baldwin, I. Gurevych, Are multilingual llms culturally-diverse reasoners? an investigation into multicultural proverbs and sayings, in: Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), 2024, pp. 2016–2039.
- [24] M. Wang, V.-T. Pham, F. Moghimifar, T.-T. Vu, Proverbs run in pairs: Evaluating proverb translation capability of large language model, arXiv preprint arXiv:2501.11953 (2025).
- [25] F. Caramagna, I 200 proverbi italiani più belli e famosi (con significato), 2025. URL: <https://aforisticamente.com/i-200-proverbi-italiani-piu-belli-e-famosi-con-significato/>.
- [26] OpenRouter, Openrouter: A unified interface for llms, 2024. URL: <https://openrouter.ai/>, accessed: 2025-06-15.
- [27] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al., Gpt-4o system card, arXiv preprint arXiv:2410.21276 (2024).
- [28] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, et al., Deepseek-v3 technical report, arXiv preprint arXiv:2412.19437 (2024).
- [29] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al., Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, arXiv preprint arXiv:2501.12948 (2025).
- [30] Q. Team, Qwen3 technical report, 2025. URL: <https://arxiv.org/abs/2505.09388>. arXiv: 2505.09388.
- [31] G. Team, A. Kamath, J. Ferret, S. Pathak, N. Vieillard, R. Merhej, S. Perrin, T. Matejovicova, A. Ramé, M. Rivière, et al., Gemma 3 technical report, arXiv preprint arXiv:2503.19786 (2025).
- [32] R. Orlando, L. Moroni, P.-L. H. Cabot, S. Conia, E. Barba, S. Orlandini, G. Fiameni, R. Navigli, Minerva llms: The first family of large language models trained from scratch on italian data, in: Proceedings of the 10th Italian Conference on Computational Linguistics (CLiC-it 2024), 2024, pp. 707–719.
- [33] mistralai, Model card for mistral-small-3.1-24b-instruct-2503, 2025. URL: <https://huggingface.co/mistralai/Mistral-Small-3.1-24B-Instruct-2503>.
- [34] Y. Sui, Y.-N. Chuang, G. Wang, J. Zhang, T. Zhang, J. Yuan, H. Liu, A. Wen, S. Zhong, H. Chen, et al., Stop overthinking: A survey on efficient reasoning for large language models, arXiv preprint arXiv:2503.16419 (2025).
- [35] J. Etxaniz, G. Azkune, A. Soroa, O. Lopez de Lacalle, M. Artetxe, Do multilingual language models think better in English?, in: K. Duh, H. Gomez, S. Bethard (Eds.), Proceedings of the

- 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), Association for Computational Linguistics, Mexico City, Mexico, 2024, pp. 550–564. URL: <https://aclanthology.org/2024.naacl-short.46/>. doi:10.18653/v1/2024.naacl-short.46.
- [36] L. Ranaldi, G. Pucci, F. Ranaldi, E. S. Ruzzetti, F. M. Zanzotto, The limits of Italian in reasoning tasks, in: F. Dell’Orletta, A. Lenci, S. Montemagni, R. Sprugnoli (Eds.), *Proceedings of the 10th Italian Conference on Computational Linguistics (CLiC-it 2024)*, CEUR Workshop Proceedings, Pisa, Italy, 2024, pp. 781–795. URL: <https://aclanthology.org/2024.clicit-1.85/>.
 - [37] E. S. Salido, J. Gonzalo, G. Marco, None of the others: a general technique to distinguish reasoning from memorization in multiple-choice llm evaluation benchmarks, *arXiv preprint arXiv:2502.12896* (2025).
 - [38] Y. Chen, J. Benton, A. Radhakrishnan, J. Uesato, C. Denison, J. Schulman, A. Somani, P. Hase, M. Wagner, F. Roger, et al., Reasoning models don’t always say what they think, *arXiv preprint arXiv:2505.05410* (2025).
 - [39] X. Wang, B. Ma, C. Hu, L. Weber-Genzel, P. Röttger, F. Kreuter, D. Hovy, B. Plank, " my answer is c": First-token probabilities do not match text answers in instruction-tuned language models, *arXiv preprint arXiv:2402.14499* (2024).

Supplementary Materials

Proverb Examples

In the following, we report a few examples in Italian from the PROVERBIT benchmark. Each example presents the beginning of a proverb followed by five possible continuations. Only one option is correct, which is *E) Nessuna delle altre risposte* (None of the other answers).

Example 1

- Non è tutto oro...
- A) ...quel che lucida
 - B) ...ciò che brilla
 - C) ...ciò che è opaco
 - D) ...quel che è metallico
 - E) Nessuna delle altre risposte

Example 2

- Vivi e ...
- A) ...lascia ridere
 - B) ...fai decorrere
 - C) ...giudica tutti
 - D) ...morirai più tardi
 - E) Nessuna delle altre risposte

Example 3

- Il gioco è bello quando...
- A) ...cura dopo
 - B) ...è breve
 - C) ...dura molto
 - D) ...tutti si divertono
 - E) Nessuna delle altre risposte

Example 4

- A caval donato ...
- A) ...non si guarda la brocca
 - B) ...non si ispezionan le fauci
 - C) ...si guarda in bocca
 - D) ...si mette la sella
 - E) Nessuna delle altre risposte

Example 5

- Chi non beve in compagnia...
- A) ...o è un lardo o è una spa
 - B) ...o è un delinquente o è uno spione
 - C) ...non è un ladro né una spia
 - D) ...guida la macchina
 - E) Nessuna delle altre risposte

Example 6

- Chi dorme...
- A) ...non piglia cesti
 - B) ...non acchiappa trote
 - C) ...piglia pesci
 - D) ...non è sveglio
 - E) Nessuna delle altre risposte

Complete CoTs

In the following, we report two examples of CoT. In the first, the use of English leads to an incorrect answer, while in the second, an answer mismatch is observed: although the CoT concludes with option C, the model selects option B.

