# Thunder-NUBench: A Benchmark for LLMs' Sentence-Level Negation Understanding

**Yeonkyoung So, Gyuseong Lee, Sungmok Jung, Joonhak Lee,**
**JiA Kang, Sangho Kim, Jaejin Lee**

Graduate School of Data Science,
Seoul National University

{kathy1028,ksnannaya,tjdahrwjd,hmjelee,jia6776,ksh4931,jaejin}@snu.ac.kr

## Abstract

Negation is a fundamental linguistic phenomenon that poses persistent challenges for Large Language Models (LLMs), particularly in tasks requiring deep semantic understanding. Existing benchmarks often treat negation as a side case within broader tasks like natural language inference, resulting in a lack of benchmarks that exclusively target negation understanding. In this work, we introduce **Thunder-NUBench**, a novel benchmark explicitly designed to assess sentence-level negation understanding in LLMs. Thunder-NUBench goes beyond surface-level cue detection by contrasting standard negation with structurally diverse alternatives such as local negation, contradiction, and paraphrase. The benchmark consists of manually curated sentence-negation pairs and a multiple-choice dataset that enables in-depth evaluation of models' negation understanding.

## 1 Introduction

Negation is a fundamental and universal phenomenon found in languages worldwide. It is closely linked to various human communicative abilities, including denial, contradiction, deception, misrepresentation, and irony. Although affirmative statements are more common, negation is still prevalent in language; approximately 25% of sentences in English texts include some form of negation (Sarabi and Blanco, 2016; Hossain et al., 2020; Horn and Wansing, 2025).

Negation plays a crucial role in various natural language processing (NLP) tasks, including sentiment analysis, question answering, knowledge base completion, and natural language inference (NLI). Accurately interpreting negation is vital for understanding semantic oppositions (Khandelwal and Sawant, 2020; Hosseini et al., 2021; Singh et al., 2023). Recent research has shown that the importance of correctly handling negation extends even to multimodal language models (Quantmeyer et al.,

2024; Alhamoud et al., 2025; Park et al., 2025), underscoring its widespread relevance across different domains.

Meanwhile, negation poses significant challenges for both humans and language models. Research shows that people often find negated statements more difficult to process and understand compared to affirmative ones (Wales and Grieve, 1969; Sarabi and Blanco, 2016). Similarly, multiple studies have found that pretrained language models (PLMs) struggle to accurately interpret negation. For example, models like BERT (Devlin et al., 2019) and even large language models (LLMs) such as GPT-3 (Radford et al.) frequently fail to differentiate between negated and affirmative statements. These models often rely on superficial cues, which can lead to incorrect outputs in the presence of negation (Kassner and Schütze, 2020; Hossain et al., 2022a; Truong et al., 2023).

Despite its significance, there is a notable lack of dedicated evaluation benchmarks for understanding negation. Most existing resources treat negation as a minor aspect within broader tasks or focus solely on narrow syntactic detection. Consequently, evaluations have primarily been limited to encoder-based models (Hossain et al., 2020; Geiger et al., 2020; Truong et al., 2022; Anschütz et al., 2023).

To address these shortcomings, we introduce *Thunder-NUBench* (Negation Understanding Benchmark), a dataset designed to assess large language models' ability to interpret negation. The contributions of this paper are summarized as follows:

- We define and categorize various negation phenomena, highlighting their differences from contradiction and paraphrase.
- We introduce a manually curated benchmark to assess the ability of LLMs to understand these distinctions.
- We perform systematic evaluations of several decoder-based LLMs using both prompting

and fine-tuning approaches.

Our benchmark offers valuable insights into the semantic reasoning abilities of language models and serves as a robust evaluation standard for future advancements in understanding negation.

## 2 Related Work

This section reviews existing studies on how language models understand and process negation.

**Negation detection and scope resolution.** Early negation detection and scope resolution work focuses on rule-based systems and handcrafted heuristics, particularly in domain-specific settings like clinical texts. These systems are effective but lack flexibility across domains (Chapman et al., 2001; Carrillo de Albornoz et al., 2012; Ballesteros et al., 2012; Basile et al., 2012). Traditional machine learning methods, such as SVMs (Hearst et al., 1998) and CRFs (Sutton et al., 2012), are later introduced, though they also remain limited to narrow domains (Morante et al., 2008; Morante and Daelemans, 2009; Read et al., 2012; Li and Lu, 2018).

More recently, deep learning approaches employing CNNs (O'shea and Nash, 2015) and BiLSTM networks (Siami-Namini et al., 2019) have improved performance through better contextual embedding and sequence modeling (Fancellu et al., 2016; Bhatia et al., 2020). Pretrained transformer models like BERT (Devlin et al., 2019) have been leveraged through transfer learning (e.g., Neg-BERT (Khandelwal and Sawant, 2020)), significantly enhancing the accuracy of negation detection tasks. However, these methods still primarily address syntactic span detection, with deeper semantic comprehension of negation remaining challenging.

**Negation-sensitive subtasks of NLU.** Negation understanding has become increasingly important in Natural Language Understanding (NLU) tasks (Hosseini et al., 2021). However, existing NLU benchmarks, such as SNLI (Bowman et al., 2015) for NLI, CommonsenseQA (Talmor et al., 2019) for QA, SST-2 (Socher et al., 2013) for sentiment analysis, STS-B (Cer et al., 2017) for textual similarity and paraphrasing, have been criticized for insufficiently accounting for the semantic impact of negation (Hossain et al., 2022a; Rezaei and Blanco, 2024). These datasets contain relatively few negation instances or include negations that are rarely critical to task performance, enabling language models to achieve high accuracy even when ignoring negation entirely.

Recent studies, such as NegNLI (Hossain et al., 2020), MoNLI (Geiger et al., 2020), NaN-NLI (Truong et al., 2022), have introduced negation-sensitive NLU benchmarks, demonstrating that model performance significantly declines when negation meaningfully affects the outcome (Naik et al., 2018; Yanaka et al., 2019; Hartmann et al., 2021; Hossain et al., 2022b; Hossain and Blanco, 2022; She et al., 2023; Anschütz et al., 2023). Such findings indicate that current language models heavily rely on superficial linguistic patterns rather than genuine semantic comprehension.

**Limitations of distributional semantics.** Distributional semantics, the theoretical basis for many PLMs, is built on the distributional hypothesis. That is, words with similar meanings tend to occur in similar contexts (Harris, 1954; Sahlgren, 2008). This assumption enables models to learn semantic representations from textual co-occurrence patterns, making unsupervised training possible (Boleda, 2020; Lenci et al., 2022).

Although powerful in capturing broad semantic relationships, distributional semantics struggles significantly with negation because negated expressions (e.g., "not good") frequently occur in similar contexts as their affirmative forms ("good"), leading models to produce similar vector representations despite their opposite meanings. Previous studies have shown this limitation, highlighting how PLMs fail to capture semantic nuances introduced by antonyms and reversing polarity (Rimell et al., 2017; Jumelet and Hupkes, 2018; Niwa et al., 2021; Jang et al., 2022; Vahtola et al., 2022). Moreover, studies suggest that PLMs like BERT struggle to differentiate between affirmative and negated contexts (Kassner and Schütze, 2020; Ettinger, 2020).

**Negations in generative language models.** Recent research on negation understanding has primarily focused on bidirectional models like BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) due to their strong performance on NLU and negation detection tasks. However, with the rise of generative foundation models such as GPT (Radford et al.) and LLaMA (Touvron et al., 2023), attention has shifted toward evaluating their negation handling. Studies have found that these models often exhibit positive bias and struggle with pro-

| Dimension | Negation Type | Definition | Example |
|---|---|---|---|
| **Scope** | **Clausal Negation ( = Sentential Negation)** | Negation that applies to the entire clause or sentence. This typically involves the use of "not", or its contracted form "n't" with auxiliary verbs. | He **speaks** English fluently. → He **doesn't speak** English fluently. |
| | **Subclausal Negation ( = Constituent / Local Negation)** | Negation that focuses on negating a specific part of a clause, such as a word or phrase, rather than the entire clause. | He speaks English **fluently**. → He speaks English, **but not fluently**. |
| **Form** | **Morphological Negation** | Negation expressed through affixes attached to words such as prefixes like "un-", "in-", "dis-", or suffixes like "-less". | She is **happy**. → She is **unhappy**. |
| | **Syntactic Negation** | Negation expressed through separate words (particles) in the syntax, such as "not", "never", "no", etc. | She is happy. → She is not happy. |
| **Target** | **Verbal Negation** | Negation that applies directly to the verb or verb phrase. | They **have finished** the work. → They **have not finished** the work. |
| | **Non-verbal Negation** | Negation that negates elements other than the verb. | There is **milk** in the fridge. → There is **no milk** in the fridge. |

Table 1: Typology of Negation.

| Contradiction Type | Definition | Example |
|---|---|---|
| **Antonym** | Contradiction caused by opposing meanings of aligned words. | The policy was a **success**. → The policy was a **failure**. |
| **Negation** | One sentence explicitly negates a statement in the other. | She **attended** the meeting. → She **did not attend** the meeting. |
| **Numeric** | Inconsistent numbers, dates, or quantities in related statements. | Totally, **ten** people were injured. → Totally, **five** people were injured. |
| **Factive/Modal** | Conflict in implied facts or modal possibilities due to verbs or auxiliaries. | He **managed to** enter the building. → He **did not enter** the building. |
| **Structure** | Syntactic rearrangement or argument swapping causes contradiction. | **Alice** hired **Bob**. → **Bob** hired **Alice**. |
| **Lexical** | Contradiction through incompatible verbs or phrases, not strictly antonyms. | The manager **praised** her performance. → The manager **expressed disappointment in** her performance. |
| **World Knowledge** | Contradiction relies on common-sense or background knowledge. | The Eiffel Tower is in **Berlin**. → The Eiffel Tower is in **Paris**. |

Table 2: Contradiction types from (de Marneffe et al., 2008). Contradiction covers a broader scope than negation.

ducing or interpreting negated statements (Truong et al., 2023; Chen et al., 2023; García-Ferrero et al., 2023). While benchmarks, such as CON-DAQA (Ravichander et al., 2022) and ScoNe (She et al., 2023), expose these limitations, robust evaluation resources tailored to generative models remain scarce.

Building upon these prior studies, this paper evaluates whether generative models can understand negation in complex sentences and distinguish subtle semantic differences beyond surface-level patterns.

## 3 Definition and Scope of Negation

Although negation has been widely studied in NLP, its definition and scope remain loosely specified, with most studies focusing on identifying negation cues or confusing negation with contradiction. In this work, we aim to refine the definition of negation by examining its semantic boundaries, distin-

guishing it from related but distinct phenomena such as contradiction, and characterizing the types of meaning reversal.

### 3.1 Typology of Negation

Negation is a core semantic and syntactic operation in natural languages expressing a proposition's denial, contradiction, or absence. In formal logic, it reverses the truth value of a proposition: if $P$ is true, then $\neg P$ (which means the negation of $P$) is false, and vice versa. Semantically, negation introduces oppositions, often to a corresponding affirmative proposition (Horn and Wansing, 2025).

Negation can be categorized along several dimensions, including scope, form, and target (see Table 1). By scope, it may affect the entire clause (*clausal negation*) or just a part of it (*subclausal negation*). In terms of form, it can appear as bound morphemes, such as prefixes and suffixes (*morphological negation*), or as separate syntactic elements like "not" or "never" (*syntactic negation*). Finally,

depending on its target, negation can apply to the verb (*verbal negation*) or to other elements in the sentence (*non-verbal negation*) (Zanuttini, 2001; Miestamo, 2007; Truong et al., 2022; Kletz et al., 2023).

## 3.2 Negation and Contradiction

Negation and contradiction, closely related concepts, are often conflated in NLP research (Jiang et al., 2021). Contradiction refers to the incompatibility of two propositions: they cannot be both true simultaneously. Although negation frequently serves as a primary mechanism to create contradictions by reversing the truth value of a proposition, contradictions may also arise through antonymy, numeric mismatch, structural and lexical differences (more details are in Table 2) (de Marneffe et al., 2008). Previous studies have often overlooked the possibility of contradictions existing independently of explicit negation. Recognizing this gap, we specifically examine the ability of LLMs to differentiate between negations and non-negated contradictions, highlighting the nuanced semantic distinctions involved.

## 3.3 Standard Negation

In this paper, we specifically examine *standard negation*: the prototypical form of negation applied to the main declarative verbal clause. Standard negation involves negating the *main verb* in a *main clause*, where the main verb expresses the main action of a clause, and the main clause itself can independently form a complete sentence. This definition excludes negation found in subordinate clauses, which are clauses dependent on a main clause (Miestamo, 2000). This paper specifies the standard negation as reversing the truth value of the main predicate (verbal phrase). Formally, if a main predicate is denoted as $P$, standard negation corresponds precisely to $\neg P$ (i.e., the logical negation of $P$).

Standard negation can encompass various dimensions, as described in Table 1. Specifically, it includes *clausal negation*, which affects the entire sentence, and *verbal negation*, which explicitly targets the main verb. In terms of form, standard negation can be realized either syntactically or morphologically.

*Syntactic negation* typically involves inserting negation particles (e.g., "not") to directly negate the main predicate. On the other hand, *morphological negation* is more limited, applying only when the

antonym of the main predicate fully encompasses its mutually exclusive semantic space (e.g., "be alive" vs. "be dead"). Thus, morphological negation qualifies as standard negation only in cases involving *complementary antonyms*, which represent absolute binary oppositions (e.g., "true" vs. "false" and "possible" vs. "impossible").

In contrast, other types of antonyms, such as gradable antonyms—words that express opposite meanings along a spectrum of quality (e.g., "happy" vs. "unhappy/sad/depressed")—and relational antonyms—words expressing opposite relational roles (e.g., "buy" vs. "sell")—do not strictly reverse truth values (Lehrer and Lehrer, 1982). These examples fall under contradiction rather than the standard negation discussed in this paper.

The above definitions and scope of negation provide the foundation for constructing our benchmark dataset, which we describe in detail in the following section.

## 4 Thunder-NUBench Dataset

We construct the Thunder-NUBench dataset based on two datasets: (1) HoVer dataset, which is designed for multi-hop fact extraction and claim verification based on Wikipedia articles (Jiang et al., 2020), (2) Wikipedia Summary dataset, which contains concise article titles and summaries extracted from English Wikipedia (Scheepers, 2017). We select these datasets as our base corpora because their factual content and complex sentence structures make them well-suited for creating a dataset to understand standard negation in long sentences.

### 4.1 Dataset Generation

The overall process for dataset generation proposed in this paper is illustrated in Figure 1. After extracting sentences from the two sources and preprocess them, we construct two types of datasets: (1) a sentence-negation pair dataset, which includes only standard negation examples generated, and (2) a multiple choice dataset, which covers four categories, systematically constructed through a combination of manual authoring and automated generation. All data is then verified and refined through human review, with a strict protocol to ensure that no author reviews data they generated to ensure quality and consistency. Further details of each step are described below.

