

# Investigating Large Language Models for Complex Word Identification in Multilingual and Multidomain Setups

Anonymous ACL submission

## Abstract

Complex Word Identification (CWI) is an important step in the lexical simplification task and has recently become a task on its own. Some variations of this binary classification task have emerged, such as lexical complexity prediction (LCP) and complexity evaluation of multi-word expressions (MWE). Large language models (LLMs) recently became popular in the Natural Language Processing community because of their versatility and capability to solve unseen tasks in zero/few-shot settings. Our work investigates LLM usage, specifically open-source models such as Llama 2, Llama 3, and Vicuna v1.5, and closed-source, such as ChatGPT-3.5-turbo and GPT-4o, in the CWI, LCP, and MWE settings. We evaluate zero-shot, few-shot, and fine-tuning settings and show that LLMs struggle in certain conditions or achieve comparable results against existing methods. In addition, we provide some views on meta-learning combined with prompt learning. In the end, we conclude that the current state of LLMs cannot or barely outperform existing methods, which are usually much smaller.

## 1 Introduction

Complex word identification (CWI) aims to determine whether words or phrases are difficult for a target group of readers to understand. It is often used in lexical simplification, which targets replacing complex words and expressions with simplified alternatives (North et al., 2023a). CWI represents the first step, and it was treated as part of the lexical simplification task until 2012 when it became a standalone task (Shardlow, 2013).

CWI was initially addressed as a binary classification task (Paetzold and Specia, 2016), identifying whether a word is complex in a given sentence. When the task became more popular (North et al., 2023b), it was extended to the continuous domain as Lexical Complexity Prediction (LCP, also referred to as the probabilistic classification for CWI)

(Yimam et al., 2018) addressing multi-language and multi-domain settings. Then, it was extended to multi-word expressions (Shardlow et al., 2021). Recently, new datasets started to emerge in various languages and domains (Ortiz Zambrano and Montejo-Ráez, 2021; Venugopal et al., 2022; Ilgen and Biemann, 2023; Zambrano et al., 2023). Previous approaches to CWI ranged from using Support Vector Machines (S.P et al., 2016) to deep neural networks based on Bidirectional Representation from Encoder Transformers (Pan et al., 2021), multi-task learning with domain adaptation (Zaharia et al., 2022), and sequence modeling (Gooding and Kochmar, 2019).

With the recent breakthrough in large language models (LLMs), OpenAI showed that Generative Pre-trained Transformer (GPT) models (Radford et al., 2019; Brown et al., 2020) can improve performances on various natural language processing tasks as we scale up the model size and the amount of training data. After ChatGPT<sup>1</sup> was announced and showed its conversational capabilities, a race has started in developing and fine-tuning new models for general purpose and domain-specific applications using PaLM (Anil et al., 2023), LLaMA (Touvron et al., 2023), Orca (Mittra et al., 2023), Mistral (Jiang et al., 2023), GPT-4 (OpenAI et al., 2023), GPT-4o<sup>2</sup>, and many others.

Our work aims to provide the current state of LLMs in addressing CWI and LCP compared against state-of-the-art approaches. We focus on evaluating open-source (pre-trained Llama 2, Llama 3, Vicuna) and OpenAI’s close-source ChatGPT-3.5-turbo and GPT-4o. We summarize the contributions as follows: We evaluate LLMs in binary (discrete set of labels) and probabilistic classification (continuous space labels) on multi-domain and multi-lingual corpora. We employ

<sup>1</sup><https://openai.com/blog/chatgpt>

<sup>2</sup><https://openai.com/index/hello-gpt-4o/>

various techniques for prompting and fine-tuning. We show that LLMs struggle to address CWI and LCP tasks; however, in limited instances, they can achieve similar results with other more lightweight approaches. In the end, we analyze and provide an insight into where the models struggle.

## 2 Related Work

### 2.1 Complex Word Identification

Aroyehun et al. (2018) compared CNN-based models with various feature engineering methods based on tree ensembles and features, such as inverse term frequency, parts-of-speech tagging, WordNet, and word2vec, achieving comparable results. Fimmimore et al. (2019) proposed mono- and cross-lingual models based on simple features and logistic regression, achieving similar results to more complex, language-specific state-of-the-art models. Zaharia et al. (2020) employed zero- and few-shot learning techniques, along with Transformers and Recurrent Neural Networks, in a multilingual setting. Gooding and Kochmar (2019) considered CWI as a sequential task, using a bidirectional LSTM with word embeddings and character-level representations and a language modeling objective to learn the complexity of words given the context. Other approaches were used to improve CWI performances, such as graph-based (Ehara, 2019), domain adaptation (Zaharia et al., 2022), and transformer-based models (Pan et al., 2021; Cheng Sheang et al., 2022).

### 2.2 Large Language Models

LLMs were successfully utilized in various generative tasks (Pu et al., 2023; Chen et al., 2021). The new paradigm in solving other non-generative tasks is based on prompting pre-trained language models to perform the prediction task (Liu et al., 2023a; Sun et al., 2023). Fine-tuning models on instructions showed improved results in zero-shot settings, especially on unseen tasks (Wei et al., 2022a). Prompt-based methods such as the use of demonstrations (Min et al., 2022), intermediate reasoning steps by breaking down complex tasks into simpler subtasks (also known as a chain of thought) (Wei et al., 2022b), and using LLMs to optimize their prompts (Zhou et al., 2023) made zero-shot inference much more appealing due to reduced costs and more efficient than fine-tuning LLMs.

### 2.3 Prompt Tuning and Meta-Learning

Prompt Tuning (Lester et al., 2021) is a soft-prompting technique for adapting a large language model for a custom task without training the model’s parameters themselves. While successful, prompt tuning falls short when applied in the few-shot learning regime, leading to a combination with meta-learning. MetaPrompting (Hou et al., 2022) utilizes the meta-learning algorithms FOMAML (Finn et al., 2017) and Reptile (Nichol et al., 2018) to obtain optimized initialization embeddings. One shortcoming of this approach is the requirement for supervised training data during meta-learning, which is alleviated in other works (Pan et al., 2023; Huang et al., 2023) by generating the meta-training tasks in an unsupervised or self-supervised manner. Additionally, other works explored the use of an adaptable gradient regulating function (Li et al., 2023) and domain-adversarial neural networks (Fang et al., 2023), both techniques used to increase model generalization.

## 3 Method

### 3.1 Problem Formulation in the Pre-LLM Era

Word complexity can be defined as absolute and relative (North et al., 2023b). Absolute complexity is determined by the objective linguistic properties (e.g., semantic, morphological, phonological). In contrast, relative complexity is related to the subjective speaker’s point of view (e.g., familiarity with sound and meaning, culture, and context). In this work, we evaluate the relative complexity of words from non-native speakers’ points of view. Considering an annotated dataset  $D = \{(x_i, y_i)\}_{i=1}^N$  of  $N$  samples, the task can be viewed as a binary classification, CWI, where, given the pair  $x_i = (C_i, w_i)$  of a sentence  $C_i = (w_1, w_2, \dots)$  and a word  $w_i \in C_i$ , the system outputs  $y_i^{CWI} \in \{0, 1\}$  (i.e., complex or non-complex) (Paetzold and Specia, 2016). A variation of the CWI task is to evaluate the complexity  $y_i^{MWE} \in \{0, 1\}$  of a multi-word expression  $e_i = (w_1, w_2, \dots)$  containing multiple words  $w_j, j = 1 : |e|$ , from a given context  $C_i$  (i.e.,  $x_i = (C_i, e_i)$ ) (Shardlow et al., 2021). Later, CWI was considered in the continuous domain known as LCP, indicating the degree of difficulty  $y_i^{LCP} \in [0, 1]$  for the given word  $w_i \in C_i$  in the context  $C_i$  (Yimam et al., 2018).

### 3.2 Problem Formulation in the LLM Era

Starting from the previous formulation, we derive the formalism in the context of LLMs.

**Binary classification.** Given an example  $x_i = (C_i, w_i)$ , the model predicts if a given phrase  $w_i$  from the sentence  $C_i$  is complex. Especially in closed-source models, the access to the tokens’ logits is limited (e.g., OpenAI’s ChatGPT or GPT-4o). Therefore, we consider having access to only the model’s predicted text labels “true” or “false” (or any equivalent form) without a confidence estimation.

**Probabilistic classification.** In the existing approaches, the model produces a real value between 0 and 1, representing the degree of complexity for  $(C_i, w_i)$ . LLMs are known to suffer from hallucinations (OpenAI et al., 2023), and reliably predicting real values is challenging. We abide by Liu et al. (2023b)’s solution for estimating the scoring function. In a nutshell, instead of letting the model predict the probability as a number, we let the model predict discrete signals, and then we estimate the score through averaging. Specifically, we ask the model to predict on the 5-point Likert scale, in natural language, one of “very easy”, “easy”, “neutral”, “difficult”, or “very difficult”. This scale is converted into a numerical representation using the following mapping: very easy - 0, easy - 0.25, neutral - 0.5, difficult - 0.75, and very difficult - 1. Since LLMs output tokens from a probability distribution, we set the temperature to a higher value (in our experiments, we use 0.8) to increase variability in the responses. The numerical representation from LLM’s output is denoted as  $s_k \in S$  for a sampling step  $k$ , with  $S = \{0, 0.25, 0.5, 0.75, 1\}$ . The model’s probability of outputting one 5-point Likert score is  $p(s_k)$ . The final score is:

$$score = \mathbb{E}_p[S] = \sum_{s \in S} p(s) \cdot s \quad (1)$$

For experiments, we use the sample mean estimator  $\overline{score} = \frac{1}{K} \sum_{k=1}^K s_k$  of  $K$  sampling steps.