CoT of Qwen

The use of English leads to a wrong answer

Okay, let's see. The user wants me to complete the proverb "Lontano dagli occhi,..." by choosing the correct option from A to E. The options are: A) ...lontano dalle suore, B) ...distante dal miocardio, C) ...vicino al cuore, D) ...si vede poco, E) Nessuna delle altre risposte.

First, I need to recall the actual Italian proverb. The common one I know is "Lontano dagli occhi, vicino al cuore," which translates to "Out of sight, out of mind" but literally means "Far from the eyes, close to the heart." Wait, that doesn't make sense in terms of translation, but the structure is similar to English proverbs where the second part contrasts or complements the first. So the original proverb should have a contrasting part. In Italian, "vicino al cuore" is the correct completion, meaning that even if someone is far from your eyes (not seen), they are still near your heart (in your thoughts or loved).

Now, checking the options. Option C is "...vicino al cuore," which matches what I know. The other options: A) "lontano dalle suore" (far from the nuns) - doesn't make sense in this context. B) "distante dal miocardio" (distant from the myocardium) - that's a medical term, not a proverb. D) "si vede poco" (you see little) - logical but not the traditional saying. E) would only be correct if none of the options are right, but since C seems correct, E isn't the choice here. So the answer must be C.

C.

CoT of DeepSeek

Answer Mismatch w.r.t the CoT

Okay, ho bisogno di completare correttamente il proverbio "Non è tutto oro..." scegliendo tra le opzioni date. Le opzioni sono: A) ...quel che lucida B) ...ciò che brilla C) ...ciò che è opaco D) ...quel che è metallico E) Nessuna delle altre risposte

Prima, mi ricordo come va il proverbio completo. So che il proverbio classico è "Non è tutto oro quel che luccica". Tuttavia, nessuna delle opzioni offre esattamente "luccica".

Guardando le opzioni: A) "quel che lucida" - "Lucida" è simile a "luccica", dato che entrambi significano brillare o splendere. B) "ciò che brilla" - Anche questo significa brillare, quindi è molto vicino al significato originale. C) "ciò che è opaco" - Opaco è l'opposto di brillante, quindi non ha senso in questo contesto. D) "quel che è metallico" - Mentre l'oro è un metallo, il proverbio non si riferisce semplicemente ai metalli ma a ciò che appare prezioso ma non lo è.

L'essenza del proverbio è che non tutto ciò che sembra prezioso o attraente (come l'oro che brilla) è effettivamente di valore. Quindi le opzioni A e B catturano entrambe questo significato.

Confrontando A e B: "quel che lucida" usa "lucida", che è un sinonimo più diretto di "luccica" nel proverbio originale. "Lucida" significa lucidare o brillare, quindi ha lo stesso senso. Invece, "ciò che brilla" (B) è più generico per brillare, ma comunque corretto nel contesto.

Tuttavia, il proverbio originale italiano è tipicamente "Non è tutto oro quel che luccica", quindi "quel che lucida" (A) è la versione più fedele perché "lucida" funziona come sostituto di "luccica". L'opzione E ("Nessuna delle altre risposte") non è necessaria perché A o B sono accettabili. Ma devo vedere quale è la scelta migliore.

Analizzo: Tecnicamente, la versione corretta è con "luccica", ma poiché non è tra le opzioni, "lucida" è un'alternativa valida e comune. Inoltre, molte fonti a volte usano varianti come "brilla", ma "lucida" è più vicino in termini di espressione.

Ho sentito usare sia "luccica" che "lucida" informalmente, quindi A mi sembra corretta. B è anche corretta nel significato, ma non è la formulazione tradizionale. Quindi scelgo A come risposta.

B.

Model	A) Assonant	B) Synonym	C) Inverse	D) Trivial	E) None of the others
GPT o3	2	34	3	0	261
Gemini 2.5 Pro	3	62	5	1	229
DeepSeek R1	2	80	2	3	213
Claude Sonnet 4	3	77	10	0	210
Qwen 3	12	63	32	6	187
GPT 4o	3	78	18	5	196
Claude Sonnet 4	3	101	50	8	138
DeepSeek V3	25	89	51	20	115
Grok 3	16	162	36	4	82
Gemini 2.5 Flash	19	136	82	25	38
LLama 4 Maverick	50	149	60	26	15
Mistral Small 3.1	19	86	56	46	93
Gemma 3	102	96	65	24	13

Table 5

Reasoning model Local models

Absolute number of responses for each error type in the PROVERBIT task.

Declaration on Generative AI

During the preparation of this work, the author(s) used Gemini (Google), Other, and Claude in order to: Grammar and spelling check and Peer review simulation. After using these tool(s)/service(s), the author(s) reviewed and edited the content as needed and take(s) full responsibility for the publication's content.