**Sentence-negation pair dataset.** We begin by randomly sampling sentences labeled as "supported
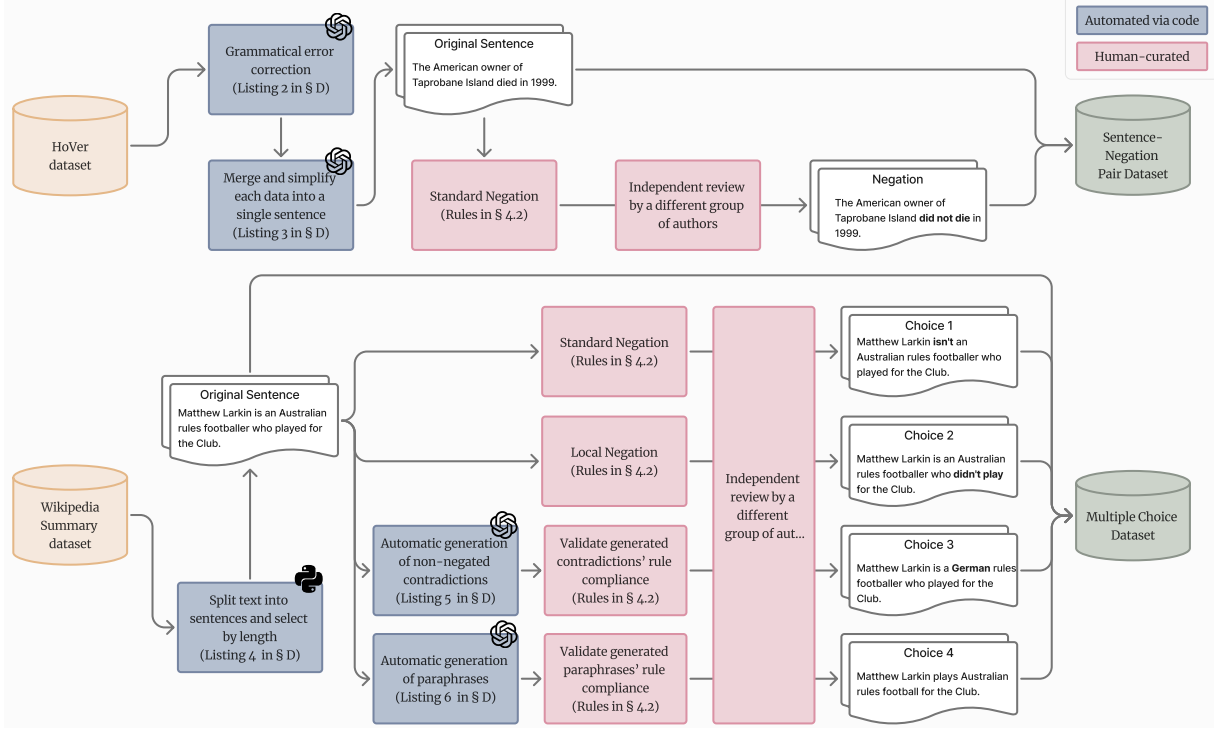
4

Figure 1: Dataset generation process.

facts" from the HoVer dataset. Since the original data often contains grammatical errors, we utilize OpenAI's API (OpenAI, 2025) to automatically correct these issues. In cases where the selected text consists of multiple sentences, we merge or split them as needed to ensure that each example is a single sentence, aligning with our sentence-level task objective (see Appendix D for details). Each sentence is manually negated according to the standard negation criteria described earlier, followed by a thorough review process.

**Multiple choice dataset.** To construct the multiple-choice dataset, we first segment the "summary" column of the Wikipedia Summary dataset into individual sentences, which often contain multiple sentences in a single entry. To focus on the challenges of negation in complex sentences, we filter out sentences that are too short. We generate multiple candidate options for each selected sentence, including standard negation, local negation, contradiction, and paraphrase.

Negation examples are manually written because LLMs often fail to generate accurate negations, frequently producing subclausal negations when standard negation is required, or generating incorrect local negations even when explicitly prompted. As a result, automated generation is not used for these cases. In contrast, non-negated contradiction and

paraphrase examples are first generated automatically using carefully designed prompts with the OpenAI API. All data are further reviewed and refined by the authors.

## 4.2 Human Curation Rules

Below, we describe the principles for manual generation and dataset review. Additional details and examples are provided in Appendix E.

**Standard negation.** The rules for standard negation are as follows:

- The standard negation is intended to reverse the sentence's overall meaning and is achieved by negating the main clause's main verb, reversing the sentence's truth value. All other elements of the main clause (subject, object, temporal context, etc.) are preserved unchanged.
- Standard negation is implemented by inserting negative particles such as "not" into the main verb or substituting the main verb with its complementary antonym, when appropriate.
- The main verb may be replaced with a synonymous verb only if the overall meaning remains strictly identical. Other components may be paraphrased with synonyms as long as the tense, sentence structure, and meaning remain equivalent. Standard negation is then applied

| Type | Structure Explanation | Local Negation Example |
|---|---|---|
| **Relative Clause Negation** | A relative clause is a type of dependent clause that gives extra details about a noun or noun phrase in the main sentence. It usually begins with a relative pronoun such as *who, which, that, whom,* or *whose*. | The man who **owns** the car is my neighbor. → The man who **does not own** the car is my neighbor. |
| **Participle Clause Negation** | A participle clause is a type of dependent clause that begins with a participle (a verb form ending in *-ing* or a past participle). It acts like an adverb, giving extra details about the main clause, often showing time, reason, result, or sequence of actions. | **Walking** through the park, she found a lost wallet. → **Not walking** through the park, she found a lost wallet. |
| **Adverbial Clause Negation** | An adverbial clause is a dependent clause that acts like an adverb, modifying a verb, adjective, or adverb. It gives information such as time, reason, condition, or contrast. These clauses are introduced by subordinating conjunctions like *because, although*, or *while*. | She stayed inside because it **was raining**. → She stayed inside because it **was not raining**. |
| **Compound Sentence with Partial Negation** | A compound sentence consists of two or more main clauses joined by coordinating conjunctions such as *and, but*, or *or*. If only one of these clauses is negated, the negation applies only locally to that clause. | He submitted the report and **attended** the meeting. → He submitted the report and **did not attend** the meeting. |

Table 3: Typology of Local Negation.

to the paraphrased sentence.
- For compound sentences joined by coordinating conjunctions *(and, or, and but)*, negation follows logical rules, such as De Morgan's laws (e.g., A and B → ¬A or ¬B, where A and B are clauses). For unnatural outputs, sentences may be split for fluency, as long as logical negation is preserved.

**Local negation.** We define local negation as a negation that targets a verb phrase outside the main clause. This work applies local negation to four types of sentence structure: relative clause, participle clause, adverbial clause, and compound sentence (see Table 3 for more details). In compound sentences, standard negation requires that all main clauses be negated to achieve sentence-level negation; if only a subset of clauses is negated, it is considered local negation. The mechanism for constructing local negation follows that of standard negation, but the negation is applied only to the intended subpart rather than the entire main clause.

This design allows us to test whether models can distinguish between full-sentence negation and local negation. Although the scope of local negation is restricted, it still contains explicit negation cues (e.g., "not"), which may mislead models that rely on shallow cue detection rather than deeper semantic understanding.

**Contradiction.** Contradiction examples in this work are constructed using mechanisms such as antonymy, numeric changes, or structural alterations, as long as the resulting sentence cannot

be true at the same time as the original and does not simply apply standard or local negation. Unlike standard negation, which reverses the truth value of the main predicate, contradiction can arise from modifying adjectives, quantities, named entities, or other semantic elements. Both negation and contradiction can involve antonyms; however, only complementary antonyms are permitted for standard and local negation, while gradable or relational antonyms are allowed for contradiction. During validation and review, authors ensure that no pair of original and contradictory sentences can be simultaneously true.

It is important to note that standard negation is a subset of contradiction: every negation is a contradiction, but not every contradiction is a standard negation. The goal of this category is to assess whether models can reliably distinguish standard negation from other forms of contradiction, as both alter the meaning of a sentence, but standard negation reverses the entire proposition, whereas contradiction, as defined in this work, does not necessarily do so.

**Paraphrase.** A paraphrase rewrites the sentence using different wording or structure while preserving the original meaning. No additional information may be introduced, and the main verbs and core content must remain unchanged. Identical or near-identical sentences of the original sentence, which often occur in automatically generated paraphrases, are carefully screened and omitted as well.

The reason for including paraphrase examples is to test whether models incorrectly interpret sen-

6

| Dataset | Split | Count |
|---|---|---|
| **Sentence-Negation Pair Dataset** | Train | 3,772 |
| **Multiple Choice Dataset** | Validation | 100 |
| | Test | 1,002 |
| | Total | 4,874 |

Table 4: Thunder-NUBench Dataset Statistics.

tences with different surface forms (e.g., synonyms, rephrased structures) as having reversed meanings. This allows us to examine the robustness of language models in distinguishing genuine reversals from similar meaning, a distinction that has been highlighted as a challenge in previous research on distributional semantics (see Section 2).

### 4.3 Dataset Statistics

The final dataset consists of a sentence-negation pair training set and a multiple-choice evaluation set (see Table 4). The multiple choice dataset presents each original sentence with four options: standard negation (`choice1`, always the answer), local negation (`choice2`), contradiction (`choice3`), and paraphrase (`choice4`).

To construct the validation set of the multiple-choice dataset, we first select 100 examples whose Wikipedia page indices are unique within the dataset, preventing any duplication with the test set. We also matched the distribution of local negation types to the overall dataset as closely as possible, ensuring that the validation set serves as a representative subset. Thunder-NUBench is available online.[1]

### 5 Experiments

We conduct experiments on Thunder-NUBench using an instruction-based prompt that explicitly includes logical rules (as illustrated in Listing 1; more details of prompt selection in Appendix J).

```
1 "Logically negate the sentence below. If the
      sentence includes 'A and B', use 'not A or not
      B'. If it includes 'A or B', use 'not A and not
      B'. Also apply 'not' or use complementary
      antonyms on the main verb(s) of the entire
      sentence.
2 Sentence: {doc['sentence']}
3 Negation:"
```

Listing 1: Instruction-based prompt format.

|  |  | Baseline | | SFT on Thunder-NUBench |
|---|---|---|---|---|
|  |  | zeroshot | fewshot (±SD) | zeroshot |
| **Qwen2.5-3B** | acc | 0.495 | 0.600 (±0.009) | 0.747 |
| | acc_norm | 0.526 | 0.622 (±0.009) | 0.753 |
| **Qwen2.5-3B-Instruct** | acc | 0.667 | 0.758 (±0.006) | 0.765 |
| | acc_norm | 0.686 | 0.771 (±0.006) | 0.775 |
| **Qwen2.5-7B** | acc | 0.499 | 0.623 (±0.009) | 0.820 |
| | acc_norm | 0.521 | 0.633 (±0.012) | 0.825 |
| **Qwen2.5-7B-Instruct** | acc | 0.545 | 0.700 (±0.004) | 0.783 |
| | acc_norm | 0.576 | 0.713 (±0.002) | 0.784 |
| **Average** | **acc** | **0.551** | **0.670** | **0.779** |
| | **acc_norm** | **0.577** | **0.685** | **0.784** |

Table 5: Evaluation results of Qwen2.5 models across different settings: baseline zero-shot, baseline few-shot (5-shot), and zero-shot after supervised fine-tuning (SFT) on Thunder-NUBench. Both accuracy (acc) and normalized accuracy (acc_norm) are reported. Few-shot results are averaged over five random seeds with the standard deviation in parentheses.

**Zero-shot and few-shot.** For each model, we evaluate both zero-shot and few-shot settings using the Language Model Evaluation Harness (Gao et al., 2024). In the few-shot scenario, we use examples from the validation set as in-context demonstrations. In few-shot experiments, results are averaged over five random seeds (42, 1234, 3000, 5000, and 7000) and five examples of validation set data (5 shots). We report performance on the test set for each model and prompt configuration. All the results are provided in Appendix M.

**SFT.** We perform Supervised Fine-Tuning (SFT) using the LLaMA-Factory framework (Zheng et al., 2024) on the Sentence-Negation Pair Dataset from Thunder-NUBench. The data is formatted in the Alpaca instruction style (Taori et al., 2023). For parameter-efficient training, we apply Low-Rank Adaptation (LoRA) (Hu et al., 2022) with a rank of 8, targeting all linear layers. Fine-tuning is conducted for three epochs with a batch size of 1, a gradient accumulation step of 8, cosine learning rate scheduling, and bfloat16 precision. After SFT, we evaluate zero-shot performance to directly measure the model's ability to generalize from instruction tuning without the influence of in-context examples. All the results are provided in Appendix N.

| | | | Error rate (1-acc) | Incorrect Choice Distribution | | | choice2 confusion rate | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | choice2 (%) | choice3 (%) | choice4 (%) | relative_part confuse (%) | pp_part confuse (%) | compound_part confuse (%) | adverb_part confuse (%) |
| **Llama-3.1-8B** | **baseline** | zeroshot | 0.523 | 70.80 | 22.90 | 6.30 | 29.35 | 34.83 | 53.43 | 34.29 |
| | | fewshot | 0.357 | 81.01 | 16.20 | 2.79 | 22.53 | 21.38 | 47.29 | 29.52 |
| | **After SFT** | zeroshot | 0.149 | 80.54 | 16.11 | 3.36 | 9.56 | 10.34 | 18.77 | 9.52 |
| **gemma-7b** | **baseline** | zeroshot | 0.508 | 68.17 | 28.88 | 2.95 | 28.67 | 33.10 | 46.93 | 35.24 |
| | | fewshot | 0.382 | 83.81 | 14.62 | 1.57 | 26.28 | 23.45 | 49.10 | 38.10 |
| | **After SFT** | zeroshot | 0.153 | 77.78 | 20.26 | 1.96 | 9.56 | 10.34 | 19.49 | 6.67 |
| **Qwen2.5-7B** | **baseline** | zeroshot | 0.501 | 66.14 | 30.08 | 3.78 | 31.40 | 33.45 | 40.07 | 30.48 |
| | | fewshot | 0.387 | 74.23 | 24.23 | 1.55 | 24.91 | 24.48 | 41.52 | 27.62 |
| | **After SFT** | zeroshot | 0.180 | 85.56 | 13.33 | 1.11 | 14.33 | 15.52 | 20.22 | 10.48 |
| **Mistral-7B-v0.3** | **baseline** | zeroshot | 0.521 | 72.03 | 22.61 | 5.36 | 26.96 | 34.83 | 58.12 | 33.33 |
| | | fewshot | 0.372 | 79.09 | 18.23 | 2.68 | 22.87 | 22.76 | 48.01 | 27.62 |
| | **After SFT** | zeroshot | 0.161 | 83.85 | 13.66 | 2.48 | 8.87 | 11.03 | 25.27 | 6.67 |

Table 6: Incorrect choice distribution and confusion analysis in negation benchmark across various 7-8B size pretrained models. Few-shot results are reported with a fixed seed (1234) to keep error patterns clear because averaging over multiple seeds could make them harder to interpret.