In essence, we want to simulate the data annotation process involved. Instead of employing multiple human annotators, we use the same model that produces a more randomized output by setting a higher temperature. Note that randomized LLM outputs do not translate to randomized output labels. Variance in the results can be used as a confidence estimator, and we show (see Section 6.4) that models act more deterministically despite

setting a high-temperature value. However, even for fine-tuned models, the variance is not 0; thus, it does not collapse to the multi-class classification setup (see Section 6.4).

### 3.3 Prompting LLMs

We instruct the model via a system prompt with the task and how to format the output. Then, we ask it through the user prompt to predict the example. Each example is prompted individually to avoid leaking knowledge from other examples.

All system prompts are listed in Appendix A. These are obtained after prompt engineering, i.e., multiple trial-and-errors. We focus on optimizing the prompt length and the LLM’s performance. Because performing prompt engineering on every LLM is time-consuming, we optimize the prompts on ChatGPT-3.5-turbo and LLama 2 7B models on LCP and CWI English, each as a whole task, respectively. The prompts for German and Spanish are translations of the English prompts.

We investigate prompting strategies for zero-shot and few-shot settings. In every setting, we evaluate with and without employing the chain of thoughts approach (Wei et al., 2022b) to reduce hallucination and keep the model focused on the task. Therefore, we ask the model to confirm the sentence and the word and then, before the final answer, provide a short proof about the reason for the response.

A similar prompting procedure is applied for open-source and close-source models. The difference mainly lies in how we format the query using either the Chat Message API<sup>3</sup> for OpenAI’s models or the chat template available in the HuggingFace tokenizer<sup>4</sup>. We provide details regarding the prompting format in Appendix A, and the evaluation protocol in Appendix C.

### 3.4 Fine-tuning LLMs

For fine-tuned models, we prepare the dataset to include a minimal system prompt and the query with the answer. First, we discretize the probabilities similar to Shardlow et al. (2021): scores between 0 and 0.2 are very easy, between 0.2 and 0.4 are easy, between 0.4 and 0.6 are neutral, between 0.6 and 0.8 are difficult, and between 0.8 and 1 are very difficult. Next, we apply the prompt for open-source models using the template specific to

<sup>3</sup><https://platform.openai.com/docs/api-reference/chat/create>

<sup>4</sup>[https://huggingface.co/docs/transformers/en/chat\\_templating](https://huggingface.co/docs/transformers/en/chat_templating)

the model available in the HuggingFace tokenizer. Fine-tuning OpenAI’s ChatGPT models involve uploading the training and validation files and starting the training job. After the fine-tuning step finishes, we follow a similar procedure to Section 3.3. However, in this setting, we do not ask the model to produce a demonstration; we only ask to confirm the task and directly provide the answer.

### 3.5 Meta-Learning

We propose using FOMAML in conjunction with Prompt tuning and P-tuning (Liu et al., 2022) in order to optimize the initial parameters of our adapters. The algorithm is described in Figure 1. For the loss function, we use the same causal loss as in fine-tuning. For the data, we select 45 tasks from the BIG-bench suite, as detailed in Subsection 4.1 and Appendix G.

---

#### Algorithm 1: FOMAML algorithm

---

**Data:**  $\alpha, \beta$  learning rates,  $n$  inner steps  
 Randomly initialize  $\theta$ ;  
**while** *not done* **do**  
   //Support and query sets  
   Sample task  $\mathcal{T} = (\mathcal{T}_s, \mathcal{T}_q)$ ;  
    $\theta'_0 \leftarrow \theta$ ;  
   **for**  $i = 0$  **to**  $n$  **do**  
      $\theta'_{i+1} \leftarrow \theta'_i - \alpha \nabla_{\theta'_i} \mathcal{L}(\theta'_i, \mathcal{T}_s)$ ;  
   **end**  
   //FOMAML optimization  
    $\nabla_{\theta} \mathcal{L}(\theta'_n, \mathcal{T}_q) \approx \nabla_{\theta'_n} \mathcal{L}(\theta'_n, \mathcal{T}_q)$   
    $\theta \leftarrow \theta - \beta \nabla_{\theta'_n} \mathcal{L}(\theta'_n, \mathcal{T}_q)$ ; //or Adam  
**end**

---

## 4 Experimental Setup

### 4.1 Datasets

**CompLex LCP 2021.** Proposed at SemEval 2021 Task 1 (Shardlow et al., 2021), CompLex LCP 2021 comprises around 10,000 sentences in English from three domains: European Parliament proceedings, the Bible, and biomedical literature. The data is split across two tasks: single-word (Single) and multi-word expressions (MWE). The complexity is provided as continuous values between 0 and 1, addressed as the probabilistic classification task. The average complexity is 0.3 for single and 0.42 for MWE. Each task dataset is split into trial/train/test splits, with 421/7662/917 samples for single words and 99/1517/184 for multi-word expressions.

**CWI Shared Dataset.** It was proposed at the CWI Shared Task in 2018 (Yimam et al., 2018) and addressed English multi-domain and multi-lingual settings. The English split contains samples from three sources (News, WikiNews, and Wikipedia) totaling approx. 35,000 samples. In the multi-lingual setting, the dataset features German and Spanish with approx. 8,000 and 17,600 samples, respectively, and a French test set containing 2,251 samples. The dataset was developed to address binary and probabilistic classification tasks by assigning probabilities and labels such that samples with 0% probability are non-complex and others as complex. We consider only the binary classification tasks (see Limitations 8). The English News is split as in 14k for training, 1.76k for validation and 2.1k for testing; English WikiNews has 7.75k for training, 870 for validation, and 1.29k for testing; and English Wikipedia has 5.55k samples for training, 694 for validation, and 870 samples for testing. In addition, German is split into train/validation/test as 6.15k/795/959, and Spanish is split into 13.8k/1.62k/2.23k samples.

**BIG-bench.** This recently proposed benchmark (Srivastava et al., 2023) contains over 200 tasks meant for evaluating large language models. We use this collection as part of the meta-learning stage in our pipeline. Since not all tasks are suitable for our use case, we select 45 tasks (detailed in Appendix G) and only pre-train on those. The unsuitable tasks may include non-categorical responses, requirements for external knowledge or may simply be too dissimilar to the target task. Our main task selection criteria were prompt length and similarity to the complex word identification task, since we want as much intrinsic knowledge to be transferred on fine-tuning.

### 4.2 Baselines

We compare against top-performing methods at CWI Shared task and LCP 2021. Camb (Gooding and Kochmar, 2018) employs heterogeneous features combined with an ensemble of AdaBoost classifiers. TMU system (Kajiwara and Komachi, 2018) uses a random forest classifier on multiple hand-crafted features. ITEC (De Hertog and Tack, 2018) combines CNN and LSTM layers. SB@GU (Alfter and Pilán, 2018) employs Random Forest and Extra Tree on top of multiple hand-crafted features. In addition, we include the XLM-RoBERTa-based approach combined with text simplification and domain adaptation (Zaharia et al., 2022), the

MLP combined with Sent2Vec solution Almeida et al. (2021), and RoBERTa<sub>LARGE</sub> with an ensemble of RoBERTa-based models (LR-Ensemble) (Pan et al., 2021).

### 4.3 Models

We evaluate several open- and closed-source LLMs. Specifically, for open-source models, we choose Llama 2 (7B and 13B) (Touvron et al., 2023), Vicuna (7B and 13B) (Zheng et al., 2024), and Llama 3 8B (AI@Meta, 2024). For closed-source models, we employ OpenAI’s ChatGPT-3.5-turbo and GPT-4o (OpenAI et al., 2023) specifically for their relatively lower prices compared with GPT-4 (see also Limitations 8). The chat model is used in the zero- and few-shot settings. Details regarding specific checkpoints for all models are listed in Appendix E.

### 4.4 Hyperparameters

**Inference.** During inference, we set the LLM to use a maximum of 4096 tokens, the repetition penalty was set to 1.2, and the temperature to 0.8. We set the top-k parameter to 10 and the top-p to 0.95. The open-source models were loaded with quantized parameters using the nf4 format through bitsandbytes (Cannizzo, 2018). For LCP, we set the number of inference steps  $N = 20$  for all open-source models, while for OpenAI models, we evaluate on  $N = 10$  inferences (see also the discussion from Appendix B).

**Fine-tuning open-source models.** The open-source models are downloaded from HuggingFace. For fine-tuning, we employ QLoRA (Detmners et al., 2023) with 4-bit quantization to reduce GPU memory usage.  $R$  was set to 16,  $\alpha$  to 32, and dropout to 0.05. The batch size was set between 10 and 32, and the learning rate using a linear scheduler with 10% warmup and a maximum value of  $1e-4$ . We trained the models using AdamW (Kingma and Ba, 2015) for three epochs.

**Fine-tuning GPT models.** We use OpenAI’s platform with the default hyperparameters for fine-tuning OpenAI models. We limited the training size to 250 samples uniformly sampled among labels from the train set specific to the dataset task and language.

**Meta-Learning.** For meta-learning, the inner learning rate is 0.1 and 0.03 for prompt learning and P-tuning, respectively, while the outer learning rates are 0.01 and 0.003. The batch size is set to 1, with the number of inner steps set to 5. We ran each

experiment for a total of 3000 steps. Additionally, during P-tuning, we follow the original work by choosing an LSTM (Hochreiter and Schmidhuber, 1997) as our model architecture. We use a total of 16 virtual tokens, and the support and query sets contain 6 examples from the same task each.