**Benchmarking with Thunder-NUBench.** Table 5 presents evaluation results for Qwen2.5-3B, Qwen2.5-3B-Instruct, Qwen2.5-7B, and Qwen2.5-7B-Instruct models (Qwen et al., 2025) in three settings: zero-shot baseline, few-shot baseline, and zero-shot after SFT using Thunder-NUBench.

Instruction-tuned models consistently outperform their pretrained counterparts. Few-shot prompting significantly improves performance, highlighting the benefit of concrete examples. Supervised fine-tuning on Thunder-NUBench further boosts accuracy, with pretrained models showing the most significant gains. Notably, Qwen2.5-3B-Instruct performs exceptionally well in zero-shot and few-shot settings, even outperforms larger models, suggesting strong alignment with the logical reasoning demands of Thunder-NUBench.

**Negation understanding analysis.** We analyze model errors to assess their ability to distinguish standard negation from similar semantic variants. Each local negation subtype in our dataset is explicitly labeled according to its sentence structure: relative clauses (relative_part), participle clauses (pp_part), compound sentences (compound_part), and adverbial clauses (adverb_part). To identify which subtypes are most often confused with standard negation, we compute the *confusion rate*, defined as the proportion of examples within each subtype where the model incorrectly selects the local negation option (choice2) instead of the correct standard negation (choice1).

For instance, among 1,002 test examples, 290 of choice2 are labeled as pp_part; if the model erroneously chooses choice2 in 29 of these, the confusion rate for pp_part is 10%. For all the details, please see Appendix P.

We conduct a comparative analysis across four 7-8B scale pretrained language models (LLaMA-3.1-8B, Gemma-7B, Qwen2.5-7B, and Mistral-7B-v0.3) (Grattafiori et al., 2024; Team et al., 2025; Qwen et al., 2025; Jiang et al., 2023) under three evaluation settings: baseline (zero-shot and few-shot) and zero-shot after SFT. All models show a consistent tendency to incorrectly select local negation (choice2), even though SFT reduces overall errors, showing that distinguishing local from full-sentence negation remains difficult. Confusion rate is especially high for compound sentence structures, highlighting specific areas where models systematically struggle with negation understanding.

## 6 Conclusion

In this work, we introduce Thunder-NUBench, a benchmark designed to evaluate large language models(LLMs)' sentence-level understanding of negation. Unlike prior sources that isolate negation in simple syntactic contexts, Thunder-NUBench incorporates structurally diverse alternatives, such as local negation, contradictions, and paraphrases, to assess models' deeper semantic comprehension. Our evaluations demonstrate that while LLMs benefit from few-shot prompting and fine-tuning, they struggle to distinguish standard negation from closely related semantic variations. Thunder-NUBench provides a valuable diagnostic tool for analyzing models' limitations regarding negation understanding, and serves as a robust benchmark

for future research.

## Limitations

Thunder-NUBench is exclusively constructed in English, despite negation being a universal linguistic phenomenon demonstrated in diverse forms across languages. The syntactic and semantic expressions of negation may vary in other languages, meaning that our current findings may not generalize to multilingual or cross-lingual settings. In future work, we aim to extend the research to a broader range of languages to enable cross-linguistic evaluation of negation understanding in language models.

Although we built the dataset using two distinct sources (HoVer and Wikipedia summaries), both are derived from encyclopedic, formal domains, which may not fully represent the variety of sentence structures and informal language found in real-world use cases. Moreover, while all examples were systematically generated and reviewed, some bias may persist due to subjective decisions in the human curation process. We attempt to mitigate this through cross-checking by an independent group of authors, but some residual bias may remain.

Thunder-NUBench primarily focuses on standard (sentence-level) negation and its distinction from local negation, contradiction, and paraphrase. Other important negation phenomena, such as double negation and negative polarity items (NPIs), are not directly addressed in this benchmark. Our current focus is establishing a strong foundation for evaluating models' understanding of standard negation. However, we aim to expand the evaluation to a broader range of negation phenomena in future work.

## Ethical Considerations

This work does not involve the use of crowdsourcing methods. Instead, all data included in the Thunder-NUBench benchmark was carefully reviewed by the authors to ensure quality, relevance, and adherence to ethical standards. The datasets and tools used for training and evaluation are publicly available and used in compliance with their respective licenses.

When leveraging OpenAI's text generation models, we take additional care to avoid generating or including any content that is harmful, biased, or violates privacy. All generated examples are manu-

ally reviewed to meet ethical and safety standards. We ensure no personally identifiable information or offensive content is present in the final dataset.

The Thunder-NUBench dataset is released under the CC BY-NC-SA 4.0 license, ensuring transparency, reproducibility, and accessibility for future research. We believe our work contributes positively to developing trustworthy and interpretable language models.

## Acknowledgments

## References

Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip Torr, Yoon Kim, and Marzyeh Ghassemi. 2025. Vision-language models do not understand negation. *arXiv preprint arXiv:2501.09425*.

Miriam Anschütz, Diego Miguel Lozano, and Georg Groh. 2023. This is not correct! negation-aware evaluation of language generation systems. In *Proceedings of the 16th International Natural Language Generation Conference*, pages 163–175, Prague, Czechia. Association for Computational Linguistics.

Miguel Ballesteros, Alberto Díaz, Virginia Francisco, Pablo Gervás, Jorge Carrillo de Albornoz, and Laura Plaza. 2012. UCM-2: a rule-based approach to infer the scope of negation via dependency parsing. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and*

*Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 288–293, Montréal, Canada. Association for Computational Linguistics.

Valerio Basile, Johan Bos, Kilian Evang, and Noortje Venhuizen. 2012. Ugroningen: Negation detection with discourse representation structures. In *\* SEM 2012: The First Joint Conference on Lexical and Computational Semantics–Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 301–309.

Parminder Bhatia, E Busra Celikkaya, and Mohammed Khalilia. 2020. End-to-end joint entity extraction and negation detection for clinical text. *Precision Health and Medicine: A Digital Revolution in Healthcare*, pages 139–148.

Gemma Boleda. 2020. Distributional semantics and linguistic theory. *Annual Review of Linguistics*, 6(1):213–234.

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.

Jorge Carrillo de Albornoz, Laura Plaza, Alberto Díaz, and Miguel Ballesteros. 2012. UCM-I: A rule-based syntactic approach for resolving the scope of negation. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 282–287, Montréal, Canada. Association for Computational Linguistics.

Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. *arXiv preprint arXiv:1708.00055*.

Wendy W Chapman, Will Bridewell, Paul Hanbury, Gregory F Cooper, and Bruce G Buchanan. 2001. A simple algorithm for identifying negated findings and diseases in discharge summaries. *Journal of biomedical informatics*, 34(5):301–310.

Jiangjie Chen, Wei Shi, Ziquan Fu, Sijie Cheng, Lei Li, and Yanghua Xiao. 2023. Say what you mean! large language models speak too positively about negative commonsense knowledge. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9890–9908, Toronto, Canada. Association for Computational Linguistics.

Marie-Catherine de Marneffe, Anna N. Rafferty, and Christopher D. Manning. 2008. Finding contradictions in text. In *Proceedings of ACL-08: HLT*, pages 1039–1047, Columbus, Ohio. Association for Computational Linguistics.

Viviane Déprez, Susagna Tubau, Anne Cheylus, and M Teresa Espinal. 2015. Double negation in a negative concord language: An experimental investigation. *Lingua*, 163:75–107.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Allyson Ettinger. 2020. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.

Federico Fancellu, Adam Lopez, and Bonnie Webber. 2016. Neural networks for negation scope detection. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 495–504, Berlin, Germany. Association for Computational Linguistics.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. The language model evaluation harness.

Iker García-Ferrero, Begoña Altuna, Javier Alvez, Itziar Gonzalez-Dios, and German Rigau. 2023. This is not a dataset: A large negation benchmark to challenge large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8596–8615, Singapore. Association for Computational Linguistics.

Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. Neural natural language inference models partially embed theories of lexical entailment and negation. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 163–173, Online. Association for Computational Linguistics.

Lila R Gleitman. 1965. Coordinating conjunctions in english. *Language*, 41(2):260–293.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Zellig S Harris. 1954. Distributional structure. *Word*, 10(2-3):146–162.

Mareike Hartmann, Miryam de Lhoneux, Daniel Hershcovich, Yova Kementchedjhieva, Lukas Nielsen,

Chen Qiu, and Anders Søgaard. 2021. A multilingual benchmark for probing negation-awareness with minimal pairs. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 244–257, Online. Association for Computational Linguistics.

Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. 1998. Support vector machines. *IEEE Intelligent Systems and their applications*, 13(4):18–28.

Kees Hengeveld. 1986. Copular verbs in a functional grammar of spanish.

Laurence R. Horn and Heinrich Wansing. 2025. Negation. In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*, Spring 2025 edition. Metaphysics Research Lab, Stanford University.

Md Mosharaf Hossain and Eduardo Blanco. 2022. Leveraging affirmative interpretations from negation improves natural language understanding. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5833–5847, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Md Mosharaf Hossain, Dhivya Chinnappa, and Eduardo Blanco. 2022a. An analysis of negation in natural language understanding corpora. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 716–723, Dublin, Ireland. Association for Computational Linguistics.

Md Mosharaf Hossain, Luke Holman, Anusha Kakileti, Tiffany Kao, Nathan Brito, Aaron Mathews, and Eduardo Blanco. 2022b. A question-answer driven approach to reveal affirmative interpretations from verbal negations. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 490–503, Seattle, United States. Association for Computational Linguistics.

Md Mosharaf Hossain, Venelin Kovatchev, Pranoy Dutta, Tiffany Kao, Elizabeth Wei, and Eduardo Blanco. 2020. An analysis of natural language inference benchmarks through the lens of negation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9106–9118, Online. Association for Computational Linguistics.

Arian Hosseini, Siva Reddy, Dzmitry Bahdanau, R Devon Hjelm, Alessandro Sordoni, and Aaron Courville. 2021. Understanding by understanding not: Modeling negation in language models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1301–1312, Online. Association for Computational Linguistics.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.

Myeongjun Jang, Frank Mtumbuka, and Thomas Lukasiewicz. 2022. Beyond distributional hypothesis: Let language models learn meaning-text correspondence. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2030–2042, Seattle, United States. Association for Computational Linguistics.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Liwei Jiang, Antoine Bosselut, Chandra Bhagavatula, and Yejin Choi. 2021. "I'm not mad": Commonsense implications of negation and contradiction. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4380–4397, Online. Association for Computational Linguistics.

Yichen Jiang, Shikha Bordia, Zheng Zhong, Charles Dognin, Maneesh Singh, and Mohit Bansal. 2020. HoVer: A dataset for many-hop fact extraction and claim verification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3441–3460, Online. Association for Computational Linguistics.

Jaap Jumelet and Dieuwke Hupkes. 2018. Do language models understand anything? on the ability of LSTMs to understand negative polarity items. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 222–231, Brussels, Belgium. Association for Computational Linguistics.

Nora Kassner and Hinrich Schütze. 2020. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7811–7818, Online. Association for Computational Linguistics.

Aditya Khandelwal and Suraj Sawant. 2020. NegBERT: A transfer learning approach for negation detection and scope resolution. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5739–5748, Marseille, France. European Language Resources Association.

David Kletz, Pascal Amsili, and Marie Candito. 2023. The self-contained negation test set. In *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 212–221, Singapore. Association for Computational Linguistics.

Adrienne Lehrer and Keith Lehrer. 1982. Antonymy. *Linguistics and philosophy*, pages 483–501.

Alessandro Lenci, Magnus Sahlgren, Patrick Jeuniaux, Amaru Cuba Gyllensten, and Martina Miliani. 2022. A comparative evaluation and analysis of three generations of distributional semantic models. *Language resources and evaluation*, 56(4):1269–1313.

Hao Li and Wei Lu. 2018. Learning with structured representations for negation scope extraction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 533–539, Melbourne, Australia. Association for Computational Linguistics.

Yinheng Li. 2023. A practical survey on zero-shot prompt design for in-context learning. In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 641–647, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Matti Miestamo. 2000. Towards a typology of standard negation. *Nordic journal of linguistics*, 23(1):65–88.

Matti Miestamo. 2007. Negation–an overview of typological research. *Language and linguistics compass*, 1(5):552–570.

Roser Morante and Walter Daelemans. 2009. A metalearning approach to processing the scope of negation. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL-2009)*, pages 21–29, Boulder, Colorado. Association for Computational Linguistics.

Roser Morante, Anthony Liekens, and Walter Daelemans. 2008. Learning the scope of negation in biomedical texts. In *Proceedings of the 2008 conference on empirical methods in natural language processing*, pages 715–724.

Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Ha Thanh Nguyen, Randy Goebel, Francesca Toni, Kostas Stathis, and Ken Satoh. 2023. A negation detection assessment of gpts: analysis with the xnot360 dataset. *arXiv preprint arXiv:2306.16638*.

Ayana Niwa, Keisuke Nishiguchi, and Naoaki Okazaki. 2021. Predicting antonyms in context using BERT. In *Proceedings of the 14th International Conference*

on *Natural Language Generation*, pages 48–54, Aberdeen, Scotland, UK. Association for Computational Linguistics.

OpenAI. 2025. Text generation and prompting. Accessed on May 16, 2025.

Keiron O'shea and Ryan Nash. 2015. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*.

Junsung Park, Jungbeom Lee, Jongyoon Song, Sangwon Yu, Dahuin Jung, and Sungroh Yoon. 2025. Know" no" better: A data-driven approach for enhancing negation awareness in clip. *arXiv preprint arXiv:2501.10913*.

Vincent Quantmeyer, Pablo Mosteiro, and Albert Gatt. 2024. How and where does CLIP process negation? In *Proceedings of the 3rd Workshop on Advances in Language and Vision Research (ALVR)*, pages 59–72, Bangkok, Thailand. Association for Computational Linguistics.

Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report. *Preprint*, arXiv:2412.15115.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training.

Abhilasha Ravichander, Matt Gardner, and Ana Marasović. 2022. Condaqa: A contrastive reading comprehension dataset for reasoning about negation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8729–8755.

Jonathon Read, Erik Velldal, Lilja Øvrelid, and Stephan Oepen. 2012. UiO1: Constituent-based discriminative ranking for negation resolution. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 310–318, Montréal, Canada. Association for Computational Linguistics.

MohammadHossein Rezaei and Eduardo Blanco. 2024. Paraphrasing in affirmative terms improves negation understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 602–615, Bangkok, Thailand. Association for Computational Linguistics.

Laura Rimell, Amandla Mabona, Luana Bulat, and Douwe Kiela. 2017. Learning to negate adjectives with bilinear models. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 71–78, Valencia, Spain. Association for Computational Linguistics.