### 4.5 Evaluation Metrics

We adopt the same evaluation methodology as in Shardlow et al. (2021) for CompLex and Yimam et al. (2018) for CWI datasets. Therefore, we use Pearson correlation (P) and Mean Average Error (MAE) on the CompLex dataset and F1-score (F1) for the CWI dataset. We also include accuracy (Acc) on the CWI dataset. We report all results on a single run for CWI and multiple runs (described by  $N$ ) for LCP.

## 5 Results

### 5.1 English Multi-Domain Setup

We present the results in Table 1. The top-performing LLMs are ChatGPT-3.5-turbo and GPT-4o, which generally achieve higher scores than other open-source LLMs, especially in the zero- and few-shot settings. When fine-tuning, we notice that open-source models achieve competitive results with ChatGPT-3.5-turbo-ft. Fine-tuned ChatGPT-3.5-turbo achieves over 80% F1-score, while the highest scores for English-News and English-Wikipedia are surpassed by Llama-3-8b-ft and Vicuna-v1.5-13b-ft, respectively, by 1-2%. However, LLMs fall behind baseline classifiers which are more lightweight and easier to run. On the Wikipedia domain, Vicuna-v1.5-13b-ft achieves the same F1-score as Camb. The main limitation we found is that LLMs don’t fully understand the task (see Section 6.1).

### 5.2 Multi-Lingual Setup

The results are presented in Table 1 for the German (De) and Spanish (Es) languages. Similar to the multi-domain setup, the fine-tuned LLMs achieve the highest score. Notably, Llama-2-7b-ft and ChatGPT-3.5-turbo-ft achieve higher scores than the submitted systems, but we cannot consider LLMs a good solution for this problem as these models achieve under 80% in F1-score. Because LLMs were trained with corpora from multiple languages, they perform similarly in German and Spanish. Zero-shot combined with the chain of

Model	En-N $\uparrow$	En-WN $\uparrow$	En-W $\uparrow$	De $\uparrow$	Es $\uparrow$
Camb	<b>87.3</b>	<b>84.0</b>	<b>81.2</b>	-	-
ITEC	86.4	81.1	78.1	-	76.3
TMU	86.3	78.7	76.1	<b>74.5</b>	<b>77.0</b>
<i>Zero-shot</i>					
Llama-2-7b-chat	32.1	19.3	37.8	49.0	30.7
Llama-2-13b-chat	11.9	12.5	20.1	44.2	56.7
Vicuna-v1.5-7b	22.6	25.3	27.8	51.5	18.2
Vicuna-v1.5-13b	13.0	12.0	16.0	11.7	53.4
Llama-3-8b-chat	43.0	29.3	43.1	10.7	50.2
ChatGPT-3.5-turbo	40.1	37.0	45.6	35.3	53.3
GPT-4o	<b>65.9</b>	<b>64.2</b>	<b>66.8</b>	<b>68.9</b>	<b>63.3</b>
<i>Zero-shot CoT</i>					
Llama-2-7b-chat	56.5	61.1	62.8	56.7	57.2
Llama-2-13b-chat	54.9	49.1	57.8	57.7	55.5
Vicuna-v1.5-7b	38.1	38.5	52.4	55.4	20.6
Vicuna-v1.5-13b	32.9	27.9	33.3	29.8	42.1
Llama-3-8b-chat	50.5	45.7	61.0	34.9	49.3
ChatGPT-3.5-turbo	64.0	64.0	66.7	63.4	59.1
GPT-4o	<b>72.9</b>	<b>74.9</b>	<b>76.1</b>	<b>70.3</b>	<b>69.5</b>
<i>Few-shot</i>					
Llama-2-7b-chat	61.4	<b>63.2</b>	<b>70.6</b>	57.9	55.8
Llama-2-13b-chat	46.2	51.2	52.2	49.5	53.5
Vicuna-v1.5-7b	43.9	46.0	48.1	45.1	50.4
Vicuna-v1.5-13b	42.3	53.4	54.9	44.6	54.6
Llama-3-8b-chat	53.0	55.3	61.7	54.1	56.6
ChatGPT-3.5-turbo	52.1	44.5	53.8	42.5	55.4
GPT-4o	<b>63.8</b>	60.7	58.3	<b>66.0</b>	<b>60.2</b>
<i>Few-shot CoT</i>					
Llama-2-7b-chat	54.9	<b>60.7</b>	<b>67.5</b>	58.6	43.5
Llama-2-13b-chat	45.2	56.3	59.3	58.5	58.0
Vicuna-v1.5-7b	39.6	44.0	56.2	56.4	32.7
Vicuna-v1.5-13b	49.1	45.8	50.9	57.0	59.2
Llama-3-8b-chat	34.8	47.4	53.6	58.2	61.6
ChatGPT-3.5-turbo	58.3	51.4	55.5	68.1	<b>64.3</b>
GPT-4o	<b>66.1</b>	53.3	61.2	<b>60.5</b>	53.5
<i>Fine-tuned</i>					
Llama-2-7b-ft	78.0	78.2	77.4	<b>74.6</b>	70.5
Llama-2-13b-ft	77.6	77.7	73.1	70.8	75.3
Vicuna-v1.5-7b-ft	80.2	76.8	77.2	72.1	70.0
Vicuna-v1.5-13b-ft	81.2	77.4	<b>81.2</b>	73.0	67.5
Llama-3-8b-ft	<b>82.1</b>	79.6	76.8	72.2	70.8
ChatGPT-3.5-turbo-ft	80.7	<b>80.9</b>	80.2	66.6	<b>78.1</b>

Table 1: The F1-scores on the test sets from CWI 2018 Shared Dataset. Notation: En - English dataset, De - German, and Es - Spanish; for English datasets, N - News domain, WN - WikiNews, W - Wikipedia. In bold, we denote the best score.

449 thought performs better than other prompting tech-  
450 niques in most cases.

### 451 5.3 Lexical Complexity Prediction Setup

452 On the CompLex LCP dataset, Pan et al. (2021)  
453 achieved the best scores. Refer to Table 2 for  
454 the results. Fine-tuned LLM-based models outper-  
455 form RoBERTa-based models on the MWE task,  
456 the best-performing model being Llama-2-13b-ft.  
457 However, RoBERTa<sub>LARGE</sub> has 355M parameters,  
458 while Llama 2 has 37 times more parameters, and  
459 the performance difference is only about 5% in per-  
460 son correlation. On the single-word expressions  
461 task, RoBERTa<sub>LARGE</sub> outperforms all models. When  
462 considering the prompting techniques, the few-shot

Model	Single-Word		Multi-Word	
	P $\uparrow$	MAE $\downarrow$	P $\uparrow$	MAE $\downarrow$
MLP+Sent2Vec	.4598	.0866	.3941	.1145
XLM-RoBERTa-based	.7744	.0652	.8285	.0708
RoBERTa <sub>LARGE</sub>	<b>.7903</b>	<b>.0648</b>	.7900	.0753
LR-Ensemble	-	-	<b>.8612</b>	<b>.0616</b>
<i>Zero-shot</i>				
Llama-2-7b-chat	.3133	.3061	.5200	.2316
Llama-2-13b-chat	.2039	.2474	.3613	.1737
Vicuna-v1.5-7b	.3108	.3684	.4502	.2680
Vicuna-v1.5-13b	.2189	.1987	.4436	.1425
Llama-3-8b-chat	.3816	.1880	.6271	.1626
ChatGPT-3.5-turbo	.5352	<b>.1447</b>	.6284	<b>.1529</b>
GPT-4o	<b>.5953</b>	.2346	<b>.7753</b>	.2377
<i>Zero-shot CoT</i>				
Llama-2-7b-chat	.3617	.1698	.5040	.1632
Llama-2-13b-chat	.4429	.1355	.5905	.1118
Vicuna-v1.5-7b	.2558	.1504	.4916	.1310
Vicuna-v1.5-13b	.4664	<b>.0922</b>	.6357	<b>.1049</b>
Llama-3-8b-chat	.4617	.1507	.6923	.1167
ChatGPT-3.5-turbo	.5901	.2012	.6836	.1624
GPT-4o	<b>.6228</b>	.2145	<b>.7389</b>	.2586
<i>Few-shot</i>				
Llama-2-7b-chat	.1409	.2021	.5016	.1781
Llama-2-13b-chat	.2010	.2178	.4412	.2118
Vicuna-v1.5-7b	.1789	.1767	.4641	.1522
Vicuna-v1.5-13b	.2686	.1871	.4401	.2157
Llama-3-8b-chat	.3102	.1730	.5843	.1796
ChatGPT-3.5-turbo	.6385	<b>.0979</b>	.6742	<b>.1197</b>
GPT-4o	<b>.7111</b>	.1859	<b>.8284</b>	.2195
<i>Few-shot CoT</i>				
Llama-2-7b-chat	.4683	.1988	.5920	.2170
Llama-2-13b-chat	.5796	<b>.1289</b>	.6468	.1615
Vicuna-v1.5-7b	.5832	.1315	.6463	.1444
Vicuna-v1.5-13b	.5576	.1477	.6832	.1524
Llama-3-8b-chat	.2723	.2048	.7148	.1146
ChatGPT-3.5-turbo	.7175	.1421	.7568	<b>.1707</b>
GPT-4o	<b>.7594</b>	.1609	<b>.8211</b>	.1850
<i>Fine-tuned</i>				
Llama-2-7b-ft	.7732	<b>.0670</b>	.7919	.0766
Llama-2-13b-ft	<b>.7815</b>	.0797	<b>.8317</b>	<b>.0717</b>
Vicuna-v1.5-7b-ft	.7613	.0840	.7862	.0782
Vicuna-v1.5-13b-ft	.7530	.0914	.8000	.0763
Llama-3-8b-chat-ft	.7497	.0909	.7800	.0834
ChatGPT-3.5-turbo-ft	.7372	.1379	.7493	.1834

Table 2: Results on the CompLex LCP 2021 dataset. In bold we denote the best score, and underlined are the second-best results.