Magnus Sahlgren. 2008. The distributional hypothesis. *Italian Journal of linguistics*, 20:33–53.

Zahra Sarabi and Eduardo Blanco. 2016. Understanding negation in positive terms using syntactic dependencies. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1108–1118, Austin, Texas. Association for Computational Linguistics.

Thijs Scheepers. 2017. Improving the compositionality of word embeddings. Master's thesis, Universiteit van Amsterdam, Science Park 904, Amsterdam, Netherlands, 11.

Jingyuan S. She, Christopher Potts, Samuel R. Bowman, and Atticus Geiger. 2023. ScoNe: Benchmarking negation reasoning in language models with fine-tuning and in-context learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1803–1821, Toronto, Canada. Association for Computational Linguistics.

Sima Siami-Namini, Neda Tavakoli, and Akbar Siami Namin. 2019. The performance of lstm and bilstm in forecasting time series. In *2019 IEEE International conference on big data (Big Data)*, pages 3285–3292. IEEE.

Rituraj Singh, Rahul Kumar, and Vivek Sridhar. 2023. NLMs: Augmenting negation in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13104–13116, Singapore. Association for Computational Linguistics.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.

Charles Sutton, Andrew McCallum, et al. 2012. An introduction to conditional random fields. *Foundations and Trends® in Machine Learning*, 4(4):267–373.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Thinh Hung Truong, Timothy Baldwin, Karin Verspoor, and Trevor Cohn. 2023. Language models are not naysayers: an analysis of language models on negation benchmarks. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, pages 101–114, Toronto, Canada. Association for Computational Linguistics.

Thinh Hung Truong, Yulia Otmakhova, Timothy Baldwin, Trevor Cohn, Jey Han Lau, and Karin Verspoor. 2022. Not another negation benchmark: The NaN-NLI test suite for sub-clausal negation. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 883–894, Online only. Association for Computational Linguistics.

Teemu Vahtola, Mathias Creutz, and Jörg Tiedemann. 2022. It is not easy to detect paraphrases: Analysing semantic similarity with antonyms and negation using the new SemAntoNeg benchmark. In *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 249–262, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ton Van der Wouden. 1996. Litotes and downward monotonicity. *Negation: a notion in focus*, 7:145.

RJ Wales and R Grieve. 1969. What is so difficult about negation? *Perception & Psychophysics*, 6:327–332.

Xingchen Wan, Ruoxi Sun, Hanjun Dai, Sercan Arik, and Tomas Pfister. 2023. Better zero-shot reasoning with self-adaptive prompting. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3493–3514, Toronto, Canada. Association for Computational Linguistics.

Hitomi Yanaka, Koji Mineshima, Daisuke Bekki, Kentaro Inui, Satoshi Sekine, Lasha Abzianidze, and Johan Bos. 2019. Can neural networks understand monotonicity reasoning? In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 31–40, Florence, Italy. Association for Computational Linguistics.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.

Raffaella Zanuttini. 2001. Sentential negation. *The handbook of contemporary syntactic theory*, pages 511–535.

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*, pages 12697–12706. PMLR.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

## A    Copular Verbs

Copular verbs, also known as linking verbs, are verbs that connect the subject of a sentence to a subject complement, which can be a noun, adjective, or other expression that describes or identifies the subject. Unlike action verbs, copular verbs do not express actions but rather states or conditions. The most common copular verb in English is "to be" in its various forms (am, is, are, was, were). Other examples include "seem," "appear," "become," "feel," "look," "sound," "taste," and "smell" when used to describe the subject's state (Hengeveld, 1986).

In the context of standard negation, as discussed in Section 3.3, negation typically targets the main predicate of a clause. For sentences with copular verbs, this means that the entire verb phrase, including the copular verb and its complement, is subject to negation. For example, in the sentence "She is a doctor," the main predicate is "is a doctor." Negating this sentence results in "She is not a doctor," where the negation applies to the entire predicate, not just the verb "is."

This approach aligns with the semantic characterization of standard negation, where the truth value of the main predicate is reversed ($P \to \neg P$). In copular constructions, the negation encompasses both the copular verb and its complement, effectively altering the state or identity attributed to the subject.

In the case of copular verbs, standard negation targets the entire predicate, including the complement. This means that negation can be realized either syntactically (e.g., "is **not** an expert") or by replacing the complement with its complementary antonym (e.g., "is a **non-expert**"), both of which result in the reversal of the main predicate's truth value. Although such constructions may superficially appear to be non-verbal negation, especially when the complement is a noun or adjective, they are, in fact, instances of verbal negation, since the negation applies to the predicate as a whole.

## B    HoVer Dataset

The HoVer (**Ho**ppy **Ver**ification) dataset is developed for the tasks of multi-hop evidence retrieval and factual claim verification. In HoVer, each claim requires supporting evidence that spans multiple English Wikipedia articles to determine whether the claim is substantiated or not. The dataset is distributed under a CC BY-SA 4.0 License, and it can be accessed via its official homepage[2]. Table 7 offers an overview of the dataset's structure. The data is split into training, validation, and test sets, containing 18,171, 4,000, and 4,000 examples respectively.

HoVer is constructed on top of the HotpotQA dataset, which is designed to evaluate multi-hop reasoning in question answering. HotpotQA itself is a large-scale collection of Wikipedia-based QA pairs created to address the limitations of prior QA datasets, which often fail to require complex reasoning or explanatory answers (Yang et al., 2018). The construction of HoVer involves rewriting HotpotQA question-answer pairs into claim statements, which are then validated and labeled by annotators. Claims are extended to require multi-hop evidence from up to four Wikipedia articles and are systematically modified to increase complexity. Final labels are assigned as SUPPORTED or NOT-SUPPORTED (Jiang et al., 2020).

---

[2]https://hover-nlp.github.io/

| Column | Detail | Example |
|---|---|---|
| id | Unique claim identifier | 0 |
| uid | User/annotator identifier | 330ca632-e83f-4011-b11b-0d0158145036 |
| claim | The statement to be verified, often requiring multi-article evidence | Skagen Painter Peder Severin Krøyer favored naturalism along with Theodor Esbern Philipsen and the artist Ossian Elgström studied with in the early 1900s. |
| supporting_facts | List of Wikipedia article titles and sentence indices providing evidence | [ { "key": "Kristian Zahrtmann", "value": 0 }, { "key": "Kristian Zahrtmann", "value": 1 }, { "key": "Peder Severin Krøyer", "value": 1 }, { "key": "Ossian Elgström", "value": 2 } ] |
| label | Whether the claim is supported | 1: SUPPORTED or 0: NOT_SUPPORTED |
| num_hops | Number of articles required for verification | 2∼4 |
| hpqa_id | Reference to the original HotpotQA pair | 5ab7a86d5542995dae37e986 |

Table 7: Details of HoVer dataset structure with examples.

| Column | Detail | Example |
|---|---|---|
| title | Article title from Wikipedia. | Alain Connes |
| description | A brief description or category for the article (when available). | French mathematician |
| summary | The extracted summary or introduction section of the article, typically more concise than the full text. | Alain Connes (; born 1 April 1947) is a French mathematician... |
| full_text | The complete article text (when included), encompassing the full body of the Wikipedia page. | Alain Connes (; born 1 April 1947) is a French mathematician... |
| __index_level_0__ | Index number for each entry in the dataset. | 3 |

Table 8: Details of Wikipeda Summary dataset structure with examples.

## C    Wikipedia Summary Dataset

The Wikipedia Summary Dataset contains the titles and introductory summaries of English Wikipedia articles, extracted in September 2017. A summary or introduction in this context refers to the content from the article title up to the content outline (i.e., before the first section heading). The dataset was originally released via GitHub[3], but is now accessible through the Hugging Face Hub[4]. The dataset license is not explicitly mentioned, but as the original Wikipedia data is distributed under the CC BY-SA 4.0, it is assumed that the dataset would be distributed under the same license. For licensing details, refer to the Wikimedia Terms of Use [5]. Table 8 offers an overview of the dataset's structure. The dataset comprises approximately 430,000 articles, only providing the training set (Scheepers, 2017).

## D    Codes for Data Generation

### D.1    Sentence-Negation Pair Dataset

To construct the sentence-negation pair dataset, we leverage OpenAI GPT models (OpenAI, 2025) for both grammar correction and sentence merging.

We select different model versions depending on the complexity of each task. For sentence merging, which demands nuanced contextual understanding and complex syntax, we use GPT-4. For grammar correction, where edits are more straightforward, GPT-3.5 is sufficient.

```
def grammar_fix(claim):
    messages = [{"role": "system", "content": "Fix grammatical errors."},
    {"role": "user", "content": f"If there are errors, please fix the sentence: {claim} \n If there aren't, return the original sentence. Provide only the resulting sentence without any additional explanation or introduction."}]
    response = client.chat.completions.create(model="gpt-3.5-turbo", messages=messages)
    fixed_text = response.choices[0].message.content.strip()
    return fixed_text
```

Listing 2: Fixing Grammar with OpenAI API.

```
def merge_sentences_with_gpt(claim):
    messages = [{"role": "system", "content": "Merge sentences into a single one."},
    {"role": "user", "content": f"Merge these sentences: {claim} \n Provide only the resulting sentence without any additional explanation or introduction."}]
    response = client.chat.completions.create(model="gpt-4-turbo-preview", messages=messages)
    merged_text = response.choices[0].message.content.strip()
    return merged_text
```

Listing 3: Merging Sentences with OpenAI API.

### D.2    Multiple Choice Dataset

To generate the multiple choice dataset, we first split the large text into single sentences using

---

[3] https://github.com/tscheepers/Wikipedia-Summary-Dataset
[4] https://huggingface.co/datasets/jordiclive/wikipedia-summary-dataset
[5] https://foundation.wikimedia.org/wiki/Policy:Terms_of_Use

Python and then automatically generate contradictions and paraphrases for each sentence via the OpenAI API (GPT-4o), followed by human review. The following scripts illustrate the procedures.

```python
import pandas as pd
import re
from datasets import load_dataset
import random

df = pd.DataFrame(load_dataset("jordiclive/wikipedia
    -summary-dataset")['train'].shuffle(seed=42).
    select(range(10000)))
df = df.drop(columns=['full_text'])

def split_into_sentences(text):
    sentences = re.split(r'(?<=[.!?]) +', text)
    return sentences

df['sentence'] = df['summary'].apply(
    split_into_sentences)
df = df.explode('sentence')
df = df[df['sentence'].apply(lambda x: len(x.split()
    ) >= 30)]
df = df.reset_index(drop=True)
df.to_csv("file/wikipedia_summary_sentences.csv",
    index=False)
```

Listing 4: Sentence extraction and preprocessing from Wikipedia summaries.

```python
def generate_contradiction(sentence):
    prompt = f"""
    You will be given a sentence. Generate a
    contradictory sentence that directly conflicts
    with the original sentence without using
    standard negation.

    Definitions:
    - Standard negation: Directly negating the main
    verb or using words like 'not', 'no', 'never',
    or negative contractions such as \"isn't\", \"
    doesn't\", or \"can't\".
    - Contradiction: A sentence that logically
    conflicts with the original statement. The
    contradiction must be such that both sentences
    cannot logically be true at the same time under
     any circumstances.

    Important:
    - Do not change the main verb from the original
    sentence.
    - Do not use 'never' or other negative words to
    form the contradiction.
    - Ensure the contradicted sentence logically
    excludes the possibility of the original
    sentence being true simultaneously.

    Examples:
    Original sentence: \"The tallest student won the
     award.\"
    Contradicted sentence: \"The shortest student
    won the award.\"

    Original sentence: \"The room was completely
    dark.\"
    Contradicted sentence: \"The room was brightly
    lit.\"

    Original sentence: \"The event took place in the
     morning.\"
    Contradicted sentence: \"The event took place in
     the evening.\"

    Original sentence: \"All people are dying.\"
    Contradicted sentence: \"Some people are dying
    .\"

    Now, generate a contradictory sentence without
    standard negation, without changing the main
    verb, and ensuring the two sentences are
    logically incompatible, for the following:

    Original sentence: \"{sentence}\"
```

```python
    Contradicted sentence:
    """

    completion = client.chat.completions.create(
        model="gpt-4o",
        messages=[
            {"role": "system", "content": "You are a
    helpful assistant tasked with generating
    logical contradictions. Do not use negation to
    make contradiction."},
            {"role": "user", "content": prompt}
        ]
    )

    return completion.choices[0].message.content
```

Listing 5: Contradiction generation with GPT-4o.

```python
def generate_paraphrase(sentence):
    prompt = f"""
        Paraphrase the following sentence using
    synonyms or slight structural variations
    without changing its meaning.
        Do not add or remove any main verbs. Keep
    the original intent of the sentence intact.

        Original sentence: "{sentence}"

        Paraphrased sentence:
        """

    completion = client.chat.completions.create(
        model="gpt-4o",
        messages=[
            {"role": "system", "content": "You are a
    helpful assistant skilled at generating
    paraphrases while keeping the meaning of
    sentences unchanged."},
            {
                "role": "user",
                "content": prompt
            }
        ]
    )

    return completion.choices[0].message.content
```

Listing 6: Paraphrase generation with GPT-4o.

# E    Detailed Principles and Examples of the Thunder-NUBench

Standard negation refers to the process of reversing the truth value of the main predicate of a sentence, typically by directly negating the main verb in the main clause, or (in some cases) by using a complementary antonym. The goal is to deny the entire proposition while preserving all other elements of the sentence (subject, objects, temporal context, etc.) unchanged.

**Paraphrasing before Negation.**  Before negating, the main verb or other components may be paraphrased with synonyms, provided that the sentence's tense, structure, and meaning remain strictly equivalent before applying standard negation. Authors refer to the Merriam-Webster Thesaurus [6]. For example,

---

[6] https://www.merriam-webster.com/

16

- **Original Sentence**: Toumour is a village and rural commune in Niger **located near** the Niger–Nigeria **border**.

  - **Paraphrased Sentence**: Toumour is a village and rural commune in Niger **that is found close to** the Niger–Nigeria **boundary**.
    → **Standard Negation after Paraphrase**: Toumour **isn't** a village and rural commune in Niger that is found close to the Niger–Nigeria boundary.
  - **Explanation**: In this example, the participle clause "located near the Niger–Nigeria border" is rephrased as a relative clause "that is found close to the Niger–Nigeria boundary." Since both constructions serve as modifiers and preserve the same semantic role, we treat them as equivalent in meaning for the purpose of standard negation.

- **Original Sentence**: The armed forces **said** Boko Haram **attacked** their military post on March 15, 2020, which they responded to by repelling the attack, killing 50 insurgents.

  - **Paraphrased Sentence**: The armed forces **stated** that Boko Haram **assaulted** their military post on March 15, 2020, which they responded to by repelling the attack, killing 50 insurgents.
    → **Standard Negation after Paraphrase**: The armed forces **didn't state** that Boko Haram assaulted their military post on March 15, 2020, which they responded to by repelling the attack, killing 50 insurgents.
  - **Explanation**: In this example, the reporting verb "said" is paraphrased as "stated," and the verb "attacked" is replaced with the synonym "assaulted." These substitutions preserve the original tense and meaning, allowing standard negation to be applied without altering the semantic content of the sentence.

**Negation of Simple Sentences.** For simple, declarative sentences, standard negation is achieved by inserting "not" after the auxiliary or main verb, or by replacing the predicate with its complementary antonym. For example, "She is happy." → "She is not happy."; "The room is occupied." → "The room is unoccupied."

**Negation in Compound Sentences.** When multiple clauses or propositions are coordinated (e.g., with "and", "or", "but"), standard negation is logically applied, governed by De Morgan's laws:

- For statements of the form "A and B", the negation is "¬A or ¬B".

- For statements of the form "A or B", the negation is "¬A and ¬B".

For example, "He passed the test and received an award." → "He did not pass the test or did not receive an award."

For compound structures involving three or more elements, the same logic applies:

- A, B, and C = (A, and B), and C
  → ¬(A, and B), or ¬C
  = (¬A, or ¬B), or ¬C
  = ¬A, ¬B, or ¬C

- A, B, or C = (A, or B), or C
  → ¬(A, or B), and ¬C
  = (¬A, and ¬B), and ¬C
  = ¬A, ¬B, and ¬C

- A and B or C = (A, and B), or C
  → ¬(A, and B), and ¬C
  = (¬A, or ¬B), and ¬C
  = ¬A, or ¬B, and ¬C

When strict application of logical negation produces unnatural language, sentences may be split or slightly rephrased for fluency, provided logical meaning is preserved. For example,

- **Original:** "He finished the report and submitted the assignment."

- **Standard Negation:** "He did not finish the report or did not submit the assignment."

- **Standard Negation, but Splitted:** "He did not finish the report. Or, he did not submit the assignment."

**Coordinated Elements in the Sentence.** When a sentence contains coordinated elements (such as subjects, objects, or predicates connected by "and" or "or"), standard negation typically follows logical principles derived from De Morgan's Laws. However, whether logical negation applies to each individual component or to the entire predicate as a whole depends on whether the coordination expresses multiple independent propositions or a single collective event.

- If the coordination introduces semantically distinct propositions, that is, each conjunct could form a complete sentence on its own, negation must be applied to each proposition individually. For example,
"My sister and I studied hard."
This sentence can be interpreted as: "My sister studied hard and I studied hard."
Therefore, the correct standard negation is: "My sister did not study hard, or I did not study hard."

- Conversely, if the coordination connects elements that jointly participate in a single action or state (e.g., a shared subject or a collective predicate), then the sentence is treated as a simple clause, and the predicate as a whole is negated. Logical decomposition is not appropriate. For example,
"My sister and I share clothes."
This expresses a single collective action involving both participants.
Therefore, the correct standard negation is: "My sister and I do not share clothes."
(NOT: "My sister does not share clothes, or I do not share clothes.")

- This distinction is crucial: even if two noun phrases are coordinated, if the sentence semantically decomposes into separate propositional content, standard negation must apply to each sub-proposition. Otherwise, it applies to the whole predicate as one unit.

- Other examples of semantically collective predicates where logical splitting is not appropriate include: "be the same", "have in common", "do something together", "combine", "unite", etc. These describe inherently joint or relational properties, not independent propositions. For example,
"Clarence Brown and Peter Glenville are from the same country." $\rightarrow$ "Clarence Brown and Peter Glenville are not from the same country."

**Use of Antonyms.** When replacing predicates with antonyms in standard negation, only complementary antonyms are appropriate, as they provide a clear binary opposition, ensuring logical consistency of negation. Gradable and relational antonyms are unsuitable for standard negation because their antonyms do not represent the logical

complement of the original predicate. In other words, replacing a predicate $P$ with its antonym does not produce $\neg P$ in a truth-conditional sense.

Specifically, unlike complementary antonyms, which form mutually exclusive pairs (i.e., $P \cup \neg P = U$ and $P \cap \neg P = \emptyset$), gradable and relational antonyms do not partition the meaning space cleanly, and thus fail to reverse the truth value reliably.

- **Complementary Antonyms**: Also called binary/contradictory antonyms. These antonyms represent mutually exclusive pairs with no intermediate states. The presence of one implies the absence of the other. Examples include:
  - alive / dead
  - true / false
  - present / absent
  - occupied / vacant

  Using complementary antonyms in negation ensures a direct and unambiguous reversal of the original proposition's truth value.

- **Gradable Antonyms**: These antonyms exist on a continuum and allow for varying degrees between the two extremes. Negating one does not necessarily affirm the other. Examples include:
  - hot / cold
  - happy / sad
  - tall / short
  - young / old

  Due to their scalar nature, gradable antonyms are inappropriate for standard negation, as they do not provide a definitive binary opposition.

- **Relational Antonyms**: Also known as converse antonyms, these pairs describe a reciprocal relationship where one implies the existence of the other. Examples include:
  - parent / child
  - teacher / student
  - buy / sell
  - employer / employee

  Relational antonyms are context-dependent and do not represent direct opposites in a binary sense, making them unsuitable for standard negation purposes.

18

**General Principles of Standard Negation.**

- The negated sentence must preserve all elements (subject, tense, objects, adjuncts, etc.) of the original, except for the truth value of the main predicate.

- When naturalness and logical negation conflict, logical correctness takes priority, but minimal rephrasing is allowed for fluency.

- If the negated clause creates a contradiction with other parts of the sentence, the contradictory clause must be removed. For example, "While the spatial size of the entire universe is unknown, it is possible to measure the size of the observable universe, which is approximately 93 billion light-years in diameter."
  → "While the spatial size of the entire universe is unknown, it isn't possible to measure the size of the observable universe."
  (The relative clause must be removed due to semantic contradiction inside one sentence.)

**Common Negation Errors and Corrections.**

- **Original sentence**: His characteristic style fuses samba, funk, rock **and** bossa nova with lyrics that blend humor and satire with often esoteric subject matter.

  - **Incorrect negation**: His distinctive style **doesn't fuse** samba, funk, rock **or** bossa nova with lyrics that blend humor and satire with often esoteric subject matter.
  - **Correct negation**: His distinctive style **doesn't fuse** samba, funk, rock **and** bossa nova with lyrics that blend humor and satire with often esoteric subject matter.
  - **Explanation**: The verb "fuse" implies a combination of all listed elements. "and" must be preserved.

- **Original sentence**: The mascot of Avon Center School is the "Koalaty Kid," **while** the mascot at Prairieview **is** an eagle **and** the mascot at Woodview **is** an owl.

  - **Incorrect negation**: Avon Center School's mascot is not the "Koalaty Kid," Prairieview's mascot **is not** an eagle, **or** Woodview's mascot **is not** an owl.

  - **Correct negation**: Avon Center School's mascot is not the "Koalaty Kid," **while** the mascot at Prairieview **is** an eagle **and** the mascot at Woodview **is** an owl.
  - **Explanation**: Two clauses connected by "while" are not coordinated propositions (they are not equally connected), but instead present contrastive information. Therefore, it is incorrect to apply logical negation across both clauses. Only the relevant clause (e.g., the first statement) should be negated.

## F  Compound Sentences and Coordinating Conjunction

A compound sentence consists of two or more independent clauses joined by a coordinating conjunction. Each clause can stand alone, but they are combined to express related ideas (Gleitman, 1965).

Coordinating conjunctions connect elements of equal grammatical rank. The seven common ones in English are: *for, and, nor, but, or, yet, so* (often remembered as FANBOYS). Among these, *and*, *or*, and *but* are indisputably used to coordinate clauses. The others can be ambiguous or function in non-coordinating roles(e.g., indicating cause or result rather than logical structure). These are the examples using *and*, *or*, and *but* to connect sentences equally.

- "She studied hard, **and** she passed the exam."

- "I wanted to go, **but** it was raining."

- "You can call me, **or** you can send an email."

Only clearly coordinating conjunctions (*and, or, but*) are considered for clause-level negation in this benchmark, indicating that two or more independent clauses are equally connected.

## G  Negation Phenomena Considered but Omitted

### G.1  Double Negation

Double negation refers to the use of two forms of grammatical negation within a single sentence. In standard English, only one negative form should be present in a subject-predicate construction; the presence of two negatives is generally considered non-standard and often results in an unintended meaning. For example, while "He's going nowhere" is

correct, "He's not going nowhere" is ungrammatical. Another example is "I won't bake no cake," which combines verb negation ("won't") with object negation ("no cake"), resulting in a grammatically incorrect construction (Déprez et al., 2015).

In English, certain double negation constructions convey affirmative meanings rather than intensifying negation, effectively paraphrasing the original positive statement (e.g., $\neg\neg p \approx p$) (Van der Wouden, 1996). This rhetorical device, known as litotes, often manifests in expressions such as "not bad," implying "good," or "not unhappy," implying "happy." Leveraging this phenomenon, we generated paraphrase candidates for our dataset using such double negation patterns. For example,

- **Sentence:** His characteristic style fuses samba, funk, rock and bossa nova with lyrics that blend humor and satire with often esoteric subject matter.
  **Double Negation:** His characteristic style **does not fail to fuse** samba, funk, rock, and bossa nova with lyrics that blend humor and satire with often esoteric topics.

- **Sentence:** It covers a broad range of fields, including the humanities, social sciences, exact sciences, applied sciences, and life sciences.
  **Double Negation:** It **does not exclude** a broad range of fields, including the humanities, social sciences, exact sciences, applied sciences, and life sciences.

- **Sentence:** Sanders was honoured to meet with many world dignitaries and representatives of UNESCO member nations, and delighted when delegates from UNESCO, visited Toowoomba in 2018 in return.
  **Double Negation:** Sanders **was not unhappy** to meet with many world dignitaries and representatives of UNESCO member nations, and not displeased when delegates from UNESCO visited Toowoomba in 2018 in return.

However, upon closer examination, these paraphrase candidates do not always preserve the exact meaning of the original sentence. The antonyms used (e.g., "exclude" for "cover," "unhappy" for "honoured") are not always true complementary antonyms, which does not effectively negate the meaning. Moreover, the litotes construction ("does not fail to fuse") tends to add an emphatic nuance, rather than being a perfect semantic equivalent.

Therefore, the boundary between paraphrasing and double negation is ambiguous, and their relationship requires more careful analysis. Given these issues, and because our primary focus is on standard, direct forms of negation, we ultimately decide to exclude double negation constructions as paraphrase candidates from our dataset.

## G.2 Conditional Statement

Negating conditional statements presents additional challenges, as the intuitive approach to negation often diverges from the logical form. For example, the negation of "If I study hard, I will pass the bar exam" is frequently (but incorrectly) interpreted as "If I do not study hard, I will not pass the bar exam."

Letting "I study hard" be A and "I will pass the bar exam" be B, the original conditional "If I study hard, I will pass the bar exam" (A ⇒ B) is logically equivalent to $\neg A \lor B$. The incorrect logical negation, "If I do not study hard, I will not pass the bar exam" (¬A ⇒ ¬B), corresponds to $A \lor \neg B$. The correct negation should be $\neg A \lor B \to A \land \neg B$, which means "I study hard and I will not pass the bar exam" (Nguyen et al., 2023).

Such distinctions reveal that logical negation and natural language intuition often do not align. Furthermore, forming proper paraphrases for the negation of conditionals is not straightforward and frequently requires extensive rewriting rather than a simple insertion of "not" or converting the verb to its complementary antonym. Since our task is centered on the correct application of verbal negation in main predicates, rather than the broader domain of logical negation in complex constructions, we choose to exclude conditional statements from our dataset.

## G.3 Local Negation Constructions Excluded

In constructing the Thunder-NUBench dataset, we consider various types of local (i.e., subclausal) negation, where negation applies to a phrase or constituent rather than the main predicate. However, several constructions are excluded due to their semantic ambiguity, syntactic irregularity, or misalignment with the benchmark's focus on predicate-level negation.

**Infinitive Phrase Negation.** Infinitive phrases (e.g., "to go") can be negated with "not" (e.g., "not to go" or "to not go"). Although grammatically correct, this construction is relatively rare and sounds

awkward depending on the context.

- **Original**: George wants to go to the park.

- **Negated (infinitive)**: George wants not to go / George wants to not go.

Since our benchmark targets the negation of predicates, both when generating standard negation and local negation, we exclude this type of infinitive phrase negation, which typically affects subordinate or embedded phrases.

**Appositive Clause Negation.** Appositive clauses are noun phrases that provide descriptive clarification. Attempting to negate an appositive typically involves lexical replacement rather than syntactic negation.

- **Original:** My brother, a talented musician, plays the guitar.

- **Negated (appositive):** My brother, not a talented musician, plays the guitar.

Such changes alter descriptive content rather than reversing the meaning of the predicate, and often fall into the domain of contradiction. Accordingly, they were excluded from the dataset.

**Prepositional Phrase Negation.** Negating a prepositional phrase often involves replacing the preposition with its antonym (e.g., "with" → "without", "in" → "outside"), which results in a sentence that differs in content, rather than reversing the meaning of the predicate.

- **Original**: She went to the park with her bird.

- **Negated (preposition)**: She went to the park without her bird.

Since such modifications do not negate the verb but instead change the nature of an adjunct or argument, they fall outside the scope of standard negation or local negation in this work and are excluded.

In all of the above cases, the negation does not target the verb but rather peripheral elements within the sentence. As the Thunder-NUBench is designed to evaluate verbal negation, these local or phrase-level forms of negation were intentionally left out.

## H  Thunder-NUBench Structure

Thunder-NUBench consists of two subsets: a sentence-negation pair dataset for supervised fine-tuning and a multiple-choice dataset for evaluation. Both datasets are built on English text and reviewed by authors following strict guidelines.

**Sentence-Negation Pair Dataset.** This subset contains pairs of affirmative and corresponding standard negation sentences. It includes the following fields:
- `index`: the index of the data.
- `premise`: the original sentence.
- `hypothesis`: its logically negated form.

**Multiple-Choice Dataset.** This evaluation set presents each original sentence with four candidate transformations.
- `wikipedia_index`: the original index of the Wikipedia Summary dataset.
- `index`: the index of the data.
- `sentence`: the original sentence.
- `choice1`: standard negation (correct answer).
- `choice2`: local negation (subclausal negation).
- `choice2_type`: specifies the type of local negation.
- `choice2_element`: a short description of the phrase or clause that was negated (e.g., "being built", "which crashed").
- `choice3`: contradiction (non-negated, semantically incompatible).
- `choice4`: paraphrase (semantically equivalent).

**Local Negation Types and Distributions.** The details of `choice2` types and distribution on validation and test sets are described in Table 9.