Model	News		WikiNews		Wikipedia	
	F1 $\uparrow$	Acc $\uparrow$	F1 $\uparrow$	Acc $\uparrow$	F1 $\uparrow$	Acc $\uparrow$
<i>P-tuning</i>						
Llama-2-7b-chat	50.3	<b>46.7</b>	<b>66.8</b>	51.2	<b>65.3</b>	<b>51.9</b>
Llama-2-7b	<b>53.8</b>	41.2	65.4	<b>52.0</b>	61.2	49.8
<i>Prompt-tuning</i>						
Llama-2-7b-chat	46.8	<b>53.0</b>	<b>66.1</b>	51.4	<b>61.9</b>	48.6
Llama-2-7b	<b>51.6</b>	43.8	64.0	<b>53.2</b>	61.3	<b>49.7</b>

Table 3: Results on the multi-domain English test set from CWI 2018 Shared Dataset in the few-shot learning regime, starting from the meta-learned models. In bold, we denote the best score.

Model	Fine-tuning inner steps					
	5	10	15	25	50	100
	<i>P-tuning</i>					
<b>Llama-2-7b-chat</b>	<b>66.8</b>	66.3	66.7	64.2	65.7	59.7
<b>Llama-2-7b</b>	60.7	62.0	<b>65.4</b>	64.8	<b>65.4</b>	64.3
	<i>Prompt-tuning</i>					
<b>Llama-2-7b-chat</b>	64.8	65.7	<b>66.1</b>	65.4	62.4	61.8
<b>Llama-2-7b</b>	<b>64.2</b>	58.8	55.8	52.2	59.6	64.0

Table 4: F1-score on the WikiNews English test set from CWI 2018 Shared Dataset. In bold, we denote the best score.

combined with the chain of thought method usually performs better.

## 5.4 Meta-Learning Setup

Due to the large computational cost of our meta-learning algorithm, we only test on Llama-2-7b, both the chat and default versions. Also, we only test using data in English, since changing the meta-training tasks requires a new suite of tasks, which can be difficult to obtain in the multilingual setting. The reasons for this are both less data availability and lack of knowledge in the other languages.

The soft prompting techniques show results comparable to the zero-shot regime, as shown in Table 3. All combinations of techniques and chat versus base versions of the LLMs show similar performances. In addition, we show the performance impact of varying the optimization steps our meta-learner goes through before the evaluation process in Table 4. The best number of steps is 5-15 for the chat versions of Llama-2, as opposed to 50-100 for the base model.

## 6 Discussions

### 6.1 LLMs’ Task Understanding

We investigate what is the LLMs’ level to understand the task firsthand before generating the prediction. That is, before providing the answer, we check if the model can reproduce what it has to solve (i.e., the sentence and target word). We report the sentence error count Z(S) and the word error count (W). In this setting, we mainly focus on the Llama2 and ChatGPT-3.5-turbo chat models, in the zero-shot chain of thought setting. In the CWI setting, we obtained an output using various packages such as Outlines (Willard and Louf, 2023), but it was not correlated with the examples, and the overall performance was not better than random. The results for the chat models on English domains are presented in Table 5. When the model

is larger, in general, the error rates decrease. The ChatGPT model obtains the lowest error counts, while the Llama 2 7b model obtains the highest. In general, the models struggle to understand what is the word they need to evaluate. Investigating the errors, we mostly see that the model considers more words than the target, for example, "America" (ground truth) vs "South America" (extracted by LLM). Other error cases we identified were completely different words to evaluate. For example, the target “years” was replaced by Llama-2-13b-chat with “Aegyptosaurus”. Text locality is not always the main reason; in the previous examples, in the first one, we have locality; in the second one, the words were in different parts of the sentence.

In the LCP setting (see Table 6), we consider all the sampling runs, and thus, we report the average and standard deviation across those runs. We report lower absolute error counts. Similar to the previous setting, we note that the sentence error count is lower than the word error count, in most cases being closer to 0. In addition, ChatGPT achieves error counts very close to 0, meaning that the model understands the task it needs to solve. In the case of Llama-2-7b, the models struggle to recall the word.

Model	News		Wikinews		Wikipedia	
	S	W	S	W	S	W
<b>Llama-2-7b-chat</b>	50	245	120	190	61	85
<b>Llama-2-13b-chat</b>	36	225	44	173	93	125
<b>ChatGPT-3.5-turbo</b>	2	47	4	17	0	10

Table 5: LLMs’ task understanding capabilities on the CWI English multi-domain dataset. The S column indicates wrong sentences, and the W column indicates wrong words.

Model	LCP-single		LCP-multi	
	S	W	S	W
<b>Llama-2-7b-chat</b>	2.4±0.7	0.1±0.2	1.0±0.2	6.5±0.5
<b>Llama-2-13b-chat</b>	0±0	0.1±0.3	0±0	3.9±1.2
<b>ChatGPT-3.5-turbo</b>	0±0	0±0	0.1±0.3	0±0

Table 6: LLMs’ task understanding capabilities on the LCP English datasets. The S column indicates wrong sentences, and the W column indicates wrong words.

### 6.2 LLMs’ Difficulty Understanding

In the zero-shot learning stage, before letting the model output the answer, we ask it to provide a brief proof regarding the choice. We ask first about the proof and then ask for the answer to enforce the model to "think before answer". If we let the model

answer and then provide proof, the proof would have been influenced by the initial answer, which would have influenced the model’s internal bias. In Table 17, we show some examples of reasoning regarding the answer provided by Llama-2-13b-chat on the CWI English dataset. The proof motivates the answer in our setting, but we notice some flaws in the reasoning. For example, the model says that "ft" (i.e., feet as a unit of measurement) is common in English. Meanwhile, it tends to contradict that being an abbreviation makes it difficult to understand. We notice this pattern quite often in the Llama models.

### 6.3 Confusion Matrices on CWI

To investigate how the predictions are affected by the domain, language, and LLM, we generate the confusion matrices, which are shown in Figures 3 and 4. The general tendency is that chat models have higher false-positive or false-negative rates. The same model checkpoints have the same bias towards one false rate in the multi-domain setting. For example, Llama-2-7b-chat has a high false-positive rate, while Llama-2-13b-chat has a high false rate. Correlated with the proofs generated by the LLMs, this is motivated by the fact that LLMs tend to either overestimate or underestimate the difficulty of a word. This is especially true if the model finds a synonym for the target word. Also, the high false rates correlate with the model’s incapacity to understand the task. On the other hand, fine-tuned models show lower false-positive/negative rates, meaning that fine-tuning makes the model focus better on the task and learn latent instructions directly from the data.

### 6.4 Fine-tuned Predictions on LCP

We analyzed the complexity probability distribution outputted by the LLMs in Figures 5, 8, 6, 7, 8 in Appendix H. This is constructed by binning the models’ real-valued estimates (on the x-axis) and generating a histogram (on the y-axis). The discrete labels were mapped equidistantly in the range 0-1, i.e., very easy (VE) in 0-0.2, easy (E) in 0.2-0.4, neutral (N) in 0.4-0.6, difficult (D) in 0.6-0.8, and very difficult (VE) in 0.8-1. In gray, we indicate the outside of the expected label (i.e., wrong labels); in the white stripe, we indicate the correctly predicted labels.

In the case of chat models, we notice a more uniform distribution among models’ predictions, especially for the low-complexity words. The ab-

solute error is more than one step in the difficulty scale. We notice that the models struggle to identify the very difficult label, regardless of whether the model was fine-tuned or not. In the fine-tuned setting, we notice that Llama-based models tend to misclassify neutral and difficult words, generally considering the words easier than the ground truth. Also, there is a tendency to label very easy words as easy. In the case of ChatGPT-3.5-turbo-ft, we notice that the outputs tend to be more deterministic – the majority of labels lie on the class scores.

## 6.5 Results Discussions

LLMs can grasp word complexity, depending on the model’s capabilities. We observed that performances across domains, language, and whether we deal with a word or a phrase, are similar if the model is fine-tuned. In the zero-shot setting, the input prompt and prediction temperature yield a high variance across the results. Also, we noticed that sometimes the models (especially Llama-2-13b-chat, in the zero-shot setting) refused to answer some examples (especially in the Biblical domain) because of racial discrimination, despite that not being the case. Models tend to consider words easier than they are, mainly because if asked to explain the choice, they could find another synonym that is not necessarily simpler. Zero-shot prompting is achieved every time poor performances are detected, and the main effect is that models tend to have a high false positive rate in the CWI task. This can be changed during fine-tuning when we notice that imbalanced datasets towards a class lead to the model being biased and producing more often the predominant label from the fine-tuning set.

## 7 Conclusions

In conclusion, we addressed CWI and LCP using LLMs, specifically Llama-based and OpenAI’s GPT models. We observed that these models can determine the word difficulty level in multiple domains and languages, although with limited performances. Meanwhile, these models struggle to label very difficult phrases correctly. Future directions imply investigating multiple models in more languages. Also, as we noticed that the prompts and example selection greatly influence the models’ performance, other future work should rely on reducing hallucination and determining which adversarial examples affect the model’s capabilities most in the context of CWI.