## I  Models We Use in Experiments

We evaluate two primary model size groups: 3–4 billion (B) parameters and 7–8 billion (B) parameters. Each group contains both pretrained models and instruction-tuned models, as detailed below (Qwen et al., 2025; Grattafiori et al., 2024; Team et al., 2025; Jiang et al., 2023; Team et al., 2024):
- **3-4B models:**
  - Pretrained: Llama-3.2-3B, Qwen2.5-3B, gemma-3-4b-pt

21

| choice2_type | definition | Validation Set | Test Set |
|---|---|---|---|
| relative_part | negation inside relative clauses (e.g., "who did not attend. . . "). | 31 | 293 |
| pp_part | negation in participle clauses (e.g., "not walking through the park. . . "). | 30 | 290 |
| compound_part | negation applied to one clause within a compound sentence. | 29 | 277 |
| adverb_part | negation in adverbial clauses (e.g., "because it was not raining"). | 8 | 105 |
| non-applicable | used when the sentence structure does not support a valid local negation variant under our definition. | 2 | 37 |
| **Total** | | 100 | 1,002 |

Table 9: Choice 2 Types and Distributions.

  - Instruction-tuned: Llama-3.2-3B-Instruct, Qwen2.5-3B-Instruct, gemma-3-4b-it
- **7-8B models:**
  - Pretrained: Mistral-7B-v0.3, Llama-3.1-8B, Qwen2.5-7B, gemma-7b
  - Instruction-tuned: Mistral-7B-Instruct-v0.3, Llama-3.1-8B-Instruct, Qwen2.5-7B-Instruct, gemma-7b-it

## J  Prompt Selection for In-Context Learning

We explore a range of prompt types, from basic to structured and reasoning-based instructions (Zhao et al., 2021; Li, 2023; Wan et al., 2023). Table 10 shows the details of all prompt templates explored. Among the prompt variants, prompt 1 (default), prompt 7 (instruction-based), and prompt 11 (reasoning-based) were selected for use in the experiments. Only the best-performing prompt on the validation set is reported in the main text (Section 5).

### J.1  Prompt Type Design

Analysis of prompt type validation results reveals that prompt 7, which provides an instruction-based, structured formulation, achieves the highest overall performance. Based on these results, we adopt this prompt as the base template for our main experiments. Additionally, among the default-type prompts, prompt 1 shows the best average performance for 7–8B models, while among reasoning-based prompts, prompt 11 shows better performance among all models. See the details in Table 11 and Table 12. Therefore, we use these prompts as the default prompt (prompt 1), instruction prompt (prompt 7), and reasoning

prompt (prompt 11) in the experiment results of Appendix M.

The optimal prompt type or format can differ depending on model series and size, highlighting that the choice of prompt significantly impacts model performance. Nevertheless, detailed instructions or reasoning-based formats generally outperform minimal "default" prompts, underscoring the importance of providing explicit context and guidance in in-context learning settings.

### J.2  Prompt Wording Variants

We further investigate the effect of prompt wording by testing minor variants in the phrasing of instructions and output cues, while maintaining the same underlying task specification. For example, all prompts in this set share the same core instruction,

*"Logically negate the sentence below. If the sentence includes 'A and B', use 'not A or not B'. If it includes 'A or B', use 'not A and not B'. Also apply 'not' or use complementary antonyms on the main verb(s) of the entire sentence."*
of the prompt 7, but differ in the wording of the instruction or the label for the output (e.g., "Negation:", "Negated:", "Output:").

See Table 13 for a full list of variants. Even such superficial changes in wording can affect performance, further emphasizing the importance of careful prompt engineering.

In our main experiments, we use the standard *"Sentence: {doc['sentence']}\nNegation:"* format as the output cue for consistency.

## K  Evaluating via Language Model Evaluation Harness

This section describes how Thunder-NUBench and ScoNe is integrated into the LM Evaluation Har-

| prompt index | prompt Type | Detail |
|---|---|---|
| **prompt 1** | Default | Negate the following sentence using standard logical negation:<br>{doc['sentence']}<br>Negated: |
| prompt 2 | Default | Logically negate the following statement:<br>{doc['sentence']}<br>Negation:" |
| prompt 3 | Default | Reverse the truth value of the statement.<br>Sentence: {doc['sentence']}<br>Answer: |
| prompt 4 | Default | Reverse the entire meaning of the following sentence using logical and standard negation.<br>Sentence: {doc['sentence']}<br>Negated: |
| prompt 5 | Definition-based + Default | Standard negation involves inserting 'not' or using a complementary antonym. Logically negate:<br>{doc['sentence']}<br>Negation: |
| prompt 6 | Structured + Definition-based | Standard negation involves negating a verbal declarative sentence either by inserting 'not' after the primary or auxiliary verb, or by replacing a predicate with its complementary antonym of the main clause.<br>Logically apply the correct standard negation for the following sentence: {doc['sentence']}<br>Standard negation: |
| **prompt 7** | structured + instruction-based | Logically negate the sentence below. If the sentence includes 'A and B', use 'not A or not B'. If it includes 'A or B', use 'not A and not B'. Also apply 'not' or use complementary antonyms on the main verb(s) of the entire sentence.<br>Sentence: {doc['sentence']}<br>Negation:" |
| prompt 8 | structured + instruction-based | Apply standard logical negation to the following sentence. This includes:<br>1. Syntactic negation using 'not' or auxiliaries on the main verb.<br>2. Logical negation (e.g., using De Morgan's laws: 'A and B' becomes 'not A or not B').<br>3. Lexical negation using complementary antonyms on the main predicate (e.g., 'alive' → 'dead').<br>Sentence: {doc['sentence']}<br>Negated version:" |
| prompt 9 | structured + definition, (comprehensive) instruction-based | Standard negation refers to the basic structural way of negating declarative verbal main clauses in English.<br>Apply standard logical negation to reverse the truth value of the following sentence.<br>This includes negating the main predicate by applying:<br>1. Syntactic negation using 'not' or auxiliaries.<br>2. Logical negation (e.g., using De Morgan's laws: 'A and B' becomes 'not A or not B').<br>3. Lexical negation using complementary antonyms (e.g., 'alive' → 'dead').<br>Sentence: {doc['sentence']}<br>Negated: |
| prompt 10 | role-based + default | You are a logic expert. Provide the logically correct standard negation of the following sentence: '{doc['sentence']}' |
| **prompt 11** | Reasoning-based | Think step-by-step to logically negate the following sentence. Use standard negation: syntactic negation ('not') or complementary antonym on main predicate, and logical rules (e.g., De Morgan's law).<br>Reverse the truth value of the statement.<br>Sentence: {doc['sentence']}<br>Negated version:" |
| prompt 12 | Reasoning-based | Let's negate this sentence step by step:<br>1. Identify the core propositions.<br>2. Apply logical negation to each part.<br>3. Construct the final negated sentence.<br>Sentence: {doc['sentence']}<br>Negation: |

Table 10: Overview of the prompt variants used for in-context learning in Thunder-NUBench experiments. Prompt 7 is the base prompt of the whole paper. Prompts 1, 7, and 11 are used in experiments (each as default, instruction-based, and reasoning-based prompt).

| | | Llama | | | | gemma | | | | Qwen | | | | Average | |
| | | Llama-3.2-3B | | Llama-3.2-3B-Instruct | | gemma-3-4b-pt | | gemma-3-4b-it | | Qwen2.5-3B | | Qwen2.5-3B-Instruct | | | |
| prompt index | prompt type | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **prompt 1** | | 0.550 | 0.550 | 0.590 | 0.600 | 0.370 | 0.410 | 0.400 | 0.430 | 0.580 | 0.590 | 0.590 | 0.640 | **0.513** | **0.537** |
| **prompt 2** | | 0.520 | 0.510 | 0.600 | 0.620 | 0.360 | 0.420 | 0.400 | 0.400 | 0.540 | 0.570 | 0.600 | 0.620 | **0.503** | **0.523** |
| **prompt 3** | default | 0.430 | 0.440 | 0.520 | 0.520 | 0.350 | 0.380 | 0.380 | 0.410 | 0.460 | 0.490 | 0.490 | 0.510 | **0.438** | **0.458** |
| **prompt 4** | | 0.500 | 0.550 | 0.590 | 0.620 | 0.440 | 0.470 | 0.440 | 0.460 | 0.530 | 0.570 | 0.590 | 0.640 | **0.515** | **0.552** |
| **prompt 5** | | 0.480 | 0.510 | 0.470 | 0.490 | 0.380 | 0.400 | 0.420 | 0.460 | 0.510 | 0.520 | 0.560 | 0.580 | **0.470** | **0.493** |
| **prompt 6** | | 0.540 | 0.560 | 0.620 | 0.650 | 0.440 | 0.480 | 0.430 | 0.460 | 0.520 | 0.550 | 0.590 | 0.610 | **0.523** | **0.552** |
| **prompt 7** | instruction | 0.530 | 0.540 | 0.590 | 0.620 | 0.470 | 0.510 | 0.460 | 0.460 | 0.590 | 0.610 | 0.770 | 0.770 | **0.568** | **0.585** |
| **prompt 8** | | 0.510 | 0.510 | 0.580 | 0.600 | 0.510 | 0.520 | 0.500 | 0.520 | 0.580 | 0.600 | 0.610 | 0.660 | **0.548** | **0.568** |
| **prompt 9** | | 0.530 | 0.530 | 0.580 | 0.610 | 0.500 | 0.530 | 0.540 | 0.550 | 0.580 | 0.570 | 0.600 | 0.630 | **0.555** | **0.570** |
| **prompt 10** | role | 0.450 | 0.480 | 0.410 | 0.440 | 0.360 | 0.420 | 0.450 | 0.510 | 0.550 | 0.560 | 0.620 | 0.660 | **0.473** | **0.512** |
| **prompt 11** | reasoning | 0.570 | 0.590 | 0.620 | 0.640 | 0.440 | 0.450 | 0.530 | 0.540 | 0.580 | 0.590 | 0.610 | 0.640 | **0.558** | **0.575** |
| **prompt 12** | | 0.550 | 0.550 | 0.630 | 0.670 | 0.450 | 0.460 | 0.510 | 0.510 | 0.600 | 0.600 | 0.600 | 0.630 | **0.557** | **0.570** |

Table 11: Prompt-wise Results on the Validation Set for 3B–4B Models. The red values indicate the highest performance in each column. *acc* denotes the accuracy, and *acc_norm* indicates the length-normalized accuracy.

Table 12:

| | | Mistral | | | | Llama | | | | gemma | | | | Qwen | | | | average | |
| | | Mistral-7B-v0.3 | | Mistral-7B-Instruct-v0.3 | | Llama-3.1-8B | | Llama-3.1-8B-Instruct | | gemma-7b | | gemma-7b-it | | Qwen2.5-7B | | Qwen/Qwen2.5-7B-Instruct | | | |
| prompt index | prompt type | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm | acc | acc_norm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| prompt 1 | | 0.540 | 0.540 | 0.620 | 0.660 | 0.490 | 0.520 | 0.610 | 0.640 | 0.510 | 0.540 | 0.630 | 0.690 | 0.630 | 0.650 | 0.680 | 0.700 | **0.589** | **0.618** |
| prompt 2 | | 0.480 | 0.500 | 0.570 | 0.630 | 0.510 | 0.520 | 0.560 | 0.590 | 0.520 | 0.560 | 0.620 | 0.630 | 0.630 | 0.660 | 0.670 | 0.670 | **0.570** | **0.595** |
| prompt 3 | default | 0.420 | 0.460 | 0.520 | 0.560 | 0.480 | 0.520 | 0.500 | 0.510 | 0.500 | 0.520 | 0.540 | 0.570 | 0.470 | 0.510 | 0.540 | 0.540 | **0.496** | **0.524** |
| prompt 4 | | 0.430 | 0.470 | 0.610 | 0.650 | 0.450 | 0.460 | 0.520 | 0.570 | 0.500 | 0.530 | 0.540 | 0.570 | 0.540 | 0.550 | 0.650 | 0.670 | **0.530** | **0.559** |
| prompt 5 | | 0.490 | 0.520 | 0.710 | 0.710 | 0.460 | 0.470 | 0.450 | 0.470 | 0.500 | 0.540 | 0.560 | 0.580 | 0.530 | 0.550 | 0.560 | 0.600 | **0.533** | **0.555** |
| prompt 6 | | 0.520 | 0.550 | 0.690 | 0.730 | 0.490 | 0.510 | 0.560 | 0.590 | 0.520 | 0.530 | 0.520 | 0.570 | 0.540 | 0.570 | 0.520 | 0.570 | **0.545** | **0.578** |
| prompt 7 | instruction | 0.540 | 0.570 | 0.740 | 0.750 | 0.500 | 0.520 | 0.670 | 0.670 | 0.570 | 0.570 | 0.640 | 0.640 | 0.600 | 0.620 | 0.670 | 0.710 | **0.616** | **0.631** |
| prompt 8 | | 0.530 | 0.540 | 0.700 | 0.730 | 0.470 | 0.500 | 0.570 | 0.600 | 0.540 | 0.540 | 0.720 | 0.720 | 0.640 | 0.660 | 0.660 | 0.690 | **0.604** | **0.623** |
| prompt 9 | | 0.570 | 0.580 | 0.730 | 0.730 | 0.460 | 0.490 | 0.560 | 0.600 | 0.540 | 0.550 | 0.610 | 0.620 | 0.600 | 0.600 | 0.650 | 0.680 | **0.590** | **0.606** |
| prompt 10 | role | 0.500 | 0.520 | 0.630 | 0.650 | 0.460 | 0.480 | 0.470 | 0.540 | 0.500 | 0.520 | 0.660 | 0.750 | 0.610 | 0.630 | 0.620 | 0.640 | **0.556** | **0.591** |
| prompt 11 | reasoning | 0.520 | 0.550 | 0.650 | 0.660 | 0.510 | 0.530 | 0.550 | 0.610 | 0.580 | 0.610 | 0.710 | 0.740 | 0.670 | 0.690 | 0.680 | 0.690 | **0.609** | **0.635** |
| prompt 12 | | 0.540 | 0.540 | 0.680 | 0.720 | 0.530 | 0.550 | 0.610 | 0.610 | 0.580 | 0.570 | 0.630 | 0.650 | 0.590 | 0.620 | 0.590 | 0.620 | **0.594** | **0.610** |

Table 12: Prompt-wise Results on the Validation Set for 7B–8B Models. The red values indicate the highest performance in each column. *acc* denotes the accuracy, and *acc_norm* indicates the length-normalized accuracy.

| prompt index | Detail | 3-4B models average | | 7-8B models average | |
| | | acc | acc_norm | acc | acc_norm |
|---|---|---|---|---|---|
| prompt 7 | {instruction}\nSentence: {doc['sentence']}\nNegation:" | 0.568 | 0.585 | 0.616 | 0.631 |
| prompt 7-1 | {instruction}\nSentence: {doc['sentence']}\nAnswer:" | 0.536 | 0.565 | 0.614 | 0.631 |
| prompt 7-2 | {instruction}\nSentence: {doc['sentence']}\nNegated:" | 0.560 | 0.578 | 0.611 | 0.631 |
| prompt 7-3 | {instruction}\nSentence: {doc['sentence']}\nNegated version:" | 0.556 | 0.573 | 0.638 | 0.646 |
| prompt 7-4 | {instruction}\nOriginal: {doc['sentence']}\nRewritten with negation:" | 0.533 | 0.542 | 0.633 | 0.646 |
| prompt 7-5 | {instruction}\nInput: {doc['sentence']}\nOutput:" | 0.528 | 0.550 | 0.616 | 0.633 |

Table 13: Prompt Wording Variants Used for Validation Set Evaluation. The base prompt is prompt 7 from Table 10. The red values indicate the highest performance in each column. *acc* denotes the accuracy, and *acc_norm* indicates the length-normalized accuracy.

ness (Gao et al., 2024) for zero-shot, few-shot evaluation. We use the prompt for ScoNe according to the original paper (Hypothesis Question Prompt format) (She et al., 2023).

```yaml
task: nubench
dataset_path: {dataset_path}/NUBench
dataset_name: null
output_type: multiple_choice
validation_split: validation
test_split: test
fewshot_split: validation
process_docs: !function utils.process_docs
doc_to_text: "{{query}}"
doc_to_target: "{{gold}}"
doc_to_choice: "choices"
metric_list:
  - metric: acc
    aggregation: mean
    higher_is_better: true
  - metric: acc_norm
    aggregation: mean
    higher_is_better: true
```

Listing 7: NUBench/NUBench.yaml

```python
        choices = [doc["choice1"], doc["choice2"], doc["choice3"], doc["choice4"]]

    return {
        "query": prompt,
        "choices": choices,
        "gold": 0  # always choice1
    }
    return dataset.map(_process_doc)
```

Listing 8: NUBench/utils.py

```yaml
task: ScoNe
dataset_path: tasksource/ScoNe
dataset_name: null
output_type: multiple_choice
training_split: train
test_split: test
process_docs: !function utils.process_docs
doc_to_text: "{{query}}"
doc_to_target: "{{gold}}"
doc_to_choice: "choices"
metric_list:
  - metric: acc
```

Listing 9: ScoNe/ScoNe.yaml

```python
import re
import datasets

def process_docs(dataset: datasets.Dataset) ->
    datasets.Dataset:
    def _process_doc(doc):
        prompt = f"Logically negate the sentence
    below. If the sentence includes 'A and B', use
    'not A or not B'. If it includes 'A or B', use
    'not A and not B'. Also apply 'not' or use
    complementary antonyms on the main verb(s) of
    the entire sentence.\nSentence: {doc['sentence
    ']}\nNegation:"

        if doc.get("choice2_type", "") == "non-
    applicable":
            choices = [doc["choice1"], doc["choice3"
    ], doc["choice4"]]
        else:
```

```python
import re
import datasets

def process_docs(dataset: datasets.Dataset) ->
    datasets.Dataset:
    def _process_doc(doc):
        prompt = f"Assume that {doc['
    sentence1_edited']}\nIs it then definitely true
    that {doc['sentence2_edited']}? Answer yes or
    no.\nAnswer:"

        choices = ["Yes", "No"]
        if doc['gold_label_edited'] == "entailment":
            gold = 0
        else:
            gold = 1
        return {
            "query": prompt,
            "choices": choices,
```

24

```
16              "gold": gold
17          }
18      return dataset.map(_process_doc)
```

Listing 10: ScoNe/utils.py

## L  Finetuning via LLaMA-Factory

We detail our supervised fine-tuning setup using LLaMA-Factory (Zheng et al., 2024) with LoRA (Hu et al., 2022) on Thunder-NUBench training data, including configuration of the fine-tuning and instruction-based examples in Alpaca format (Taori et al., 2023).

The YAML configuration provided in Listing 11 is specific to the LLaMA-3.1-8B model. Other models (e.g., Qwen or Mistral) can be fine-tuned similarly by modifying the `model_name_or_path` and `template` fields in the configuration file accordingly.

```
1  ### model
2  model_name_or_path: Llama-3.1-8B
3  trust_remote_code: true
4
5  ### method
6  stage: sft
7  do_train: true
8  finetuning_type: lora
9  lora_rank: 8
10 lora_target: all
11
12 ### dataset
13 dataset: nubench_train
14 template: llama3
15 cutoff_len: 512
16 max_samples: 5000
17 overwrite_cache: true
18 preprocessing_num_workers: 16
19 dataloader_num_workers: 4
20
21 ### output
22 output_dir: lora/sft/Llama-3.1-8B
23 logging_steps: 10
24 save_steps: 500
25 plot_loss: true
26 overwrite_output_dir: true
27 save_only_model: false
28 report_to: none  # choices: [none, wandb,
        tensorboard, swanlab, mlflow]
29
30 ### train
31 per_device_train_batch_size: 1
32 gradient_accumulation_steps: 8
33 learning_rate: 1.0e-4
34 num_train_epochs: 3.0
35 lr_scheduler_type: cosine
36 warmup_ratio: 0.1
37 bf16: true
38 ddp_timeout: 180000000
39 resume_from_checkpoint: null
```

Listing 11: Llama-3.1-8B_lora_sft.yaml

```
1  [
2    {
3      "instruction": "Logically negate the sentence
        below. If the sentence includes 'A and B', use
        'not A or not B'. If it includes 'A or B', use
        'not A and not B'. Also apply 'not' or use
        complementary antonyms on the main verb(s) of
        the entire sentence.",
4      "input": "Sentence: Eddie Vedder was born before
        Nam Woo-hyun.\nNegation:",
5      "output": "Eddie Vedder wasn't born before Nam
        Woo-hyun."
```

```
6    },
7    {
8      "instruction": "Logically negate the sentence
        below. If the sentence includes 'A and B',
        use 'not A or not B'. If it includes 'A or B',
        use 'not A and not B'. Also apply 'not' or use
        complementary antonyms on the main verb(s) of
        the entire sentence.",
9      "input": "Sentence: Halestorm is from
        Pennsylvania, while Say Anything is from
        California.\nNegation:",
10     "output": "Halestorm is not from Pennsylvania,
        while Say Anything is from California."
11   },
12   (...)
13 ]
```

Listing 12: Sentence-Negation Pair dataset for training in alpaca format.

We use 4 NVIDIA RTX 3090 GPUs (24GB) with CUDA 12.4 for all model training and evaluation. Both supervised fine-tuning and evaluation are completed in under an hour using these resources.

## M  Total Model Performance on Thunder-NUBench

This section reports zero-shot, 5-shot performance (averaged over five random seeds: 42, 1234, 3000, 5000, and 7000) of various pre-trained language models on Thunder-NUBench, prior to any fine-tuning. The default prompt is prompt 1, the instruction prompt is prompt 7, and the reasoning prompt is prompt 11 from Table 10. Table 14 shows the overall results.

Instruction-tuned models consistently outperform their pretrained counterparts across all model series. Although detailed instruction prompts provide explicit task descriptions, the substantial performance improvement observed with few-shot compared to zero-shot settings suggests that concrete examples are even more effective in helping models understand and perform the task. This trend holds even under instruction-based prompting, indicating that examples play a more critical role than instructions alone in guiding model behavior on logical negation tasks. We can also observe that performance varies substantially depending on the prompt type. For example, Mistral-7B-Instruct-v0.3 achieves an accuracy of 0.566 with the default prompt in a zero-shot setting, but this rises to 0.667 with the instruction prompt and 0.606 with the reasoning prompt. Notably, both results on instruction and reasoning prompts significantly outperform the ones on default prompt. This pattern indicates that each model has different preferences for prompt type, and selecting the optimal prompt is crucial for maximizing performance.

| | | | | default prompt | | instruction prompt | | reasoning prompt | | average | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | zeroshot | fewshot (±SD) | zeroshot | fewshot (±SD) | zeroshot | fewshot (±SD) | zeroshot | fewshot (±SD) |
| **3B-4B** | **Llama** | **Llama-3.2-3B** | acc | 0.475 | 0.603 (±0.006) | 0.463 | 0.616 (±0.005) | 0.518 | 0.634 (±0.006) | **0.485** | **0.618** |
| | | | acc_norm | 0.504 | 0.623 (±0.006) | 0.483 | 0.633 (±0.004) | 0.539 | 0.653 (±0.006) | **0.509** | **0.636** |
| | | **Llama-3.2-3B-Instruct** | acc | 0.496 | 0.640 (±0.009) | 0.490 | 0.619 (±0.006) | 0.573 | 0.612 (±0.009) | **0.520** | **0.624** |
| | | | acc_norm | 0.517 | 0.653 (±0.005) | 0.512 | 0.633 (±0.001) | 0.592 | 0.630 (±0.008) | **0.540** | **0.639** |
| | **gemma** | **gemma-3-4b-pt** | acc | 0.369 | 0.595 (±0.018) | 0.4381 | 0.614 (±0.007) | 0.498 | 0.640 (±0.006) | **0.435** | **0.616** |
| | | | acc_norm | 0.383 | 0.608 (±0.020) | 0.4581 | 0.624 (±0.010) | 0.520 | 0.653 (±0.006) | **0.454** | **0.628** |
| | | **gemma-3-4b-it** | acc | 0.375 | 0.674 (±0.005) | 0.421 | 0.675 (±0.006) | 0.608 | 0.784 (±0.003) | **0.468** | **0.711** |
| | | | acc_norm | 0.394 | 0.683 (±0.008) | 0.432 | 0.680 (±0.007) | 0.626 | 0.789 (±0.004) | **0.484** | **0.718** |
| | **Qwen** | **Qwen2.5-3B** | acc | 0.477 | 0.599 (±0.005) | 0.495 | 0.600 (±0.009) | 0.482 | 0.585 (±0.007) | **0.485** | **0.595** |
| | | | acc_norm | 0.508 | 0.619 (±0.007) | 0.526 | 0.622 (±0.009) | 0.495 | 0.605 (±0.009) | **0.510** | **0.615** |
| | | **Qwen2.5-3B-Instruct** | acc | 0.539 | 0.723 (±0.005) | 0.667 | 0.758 (±0.006) | 0.568 | 0.6832 (±0.003) | **0.591** | **0.721** |
| | | | acc_norm | 0.558 | 0.734 (±0.005) | 0.686 | 0.771 (±0.006) | 0.586 | 0.697 (±0.004) | **0.610** | **0.734** |
| | **average** | | acc | **0.455** | **0.639** | **0.496** | **0.647** | **0.541** | **0.656** | **0.497** | **0.647** |
| | | | acc_norm | **0.477** | **0.653** | **0.516** | **0.660** | **0.560** | **0.671** | **0.518** | **0.662** |

| | | | | default prompt | | instruction prompt | | reasoning prompt | | average | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | zeroshot | fewshot (±SD) | zeroshot | fewshot (±SD) | zeroshot | fewshot (±SD) | zeroshot | fewshot (±SD) |
| **7B-8B** | **Llama** | **Llama-3.1-8B** | acc | 0.464 | 0.662 (±0.007) | 0.477 | 0.643 (±0.005) | 0.464 | 0.656 (±0.007) | **0.468** | **0.654** |
| | | | acc_norm | 0.487 | 0.676 (±0.010) | 0.504 | 0.657 (±0.008) | 0.493 | 0.673 (±0.007) | **0.495** | **0.669** |
| | | **Llama3.1-8B-Instruct** | acc | 0.492 | 0.706 (±0.004) | 0.564 | 0.698 (±0.007) | 0.513 | 0.692 (±0.005) | **0.523** | **0.699** |
| | | | acc_norm | 0.518 | 0.718 (±0.005) | 0.584 | 0.706 (±0.006) | 0.535 | 0.703 (±0.006) | **0.546** | **0.709** |
| | **gemma** | **gemma-7b** | acc | 0.471 | 0.648 (±0.006) | 0.492 | 0.621 (±0.005) | 0.504 | 0.646 (±0.002) | **0.489** | **0.638** |
| | | | acc_norm | 0.495 | 0.660 (±0.006) | 0.518 | 0.635 (±0.006) | 0.530 | 0.658 (±0.003) | **0.514** | **0.651** |
| | | **gemma-7b-it** | acc | 0.559 | 0.731 (±0.009) | 0.588 | 0.700 (±0.007) | 0.604 | 0.723 (±0.007) | **0.584** | **0.718** |
| | | | acc_norm | 0.598 | 0.748 (±0.009) | 0.610 | 0.719 (±0.008) | 0.629 | 0.744 (±0.006) | **0.612** | **0.737** |
| | **Qwen** | **Qwen2.5-7B** | acc | 0.504 | 0.622 (±0.010) | 0.499 | 0.623 (±0.009) | 0.554 | 0.635 (±0.007) | **0.519** | **0.626** |
| | | | acc_norm | 0.527 | 0.637 (±0.011) | 0.521 | 0.633 (±0.012) | 0.576 | 0.650 (±0.005) | **0.541** | **0.640** |
| | | **Qwen2.5-7B-Instruct** | acc | 0.568 | 0.713 (±0.007) | 0.545 | 0.700 (±0.004) | 0.543 | 0.707 (±0.005) | **0.552** | **0.707** |
| | | | acc_norm | 0.587 | 0.727 (±0.006) | 0.576 | 0.713 (±0.002) | 0.565 | 0.721 (±0.007) | **0.576** | **0.720** |
| | **Mistral** | **Mistral-7B-v0.3** | acc | 0.463 | 0.640 (±0.004) | 0.479 | 0.636 (±0.007) | 0.443 | 0.670 (±0.004) | **0.462** | **0.649** |
| | | | acc_norm | 0.487 | 0.654 (±0.003) | 0.499 | 0.650 (±0.008) | 0.473 | 0.685 (±0.006) | **0.486** | **0.663** |
| | | **Mistral-7B-Instruct-v0.3** | acc | 0.566 | 0.767 (±0.006) | 0.667 | 0.754 (±0.008) | 0.606 | 0.785 (±0.004) | **0.613** | **0.768** |
| | | | acc_norm | 0.593 | 0.779 (±0.006) | 0.686 | 0.767 (±0.009) | 0.631 | 0.796 (±0.006) | **0.636** | **0.781** |
| | **average** | | acc | **0.511** | **0.686** | **0.539** | **0.672** | **0.529** | **0.689** | **0.526** | **0.682** |
| | | | acc_norm | **0.565** | **0.700** | **0.586** | **0.685** | **0.578** | **0.704** | **0.576** | **0.696** |

Table 14: Total zero-shot and few-shot evaluation results on Thunder-NUBench. SD denotes standard deviation across random seeds or runs. *acc* denotes the accuracy, and *acc_norm* indicates the length-normalized accuracy. In all few-shot experiments, results are averaged over five random seeds (42, 1234, 3000, 5000, and 7000) and five examples of validation set data (5 shots).

## N Total Model Performance after SFT

We present model performance on Thunder-NUBench after supervised fine-tuning, highlighting improvements from task-specific training. Because SFT is performed on instruction-based data, only the instruction-based prompt format (prompt 7 from Table 10) is used for evaluation. We assess zero-shot performance exclusively to directly measure the model's ability from fine-tuning without the influence of in-context examples.