## 8 Limitations

Our approach has some limitations regarding prompt design. During experiments, we noticed that prompt design can highly influence the results, especially in the case of zero-shot settings. Using the same prompt across all models is not optimal, but we tried to find those instructions that benefit all models. Providing the model with specific instructions helps the model to better focus on the task and reduce hallucination. One way to mitigate hallucinations was to use a specific JSON format (see Appendix A), which the model required to confirm the task.

Also, we know that random sampling is not the optimal solution for choosing fine-tuning examples for ChatGPT-3.5-turbo. The size and quality of data greatly impact the prediction performance. To reduce this effect, we created a balanced dataset among label difficulties, such that the model equally sees easy and difficult words. We also kept a uniform distribution among complexity probabilities strictly greater than zero for both tasks (CWI and LCP).

## 9 Ethical Considerations

Since we used pre-trained LLMs, all their limitations apply to our work. Developing CWI and LCP systems can be beneficial for new language learners (e.g., chat-based applications in which LLMs help new language learners to understand difficult words and even provide alternatives), however, because of hallucination and inaccuracies that such models may provide, these systems can violate codes of ethics and harm or address attacks to such individuals. We are aware of the fast-paced development in the LLM area, and we think this area of research needs some attention. Therefore, we will make the fine-tuned models publicly available for transparency and fair comparison with feature works. These models should only be used for research. All the data we used is already publicly available, and the pre-trained Llama models are available on HuggingFace<sup>5</sup>, under the Llama 2 License Agreement<sup>6</sup>. We did not use the resources for other purposes than the ones allowed.

<sup>5</sup><https://huggingface.co/meta-llama>

<sup>6</sup><https://github.com/facebookresearch/llama/blob/main/LICENSE>

## References

- AI@Meta. 2024. [Llama 3 model card](#). 677
- David Alfter and Ildikó Pilán. 2018. [SB@GU at the complex word identification 2018 shared task](#). In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 315–321, New Orleans, Louisiana. Association for Computational Linguistics. 678–682
- Raul Almeida, Hegler Tissot, and Marcos Didonet Del Fabro. 2021. [C3SL at SemEval-2021 task 1: Predicting lexical complexity of words in specific contexts with sentence embeddings](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 683–687, Online. Association for Computational Linguistics. 684–690
- Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez Abrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jeffrey Hui, Jeremy Hurwitz, Michael Isard, Abe Ittycheriah, Matthew Jagielski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun, Sneha Kudugunta, Chang Lan, Katherine Lee, Benjamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu, Frederick Liu, Marcello Maggioni, Aroma Mahendru, Joshua Maynez, Vedant Misra, Maysam Moussalem, Zachary Nado, John Nham, Eric Ni, Andrew Nystrom, Alicia Parrish, Marie Pellat, Martin Polacek, Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan Richter, Parker Riley, Alex Castro Ros, Aurko Roy, Brennan Saeta, Rajkumar Samuel, Renee Shelby, Ambrose Slone, Daniel Smilkov, David R. So, Daniel Sohn, Simon Tokumine, Dasha Valter, Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pidong Wang, Zirui Wang, Tao Wang, John Wieting, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav Petrov, and Yonghui Wu. 2023. [Palm 2 technical report](#). 691–730
- Segun Taofeek Aroyehun, Jason Angel, Daniel Alejandro Pérez Alvarez, and Alexander Gelbukh. 2018. [Complex word identification: Convolutional neural network vs. feature engineering](#). In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP* 731–734

736	<i>for Building Educational Applications</i> , pages 322–327, New Orleans, Louisiana. Association for Computational Linguistics.	<i>IEEE International Conference On Machine Learning And Applications (ICMLA)</i> , pages 1982–1986.	794 795
739	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. <a href="#">Language models are few-shot learners</a> .	Feiteng Fang, Min Yang, Chengming Li, and Ruifeng Xu. 2023. Adversarial meta prompt tuning for open compound domain adaptive intent detection. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 1791–1795.	796 797 798 799 800 801
750	Fabio Cannizzo. 2018. <a href="#">A fast and vectorizable alternative to binary search in o(1) with wide applicability to arrays of floating point numbers</a> . <i>Journal of Parallel and Distributed Computing</i> , 113:37–54.	Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In <i>International conference on machine learning</i> , pages 1126–1135. PMLR.	802 803 804 805
751		Pierre Finamore, Elisabeth Fritzsche, Daniel King, Alison Sneyd, Aneeq Ur Rehman, Fernando Alva-Manchego, and Andreas Vlachos. 2019. <a href="#">Strong baselines for complex word identification across multiple languages</a> . In <i>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)</i> , pages 970–977, Minneapolis, Minnesota. Association for Computational Linguistics.	806 807 808 809 810 811 812 813 814 815
752		Sian Gooding and Ekaterina Kochmar. 2018. <a href="#">CAMB at CWI shared task 2018: Complex word identification with ensemble-based voting</a> . In <i>Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 184–194, New Orleans, Louisiana. Association for Computational Linguistics.	816 817 818 819 820 821 822
753		Sian Gooding and Ekaterina Kochmar. 2019. <a href="#">Complex word identification as a sequence labelling task</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1148–1153, Florence, Italy. Association for Computational Linguistics.	823 824 825 826 827 828
754	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. <a href="#">Evaluating large language models trained on code</a> . <i>CoRR</i> , abs/2107.03374.	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural computation</i> , 9(8):1735–1780.	829 830 831
755		Yutai Hou, Hongyuan Dong, Xinghao Wang, Bohan Li, and Wanxiang Che. 2022. Metaprompting: Learning to learn better prompts. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 3251–3262.	832 833 834 835 836
756		Yukun Huang, Kun Qian, and Zhou Yu. 2023. Learning a better initialization for soft prompts via meta-learning. In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 67–75.	837 838 839 840 841 842 843
757		Bahar Ilgen and Chris Biemann. 2023. <a href="#">Cwitr: A corpus for automatic complex word identification in turkish texts</a> . In <i>Proceedings of the 2022 6th International Conference on Natural Language Processing and Information Retrieval, NLPPIR '22</i> , page 157–163, New York, NY, USA. Association for Computing Machinery.	844 845 846 847 848 849 850
758	Kim Cheng Sheang, Anaïs Koptient, Natalia Grabar, and Horacio Saggion. 2022. <a href="#">Identification of complex words and passages in medical documents in French</a> . In <i>Actes de la 29e Conférence sur le Traitement Automatique des Langues Naturelles. Volume 1 : conférence principale</i> , pages 116–125, Avignon, France. ATALA.		
759	Dirk De Hertog and Anaïs Tack. 2018. <a href="#">Deep learning architecture for complex word identification</a> . In <i>Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 328–334, New Orleans, Louisiana. Association for Computational Linguistics.		
760			
761			
762			
763			
764			
765			
766			
767			
768			
769			
770			
771			
772			
773			
774			
775			
776			
777			
778			
779			
780			
781			
782			
783			
784			
785			
786			
787			
788	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. <a href="#">Qlora: Efficient finetuning of quantized llms</a> .		
789			
790			
791	Yo Ehara. 2019. <a href="#">Graph-based analysis of similarities between word frequency distributions of various corpora for complex word identification</a> . In <i>2019 18th</i>		
792			
793			

851	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023. <a href="#">Mistral 7b</a> .	
852		
853		
854		
855		
856		
857		
858	Tomoyuki Kajiwara and Mamoru Komachi. 2018. <a href="#">Complex word identification based on frequency in a learner corpus</a> . In <i>Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 195–199, New Orleans, Louisiana. Association for Computational Linguistics.	
859		
860		
861		
862		
863		
864		
865	Diederik P. Kingma and Jimmy Ba. 2015. <a href="#">Adam: A method for stochastic optimization</a> . In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	
866		
867		
868		
869		
870	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	
871		
872		
873		
874		
875		
876		
877	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3045–3059.	
878		
879		
880		
881		
882	Juncheng Li, Minghe Gao, Longhui Wei, Siliang Tang, Wenqiao Zhang, Mengze Li, Wei Ji, Qi Tian, Tat-Seng Chua, and Yueting Zhuang. 2023. Gradient-regulated meta-prompt learning for generalizable vision-language models. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 2551–2562.	
883		
884		
885		
886		
887		
888		
889	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023a. <a href="#">Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing</a> . <i>ACM Comput. Surv.</i> , 55(9).	
890		
891		
892		
893		
894	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. <a href="#">P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 61–68, Dublin, Ireland. Association for Computational Linguistics.	
895		
896		
897		
898		
899		
900		
901		
902	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. <a href="#">G-eval: NLG evaluation using gpt-4 with better human alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	
903		
904		
905		
906		
907		
908		
	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. <a href="#">Rethinking the role of demonstrations: What makes in-context learning work?</a> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	909
		910
		911
		912
		913
		914
		915
		916
	Arindam Mitra, Luciano Del Corro, Shweti Mahajan, Andres Coda, Clarisse Simoes, Sahaj Agarwal, Xuxi Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Agarwal, Hamid Palangi, Guoqing Zheng, Corby Rosset, Hamed Khanpour, and Ahmed Awadallah. 2023. <a href="#">Orca 2: Teaching small language models how to reason</a> .	917
		918
		919
		920
		921
		922
		923
	Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms. <i>arXiv preprint arXiv:1803.02999</i> .	924
		925
		926
	Kai North, Tharindu Ranasinghe, Matthew Shardlow, and Marcos Zampieri. 2023a. <a href="#">Deep learning approaches to lexical simplification: A survey</a> .	927
		928
		929
	Kai North, Marcos Zampieri, and Matthew Shardlow. 2023b. <a href="#">Lexical complexity prediction: An overview</a> . <i>ACM Comput. Surv.</i> , 55(9).	930
		931
		932
	OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Sim��n Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook	933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967