Note that we exclude both gemma-3-4b-it and gemma-3-4b-pt from our SFT experiments. During preliminary fine-tuning using LoRA, these models consistently produced unexpected errors across all configurations. Such failures are not limited to a specific checkpoint or setting, but appear to be a universal issue affecting the entire Gemma-3 series (Team et al., 2025). These errors stem from a structural incompatibility between the model's uploaded configuration and the expectations of the transformers library's auto-loading mechanism. [7]. We have initially selected these models to explore performance within the 3-4B parameter range in the Gemma family, but due to the unresolved errors during SFT, we are unable to include them in our final evaluation. Future work may revisit these models once the associated issues are addressed.

Table 15 presents model performance on Thunder-NUBench across three evaluation settings: zero-shot, few-shot (with 5-shot averaging), and after supervised fine-tuning (SFT). Supervised fine-tuning leads to substantial performance improvements compared to both baseline zero-shot and few-shot conditions. Across nearly all models and sizes, the SFT-zero-shot setting consistently outperforms few-shot baselines.

## O ScoNe Performance after Thunder-NUBench Fine-Tuning

ScoNe (**Sco**ped **Ne**gation Dataset) is a benchmark designed to evaluate whether models truly understand the semantic scope of negation in natural language, using controlled contrast sets in a natural language inference (NLI) framework with up to two negation morphemes per example. ScoNe consists of a train and test set, each of which has 5,010 and 100 instances (She et al., 2023).

In this experiment, we first measure the zero-shot baseline performance of various models on the

---

ScoNe dataset. We then evaluate the same models after supervised fine-tuning (SFT) on the Thunder-NUBench training set, without ever exposing them to ScoNe during training. This setup allows us to assess whether negation understanding learned from Thunder-NUBench generalizes to ScoNe.

Overall, we observe that performance improves moderately after fine-tuning (total average of accuracy: $0.678 \rightarrow 0.686$). This suggests that negation understanding learned from Thunder-NUBench transfers effectively to a structurally different evaluation set like ScoNe. Instruction-tuned models benefit the most from SFT on Thunder-NUBench, likely because the fine-tuning data follows an instruction format, to which they are already well adapted. On the other hand, some pretrained models exhibit performance degradation.

## P Total Analysis of Model Predictions on the Thunder-NUBench

We provide a detailed analysis of model predictions across negation types and sentence structures, examining common error patterns and behavioral trends.

Table 17 presents the incorrect choice distribution and confusion rates for local negation types across a range of 3-4B and 7-8B pretrained and instruction-tuned models under zero-shot and few-shot settings. We report few-shot results using a fixed random seed (1234), which corresponds to the default seed used in the LM Evaluation Harness framework. Averaging over multiple seeds was avoided, as it could obscure specific error patterns and make fine-grained confusion analysis less interpretable.

A notable observation is that `choice4` (paraphrase) was rarely selected as the incorrect answer, suggesting that models can generally distinguish paraphrases from negation. In particular, the Qwen2.5-7B-Instruct model under a few-shot setting achieved a remarkable result, never selecting `choice4` as the incorrect answer. This implies that even when structural transformations or synonyms are introduced, experimented models do not confuse them with negation, possibly because paraphrases tend to preserve the positive polarity of the original sentence, in contrast to the strong semantic shift introduced by local negation or contradiction.

Across all models, the most frequent source of confusion was `choice2` (local negation), often selected as the answer instead of standard negation.

---

| | | | | zeroshot | fewshot (±SD) | zeroshot sft |
|---|---|---|---|---|---|---|
| **3B-4B** | Llama | Llama-3.2-3B | acc | 0.463 | 0.616 (±0.005) | 0.792 |
| | | | acc_norm | 0.483 | 0.633 (±0.004) | 0.800 |
| | | Llama-3.2-3B-Instruct | acc | 0.490 | 0.619 (±0.006) | 0.808 |
| | | | acc_norm | 0.512 | 0.633 (±0.001) | 0.815 |
| | Qwen | Qwen2.5-3B | acc | 0.495 | 0.600 (±0.009) | 0.747 |
| | | | acc_norm | 0.526 | 0.622 (±0.009) | 0.753 |
| | | Qwen2.5-3B-Instruct | acc | 0.667 | 0.758 (±0.006) | 0.765 |
| | | | acc_norm | 0.686 | 0.771 (±0.006) | 0.775 |
| | | average | acc | **0.529** | **0.648** | **0.789** |
| | | | acc_norm | **0.552** | **0.665** | **0.778** |

| | | | | zeroshot | fewshot (±SD) | zeroshot sft |
|---|---|---|---|---|---|---|
| **7B-8B** | Llama | Llama-3.1-8B | acc | 0.477 | 0.643 (±0.005) | 0.851 |
| | | | acc_norm | 0.504 | 0.657 (±0.008) | 0.860 |
| | | Llama3.1-8B-Instruct | acc | 0.564 | 0.698 (±0.007) | 0.855 |
| | | | acc_norm | 0.584 | 0.706 (±0.006) | 0.868 |
| | gemma | gemma-7b | acc | 0.492 | 0.621 (±0.005) | 0.847 |
| | | | acc_norm | 0.518 | 0.635 (±0.006) | 0.849 |
| | | gemma-7b-it | acc | 0.588 | 0.700 (±0.007) | 0.838 |
| | | | acc_norm | 0.610 | 0.719 (±0.008) | 0.845 |
| | Qwen | Qwen2.5-7B | acc | 0.499 | 0.623 (±0.009) | 0.820 |
| | | | acc_norm | 0.521 | 0.633 (±0.012) | 0.825 |
| | | Qwen2.5-7B-Instruct | acc | 0.545 | 0.700 (±0.004) | 0.783 |
| | | | acc_norm | 0.576 | 0.713 (±0.002) | 0.784 |
| | Mistral | Mistral-7B-v0.3 | acc | 0.479 | 0.636 (±0.007) | 0.839 |
| | | | acc_norm | 0.499 | 0.650 (±0.008) | 0.851 |
| | | Mistral-7B-Instruct-v0.3 | acc | 0.667 | 0.754 (±0.008) | 0.815 |
| | | | acc_norm | 0.686 | 0.767 (±0.009) | 0.823 |
| | | average | acc | **0.539** | **0.672** | **0.831** |
| | | | acc_norm | **0.586** | **0.685** | **0.838** |

Table 15: Comparison of negation understanding performance on Thunder-NUBench across 3B–4B and 7B–8B models. Each model is evaluated in zero-shot, few-shot (5-shot average with standard deviation), and post-supervised fine-tuning (SFT) zero-shot settings. Metrics include raw accuracy (acc) and normalized accuracy (acc_norm). All evaluations use the instruction-based prompt format, prompt 7 from Table 10, consistent with the SFT data format.

This suggests that models struggle to differentiate sentence-level negation from subclausal negation. Within local negation types, the compound_part category exhibited the highest confusion rates across many models (e.g., 58.12% for Mistral-7B-v0.3 zero-shot, 56.68% for Llama-3.2-3B zero-

| model size | model series | Model Name | Baseline zeroshot | SFT (on Thunder-NUBench) zeroshot |
|---|---|---|---|---|
| 3-4B | Llama | Llama-3.2-3B | 0.585 | 0.546 |
| | | Llama-3.2-3B-Instruct | 0.577 | 0.676 |
| | Qwen | Qwen2.5-3B | 0.765 | 0.760 |
| | | Qwen2.5-3B-Instruct | 0.744 | 0.792 |
| | | **3-4B models average** | **0.668** | **0.694** |
| 7-8B | Llama | Llama-3.1-8B | 0.620 | 0.544 |
| | | Llama-3.1-8B-Instruct | 0.696 | 0.759 |
| | gemma | gemma-7b | 0.586 | 0.601 |
| | | gemma-7b-it | 0.740 | 0.743 |
| | Qwen | Qwen2.5-7B | 0.758 | 0.769 |
| | | Qwen2.5-7B-Instruct | 0.778 | 0.794 |
| | Mistral | Mistral-7B-v0.3 | 0.544 | 0.502 |
| | | Mistral-7B-Instruct-v0.3 | 0.740 | 0.740 |
| | | **7-8B models average** | **0.683** | **0.682** |
| | | **total average** | **0.678** | **0.686** |

Table 16: Zero-shot performance on the ScoNe benchmark before and after supervised fine-tuning (SFT) on the Thunder-NUBench training set. Each model is evaluated without seeing ScoNe data during training.

shot), indicating that when negation applies to one clause within a compound sentence, it is particularly challenging for models to resolve its scope relative to full-sentence negation.

| | | | | Error rate (1-acc) | Incorrect Choice Distribution | | | choice2 confusion rate | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | choice2 (%) | choice3 (%) | choice4 (%) | relative_part confuse rate | pp_part confuse rate | compound_part confuse rate | adverb_part confuse rate |
| **3-4B** | **Llama** | **Llama-3.2-3B** | **Baseline** zeroshot | 0.537 | 68.59 | 23.05 | 8.36 | 25.26 | 35.52 | 56.68 | 33.33 |
| | | | fewshot | 0.387 | 78.35 | 19.33 | 2.32 | 23.55 | 23.79 | 48.01 | 31.43 |
| | | | **After SFT** zeroshot | 0.208 | 82.21 | 14.90 | 2.88 | 11.26 | 17.24 | 27.08 | 12.38 |
| | | **Llama-3.2-3B-Instruct** | **Baseline** zeroshot | 0.510 | 67.51 | 30.14 | 2.35 | 29.01 | 33.45 | 47.29 | 30.48 |
| | | | fewshot | 0.379 | 70.79 | 28.42 | 0.79 | 26.62 | 22.07 | 33.21 | 33.33 |
| | | | **After SFT** zeroshot | 0.192 | 84.38 | 12.5 | 3.12 | 9.56 | 12.76 | 29.96 | 13.33 |
| | **Qwen** | **Qwen2.5-3B** | **Baseline** zeroshot | 0.501 | 66.14 | 30.08 | 3.78 | 31.40 | 33.45 | 40.07 | 30.48 |
| | | | fewshot | 0.392 | 71.76 | 23.16 | 5.09 | 25.26 | 22.76 | 38.63 | 33.33 |
| | | | **After SFT** zeroshot | 0.254 | 84.65 | 14.17 | 1.18 | 14.68 | 23.45 | 33.57 | 10.48 |
| | | **Qwen2.5-3B-Instruct** | **Baseline** zeroshot | 0.455 | 60.75 | 37.72 | 1.54 | 25.60 | 32.07 | 28.16 | 29.52 |
| | | | fewshot | 0.241 | 75.10 | 24.48 | 0.41 | 16.04 | 14.83 | 24.55 | 21.90 |
| | | | **After SFT** zeroshot | 0.236 | 80.93 | 16.95 | 2.12 | 12.97 | 21.03 | 29.24 | 10.48 |
| **7-8B** | **Llama** | **Llama-3.1-8B** | **Baseline** zeroshot | 0.523 | 70.80 | 22.90 | 6.30 | 29.35 | 34.83 | 53.43 | 34.29 |
| | | | fewshot | 0.357 | 81.01 | 16.20 | 2.79 | 22.53 | 21.38 | 47.29 | 29.52 |
| | | | **After SFT** zeroshot | 0.149 | 80.54 | 16.11 | 3.36 | 9.56 | 10.34 | 18.77 | 9.52 |
| | | **Llama-3.1-8B-Instruct** | **Baseline** zeroshot | 0.436 | 67.05 | 32.04 | 0.92 | 29.01 | 30.69 | 30.69 | 32.38 |
| | | | fewshot | 0.297 | 77.18 | 22.15 | 0.67 | 22.53 | 19.31 | 28.88 | 26.67 |
| | | | **After SFT** zeroshot | 0.145 | 84.83 | 13.10 | 2.07 | 9.22 | 11.72 | 18.05 | 11.43 |
| | **gemma** | **gemma-7b** | **Baseline** zeroshot | 0.508 | 68.17 | 28.88 | 2.95 | 28.67 | 33.10 | 46.93 | 35.24 |
| | | | fewshot | 0.382 | 83.81 | 14.62 | 1.57 | 26.28 | 23.45 | 49.10 | 38.10 |
| | | | **After SFT** zeroshot | 0.153 | 77.78 | 20.26 | 1.96 | 9.56 | 10.34 | 19.49 | 6.67 |
| | | **gemma-7b-it** | **Baseline** zeroshot | 0.412 | 61.99 | 34.14 | 3.87 | 17.06 | 30.00 | 34.66 | 21.90 |
| | | | fewshot | 0.298 | 65.22 | 32.44 | 2.34 | 14.33 | 16.55 | 29.96 | 20.95 |
| | | | **After SFT** zeroshot | 0.162 | 77.16 | 20.37 | 2.47 | 10.24 | 11.72 | 18.77 | 8.57 |
| | **Qwen** | **Qwen2.5-7B** | **Baseline** zeroshot | 0.501 | 66.14 | 30.08 | 3.78 | 31.40 | 33.45 | 40.07 | 30.48 |
| | | | fewshot | 0.387 | 74.23 | 24.23 | 1.55 | 24.91 | 24.48 | 41.52 | 27.62 |
| | | | **After SFT** zeroshot | 0.180 | 85.56 | 13.33 | 1.11 | 14.33 | 15.52 | 20.22 | 10.48 |
| | | **Qwen2.5-7B-Instruct** | **Baseline** zeroshot | 0.333 | 57.19 | 40.42 | 2.40 | 16.38 | 22.07 | 20.58 | 20.95 |
| | | | fewshot | 0.299 | 68.67 | 31.33 | 0.00 | 19.45 | 20.00 | 25.99 | 18.10 |
| | | | **After SFT** zeroshot | 0.217 | 84.79 | 14.29 | 0.92 | 14.68 | 15.86 | 28.88 | 14.29 |
| | **Mistral** | **Mistral-7B-v0.3** | **Baseline** zeroshot | 0.521 | 72.03 | 22.61 | 5.36 | 26.96 | 34.83 | 58.12 | 33.33 |
| | | | fewshot | 0.372 | 79.09 | 18.23 | 2.68 | 22.87 | 22.76 | 48.01 | 27.62 |
| | | | **After SFT** zeroshot | 0.161 | 83.85 | 13.66 | 2.48 | 8.87 | 11.03 | 25.27 | 6.67 |
| | | **Mistral-7B-Instruct-v0.3** | **Baseline** zeroshot | 0.333 | 57.49 | 41.92 | 0.60 | 20.48 | 21.72 | 18.05 | 18.10 |
| | | | fewshot | 0.236 | 70.34 | 29.24 | 0.42 | 13.65 | 14.14 | 25.63 | 13.33 |
| | | | **After SFT** zeroshot | 0.185 | 83.78 | 14.59 | 1.62 | 9.56 | 12.41 | 29.24 | 9.52 |

Table 17: Incorrect choice distribution and confusion analysis in Thunder-NUBench across 3-4B/7-8B, pretrained/instruction-tuned models.