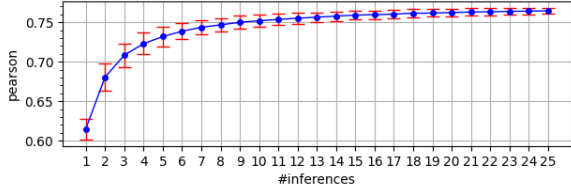
968	Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. <a href="#">Gpt-4 technical report</a> .	
969		
970		
971		
972		
973		
974		
975		
976		
977		
978		
979		
980		
981		
982		
983		
984		
985		
986		
987		
988		
989		
990		
991		
992		
993		
994		
995		
996		
997		
998		
999		
1000		
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010		
1011		
1012		
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020	Jenny A. Ortiz Zambrano and Arturo Montejó-Ráez. 2021. <a href="#">CLexIS2: A new corpus for complex word identification research in computing studies</a> . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 1075–1083, Held Online. INCOMA Ltd.	
1021		
1022		
1023		
1024		
1025		
1026		
1027	Gustavo Paetzold and Lucia Specia. 2016. <a href="#">SemEval 2016 task 11: Complex word identification</a> . In <i>Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)</i> , pages 560–569, San Diego, California. Association for Computational Linguistics.	1030
1028		1031
1029		1032
	Chunguang Pan, Bingyan Song, Shengguang Wang, and Zhipeng Luo. 2021. <a href="#">DeepBlueAI at SemEval-2021 task 1: Lexical complexity prediction with a deep ensemble approach</a> . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 578–584, Online. Association for Computational Linguistics.	1033
		1034
		1035
		1036
		1037
		1038
		1039
	Kaihang Pan, Juncheng Li, SONG Hongye, Jun Lin, Xiaozhong Liu, and Siliang Tang. 2023. Self-supervised meta-prompt learning with meta-gradient regularization for few-shot generalization. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	1040
		1041
		1042
		1043
		1044
		1045
	Xiao Pu, Mingqi Gao, and Xiaojun Wan. 2023. <a href="#">Summarization is (almost) dead</a> .	1046
		1047
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	1048
		1049
		1050
		1051
	Matthew Shardlow. 2013. <a href="#">A comparison of techniques to automatically identify complex words</a> . In <i>51st Annual Meeting of the Association for Computational Linguistics Proceedings of the Student Research Workshop</i> , pages 103–109, Sofia, Bulgaria. Association for Computational Linguistics.	1052
		1053
		1054
		1055
		1056
		1057
	Matthew Shardlow, Richard Evans, Gustavo Henrique Paetzold, and Marcos Zampieri. 2021. <a href="#">SemEval-2021 task 1: Lexical complexity prediction</a> . In <i>Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)</i> , pages 1–16, Online. Association for Computational Linguistics.	1058
		1059
		1060
		1061
		1062
		1063
	Sanjay S.P, Anand Kumar M, and Soman K P. 2016. <a href="#">AmritaCEN at SemEval-2016 task 11: Complex word identification using word embedding</a> . In <i>Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)</i> , pages 1022–1027, San Diego, California. Association for Computational Linguistics.	1064
		1065
		1066
		1067
		1068
		1069
		1070
	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>Transactions on Machine Learning Research</i> .	1071
		1072
		1073
		1074
		1075
		1076
		1077
	Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023. <a href="#">Text classification via large language models</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8990–9005, Singapore. Association for Computational Linguistics.	1078
		1079
		1080
		1081
		1082
		1083
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay	1084
		1085

1086	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	George-Eduard Zaharia, Răzvan-Alexandru Smădu, Du-	1145
1087	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	mitru Cercel, and Mihai Dascalu. 2022. <a href="#">Domain</a>	1146
1088	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	<a href="#">adaptation in multilingual and multi-domain mono-</a>	1147
1089	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	<a href="#">lingual settings for complex word identification</a> . In	1148
1090	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	<i>Proceedings of the 60th Annual Meeting of the Associ-</i>	1149
1091	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	<i>ation for Computational Linguistics (Volume 1: Long</i>	1150
1092	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	<i>Papers)</i> , pages 70–80, Dublin, Ireland. Association	1151
1093	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	for Computational Linguistics.	1152
1094	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	Jenny Alexandra Ortiz Zambrano, César Espin-Riofrio,	1153
1095	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	and Arturo Montejo-Ráez. 2023. <a href="#">Legalec: A new</a>	1154
1096	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	<a href="#">corpus for complex word identification research in</a>	1155
1097	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	<a href="#">law studies in ecuatorian spanish</a> . <i>Proces. del Leng.</i>	1156
1098	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	<i>Natural</i> , 71:247–259.	1157
1099	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	1158
1100	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	1159
1101	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024.	1160
1102	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	Judging llm-as-a-judge with mt-bench and chatbot	1161
1103	Melanie Kambadur, Sharan Narang, Aurelien Rod-	arena. <i>Advances in Neural Information Processing</i>	1162
1104	riguez, Robert Stojnic, Sergey Edunov, and Thomas	<i>Systems</i> , 36.	1163
1105	Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-</a>	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han,	1164
1106	<a href="#">tuned chat models</a> .	Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy	1165
1107	Gayatri Venugopal, Dhanya Pramod, and Ravi Shekhar.	Ba. 2023. <a href="#">Large language models are human-level</a>	1166
1108	2022. <a href="#">CWID-hi: A dataset for complex word iden-</a>	<a href="#">prompt engineers</a> . In <i>The Eleventh International</i>	1167
1109	<a href="#">tification in Hindi text</a> . In <i>Proceedings of the Thir-</i>	<i>Conference on Learning Representations, ICLR 2023,</i>	1168
1110	<i>teenth Language Resources and Evaluation Confer-</i>	<i>Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1169
1111	<i>ence</i> , pages 5627–5636, Marseille, France. European		
1112	Language Resources Association.		
1113	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin		
1114	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-		
1115	drew M. Dai, and Quoc V. Le. 2022a. <a href="#">Finetuned</a>		
1116	<a href="#">language models are zero-shot learners</a> . In <i>The Tenth</i>		
1117	<i>International Conference on Learning Representa-</i>		
1118	<i>tions, ICLR 2022, Virtual Event, April 25-29, 2022</i> .		
1119	OpenReview.net.		
1120	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten		
1121	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,		
1122	and Denny Zhou. 2022b. <a href="#">Chain-of-thought prompt-</a>		
1123	<a href="#">ing elicits reasoning in large language models</a> . In		
1124	<i>Advances in Neural Information Processing Systems</i>		
1125	<i>35: Annual Conference on Neural Information Pro-</i>		
1126	<i>cessing Systems 2022, NeurIPS 2022, New Orleans,</i>		
1127	<i>LA, USA, November 28 - December 9, 2022</i> .		
1128	Brandon T Willard and Rémi Louf. 2023. Effi-		
1129	cient guided generation for llms. <i>arXiv preprint</i>		
1130	<i>arXiv:2307.09702</i> .		
1131	Seid Muhie Yimam, Chris Biemann, Shervin Malmasi,		
1132	Gustavo Paetzold, Lucia Specia, Sanja Štajner, Anaïs		
1133	Tack, and Marcos Zampieri. 2018. <a href="#">A report on the</a>		
1134	<a href="#">complex word identification shared task 2018</a> . In <i>Pro-</i>		
1135	<i>ceedings of the Thirteenth Workshop on Innovative</i>		
1136	<i>Use of NLP for Building Educational Applications</i> ,		
1137	pages 66–78, New Orleans, Louisiana. Association		
1138	for Computational Linguistics.		
1139	G. Zaharia, D. Cercel, and M. Dascalu. 2020. <a href="#">Cross-</a>		
1140	<a href="#">lingual transfer learning for complex word identifica-</a>		
1141	<a href="#">tion</a> . In <i>2020 IEEE 32nd International Conference</i>		
1142	<i>on Tools with Artificial Intelligence (ICTAI)</i> , pages		
1143	384–390, Los Alamitos, CA, USA. IEEE Computer		
1144	Society.		

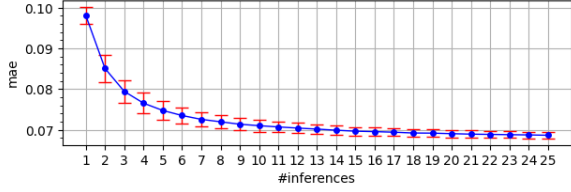
1170	<b>A Prompts</b>			
1171	<b>A.0.1 LCP English Prompt</b>		<b>A.0.3 CWI German Prompt</b>	1217
1172	You are a helpful, honest, and respect-		Sie sind ein hilfsbereiter, ehrlicher und	1218
1173	ful assistant for identifying the word		respektvoller Assistent, um die Wortkom-	1219
1174	complexity for beginner English learners.		plexität für Anfänger im Deutschen zu	1220
1175	You are given one sentence in English		identifizieren. Sie erhalten einen Satz	1221
1176	and a phrase from that sentence. Your		auf Deutsch und eine Phrase aus diesem	1222
1177	task is to evaluate the complexity of the		Satz. Ihre Aufgabe ist es zu sagen, ob	1223
1178	word. Answer with one of the following:		die Phrase komplex ist. Bewerten Sie	1224
1179	very easy, easy, neutral, difficult, very		die Antwort für die Phrase, anhand des	1225
1180	difficult. Be concise. Please, answer us-		Kontexts aus dem Satz. Seien Sie kurz.	1226
1181	ing the following JSON format:		Bitte verwenden Sie das folgende JSON-	1227
			Schema:	1228
1182	{		{	1229
1183	"sentence": "the sentence you were		"sentence": "der Satz, den Sie erhalten	1230
1184	provided",		haben",	1231
1185	"word": "the word or words you have		"word": "das Wort oder die Wörter,	1232
1186	to analyze",		die Sie analysieren müssen",	1233
1187	"proof": "explain your response in		"proof": "erklären Sie Ihre Antwort	1234
1188	maximum 50 words",		in maximal 50 Wörtern",	1235
1189	"complex": "either very easy, easy,		"complex": "entweder false (für	1236
1190	neutral, difficult, or very difficult"		einfach) oder true (für komplex)",	1237
1191	}		}	1238
1192	What is the difficulty of '{token}' from		Ist '{token}' von '{sentence}' com-	1239
1193	'{sentence}'?		plex?	1240
1194	<b>A.0.2 CWI English Prompt</b>		<b>A.0.4 CWI Spanish Prompt</b>	1241
1195	You are a helpful, honest, and respect-		Eres un asistente útil, honesto y respetu-	1242
1196	ful assistant for identifying the words		oso para identificar la complejidad de	1243
1197	complexity for beginner English learners.		las palabras para los principiantes que	1244
1198	You are given one sentence in English		aprenden español. Se te da una oración	1245
1199	and a phrase from that sentence. Your		en español y una frase de esa oración.	1246
1200	task is to say whether the phrase is com-		Tu tarea es decir si la frase es compleja.	1247
1201	plex. Assess the answer for the phrase,		Evalúa la respuesta para la frase, dada el	1248
1202	given the context from the sentence. Be		contexto de la oración. Sé conciso. Por	1249
1203	concise. Please, use the following JSON		favor, usa el siguiente esquema JSON:	1250
1204	schema:			
1205	{		{	1251
1206	"sentence": "the sentence you were		"sentence": "la oración que se te	1252
1207	provided",		proporcionó",	1253
1208	"word": "the word or words you have		"word": "la palabra o palabras que	1254
1209	to analyze",		tienes que analizar",	1255
1210	"proof": "explain your response in		"proof": "explica tu respuesta en	1256
1211	maximum 50 words",		máximo 50 palabras",	1257
1212	"complex": "either false (for simple)		"complex": "false (para simple) o true	1258
1213	or true (for complex)",		(para complejo)"	1259
1214	}		}	1260
1215	Is '{token}' complex in		¿Es '{token}' complejo en	1261
1216	'{sentence}'?		'{sentence}'?	1262

1263	<b>A.1 Fine-Tune Prompts</b>	<b>B Choice for Number of Inference Steps</b>	1310
1264	<b>A.1.1 LCP English Prompt</b>	As presented in Section 3, the estimated score in the LCP setting was an average of scores obtained after $N$ inference steps. We wanted to know what is the minimum number of inference steps required until the results do not change significantly anymore. Therefore, we set $N = 25$ for Llama-2-7b-ft, $N = 20$ for Llama-2-13b-ft, and $N = 10$ for ChatGPT-3.5-turbo-ft, and then estimated the average score per number of iterations using bootstrapping, with 100 samples. The plots are shown in Figure 1. We obtained that at least 10 to 15 runs are required, after which the scores do not change significantly.	1311 1312 1313 1314 1315 1316 1317 1318 1319 1320 1321 1322
1265	You are a helpful, honest, and respectful assistant for identifying the word difficulty for non-native English speakers. You are given one sentence in English and a word from that sentence. Your task is to evaluate the difficulty of the word. Answer only with one of the following: very easy, easy, neutral, difficult, very difficult.		
1266			
1267			
1268			
1269			
1270			
1271	sentence: ‘{sentence}’		
1272	word: ‘{token}’		
1273			
1274			
1275			
1276	<b>A.1.2 CWI English Prompt</b>	<b>C Evaluation Protocol</b>	1323
1277	You are a helpful, honest, and respectful assistant for identifying the word complexity for non-native English speakers. You are given one sentence in English and a word from that sentence. Your task is to say whether a word is complex or not. Answer only with one of the following: yes, no.	To evaluate the LLMs effectively, we employ an approach that uses optimized inference servers and a generic way to interface with LLMs. The overall protocol is showcased in Figure 2. First, we load the dataset, and for every example, we apply the system, user, and assistant prompt templates. Only in the few-shot setting we apply the assistant prompt template on the few-shot examples. The final prompt is sent to the server, which processes the input and returns the LLM’s prediction. We send the queries in parallel to use batching and other optimizations, thus reducing the execution time. For local inference endpoints, we use HuggingFace Text Generation Inference (TGI) <sup>7</sup> for most LLMs available in HuggingFace and lmsys’ FastChat <sup>8</sup> with vllm (Kwon et al., 2023) integration for improved inference throughput. OpenAI’s models are evaluated using their endpoints. In the end, we aggregate the results and evaluate them against the ground truth labels. Depending on the task, we compute and report the metrics and perform analysis.	1324 1325 1326 1327 1328 1329 1330 1331 1332 1333 1334 1335 1336 1337 1338 1339 1340 1341 1342 1343 1344 1345 1346 1347 1348 1349 1350 1351 1352 1353
1278			
1279			
1280			
1281			
1282			
1283			
1284			
1285	sentence: ‘{sentence}’		
1286	word: ‘{token}’		
1287	<b>A.1.3 CWI German Prompt</b>		
1288	Du bist ein hilfsbereiter, ehrlicher und respektvoller Assistent für die Identifizierung der Wortkomplexität für nicht-deutsche Muttersprachler. Dir wird ein Satz auf Deutsch und ein Wort aus diesem Satz gegeben. Deine Aufgabe ist es zu sagen, ob ein Wort komplex ist oder nicht. Antworten nur mit einem der Folgenden: ja, nein.		
1289			
1290			
1291			
1292			
1293			
1294			
1295			
1296			
1297	Satz: ‘{sentence}’		
1298	Wort: ‘{token}’		
1299	<b>A.1.4 CWI Spanish Prompt</b>		
1300	Eres un asistente útil, honesto y respetuoso para identificar la complejidad de las palabras para hablantes no nativos de inglés. Se te da una oración en inglés y una palabra de esa oración. Tu tarea es decir si una palabra es compleja o no. Responde solo con una de las siguientes opciones: sí, no.		
1301			
1302			
1303			
1304			
1305			
1306			
1307			
1308	oracion: ‘{sentence}’		
1309	palabra: ‘{token}’		

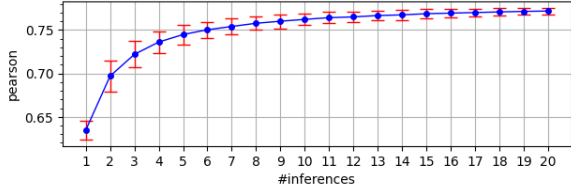
<sup>7</sup><https://huggingface.co/docs/text-generation-inference/en/index>  
<sup>8</sup><https://github.com/lm-sys/FastChat>  
<sup>9</sup><https://github.com/guidance-ai/guidance>



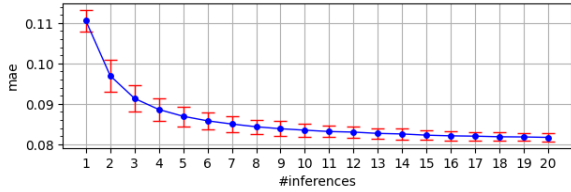
(a) Pearson score on Llama-2-7b-ft



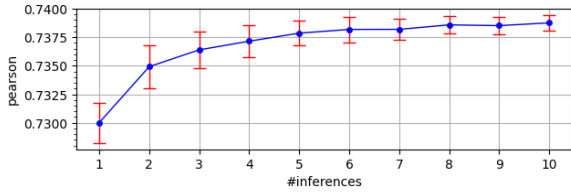
(b) MAE score on Llama-2-7b-ft



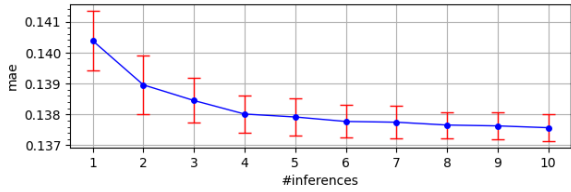
(c) Pearson score on Llama-2-13b-ft



(d) MAE score on Llama-2-13b-ft



(e) Pearson score on ChatGPT-3.5-turbo-ft



(f) MAE score on ChatGPT-3.5-turbo-ft

Figure 1: Estimated Pearson and MAE scores against the number of LLM inference steps.

## D Computational Costs

We trained and ran inferences on NVIDIA RTX 4080 and 4090 (consumer-class GPUs) and NVIDIA RTX A6000, A100 40GB-PCIe, and H100 80GB-SXM (server-class GPUs), depending on the minimal requirements to run the model

and execution time. The NVIDIA RTX 4090 (approx. \$0.5/hr), A6000 (approx. \$0.9/hr), and H100 80GB-SXM (approx. \$3.2/hr) were rented from vast.ai.

For OpenAI’s API, we used inference endpoints for ChatGPT-3.5-turbo and GPT-4o as well as training endpoints for ChatGPT-3.5-turbo, with the pricing at the time of writing this paper: \$0.0005 per 1k input tokens and \$0.0015 per 1k output tokens for ChatGPT-3.5-turbo; and \$0.0080 per training tokens, \$0.003 per 1k input tokens, and \$0.006 per 1k output tokens). For GPT-4o, we had access only to inference, with \$5 per 1M input tokens and \$15 per 1M output tokens. All experiments related to OpenAI’s models totaled about \$300. Because of the high costs involved, we limited our experiments to only classification on large test sets.

## E Model Checkpoints

In Table 7, we present the checkpoints used in this work. We indicate with “ft” where we use a different checkpoint for fine-tuning.

Model	Checkpoint
Llama 2 7B	meta-llama/Llama-2-7b-chat-hf
Llama 2 7B ft	meta-llama/Llama-2-7b-hf
Llama 2 13B	meta-llama/Llama-2-13b-chat-hf
Llama 2 13B ft	meta-llama/Llama-2-13b-hf
Vicuna 1.5 7B	lmsys/vicuna-7b-v1.5
Vicuna 1.5 7B AWQ	TheBloke/vicuna-7B-v1.5-AWQ
Vicuna 1.5 13B	lmsys/vicuna-13b-v1.5
Vicuna 1.5 13B AWQ	TheBloke/vicuna-13B-v1.5-AWQ
Llama 3 8B	meta-llama/Meta-Llama-3-8B-Instruct
Llama 3 8B ft	meta-llama/Meta-Llama-3-8B
ChatGPT-3.5-turbo	gpt-3.5-turbo-0125
GPT-4o	gpt-4o-2024-05-13

Table 7: Checkpoints used during experiments.

## F Few-Shot Examples and Proofs

The few-shot examples for CWI datasets (see Tables 8, 9, 10, 11, and 12) were chosen such that we provide two samples for false complexity (probability is 0%) and one sample for every discrete label as presented in Section 3.2. This way, even if we have a bias towards complex sentences, the distribution among discrete labels is uniform. The sentences were randomly chosen such that they fulfilled the previous criteria. The proofs are generated by GPT-4o by asking it what would be the reason, given the sentence and token, to have a particular label.



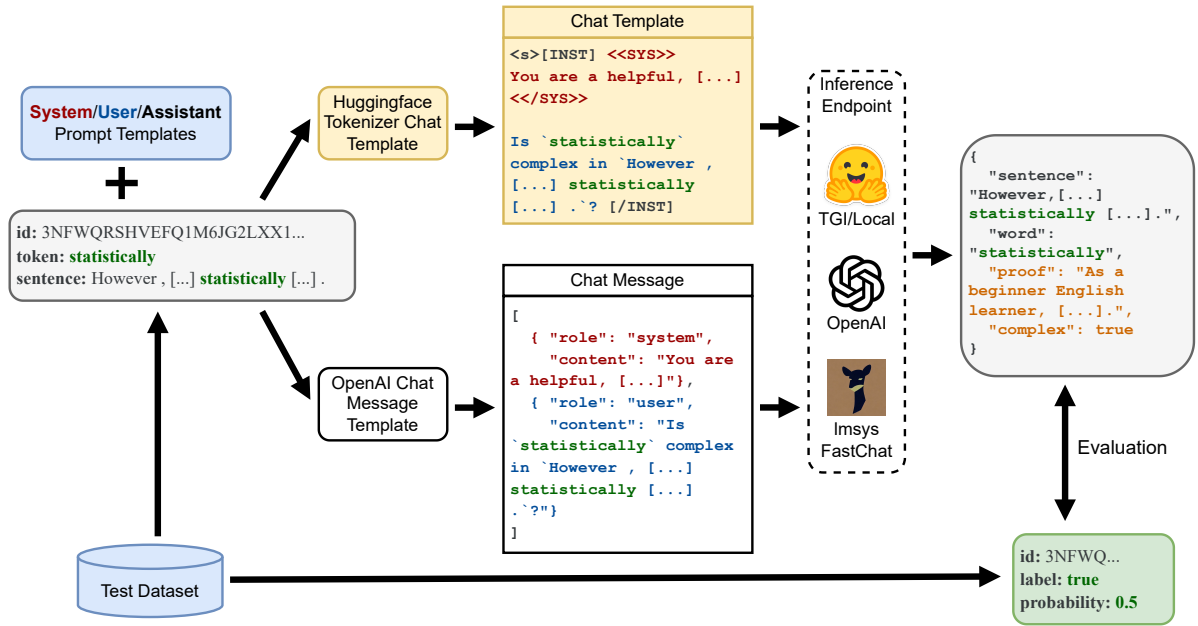


Figure 2: Evaluation protocol.

1393 Similarly, for LCP datasets (see Tables 13, 14), we  
 1394 sample an example from every category. In both  
 1395 tasks, the few-shot samples were selected from the  
 1396 train set. When generating the prompt, we shuffle  
 1397 the few-shot examples to reduce the chances of a  
 1398 position bias; that is, the model would predict a  
 1399 certain label due to the way the few-shot samples  
 1400 are ordered.

## 1401 G Meta-Learning Datasets

1402 We select 45 tasks from the BIG-bench benchmark,  
 1403 all of which being classification tasks. Some tasks  
 1404 offer the available choices in the original prompt,  
 1405 whereas others do not. For the ones that don't, we  
 1406 manually append them in the prompt. The tasks  
 1407 can be viewed in Table 15.

## 1408 H Supporting Figures for Discussions

1409 In this section, we include the figures that were  
 1410 used during discussions.

<b>ID</b>	<b>Sentence</b>	<b>Token</b>	<b>Complexity</b>	<b>Proof</b>
7329	Northern Ireland’s deputy first minister and Mid-Ulster MP, Martin McGuinness, said his heart went out to the family of the girl tragically killed.	MP	False	The abbreviation ‘MP’ for Member of Parliament is commonly understood in contexts involving government or politics, making it familiar and not complex.
7700	Activists said at least 30 people died on Sunday - mainly civilians - as violence surged at flashpoints across the country despite an increase of UN observers.	civilians	True	The word ‘civilians’ may be considered complex because it is specific to contexts involving military or emergency situations, and not everyone might be familiar with its precise meaning.
407	The regime is trying to punish these villages and to put an end to this revolution as quickly as it can, he said.	punish	False	The word ‘punish’ is a common verb used in everyday language to describe the act of imposing a penalty or suffering for an offense, making it familiar and not complex.
5182	The Philippines and Vietnam complained last year of increasingly aggressive acts by China in staking its claim to the South China Sea.	aggressive acts	True	The phrase ‘aggressive acts’ may be considered complex as it involves understanding both ‘aggressive’ and ‘acts’ together, implying a specific type of behavior which might not be immediately clear without context.
12472	Goodyear said police previously responded to the Florida ‘Space Coast’ home, about 15 miles south of Cape Canaveral, for domestic disturbance calls involving Jaxs Johnson.	domestic disturbance	True	The term ‘domestic disturbance’ can be complex as it combines ‘domestic’, related to the home or family, with ‘disturbance’, indicating trouble or conflict, requiring an understanding of both terms in context.
7131	Spain is set to intensify the clean-up of its banks on Friday after difficult last-minute talks between the government and lenders on details of planned financial system reforms.	Friday	False	The word ‘Friday’ is a basic term indicating a day of the week, universally understood and not complex.
10459	The country’s leaders have to admit that there were numerous falsifications and rigging and the results do not reflect the will of the people, Gorbachev told Interfax, according to the AFP.	rigging	True	The word ‘rigging’ can be considered complex as it refers to the act of manipulating or tampering with something, often in a fraudulent way, which may not be a familiar concept to everyone.

Table 8: Few shot examples for CWI 2018 English - news domain

ID	Sentence	Token	Complexity	Proof
4055	#29-17 He joins 139 other Republican Party presidential candidates who have done likewise.	Party	false	The word 'Party' is common and widely understood in political contexts, making it familiar to both native and non-native speakers.
5461	#11-14 The experiments were funded by national research organizations in the United States and China and the government of Brazil.	national	false	The word 'national' is a basic adjective used to describe something related to a nation, and is commonly used in many contexts, making it easy for most speakers.
4911	#42-4 The team used Formica fusca, an ant species that can form thousand-strong colonies.	Formica fusca	true	The term 'Formica fusca' is a scientific name for a specific ant species, which is likely unfamiliar to most people outside of entomology or biological sciences.
3758	#22-5 According to doctors at Bethany Hospital, Kalam was dead by 7 p.m. but they waited for the arrival of Meghalaya chief minister V. Shanmuganathan, about an hour later, before announcing the death.	announcing	true	The word 'announcing' can be challenging due to its length, the presence of a silent letter, and the necessity to understand the appropriate context for its use.
2220	#36-16 Another had been to tether the nose cone to the car; Hunter-Reay mentioned renderings developed of a boomerang-like debris-deflector positioned in front of the driver.	tether	true	The word 'tether' is less commonly used and may not be familiar to many people, leading to difficulty in understanding its meaning and usage.
1951	#24-37 Furthermore, the data of radars at Maldives airports have also been analysed and shows no indication of the said flight", said Malaysian Transport Minister Hishamuddin Hussein.	analysed	true	The word 'analysed' can be difficult due to its British English spelling (with 's' instead of 'z'), which might confuse those more familiar with American English.
1498	#3-10 Pavlensky and Oksana were detained in December at Sheremetyevo airport for questioning, which went on for seven hours.	detained	true	The word 'detained' may be difficult due to its legal context and the less frequent use in everyday language, requiring a higher level of vocabulary knowledge.

Table 9: Few shot examples for CWI 2018 English - WikiNews domain

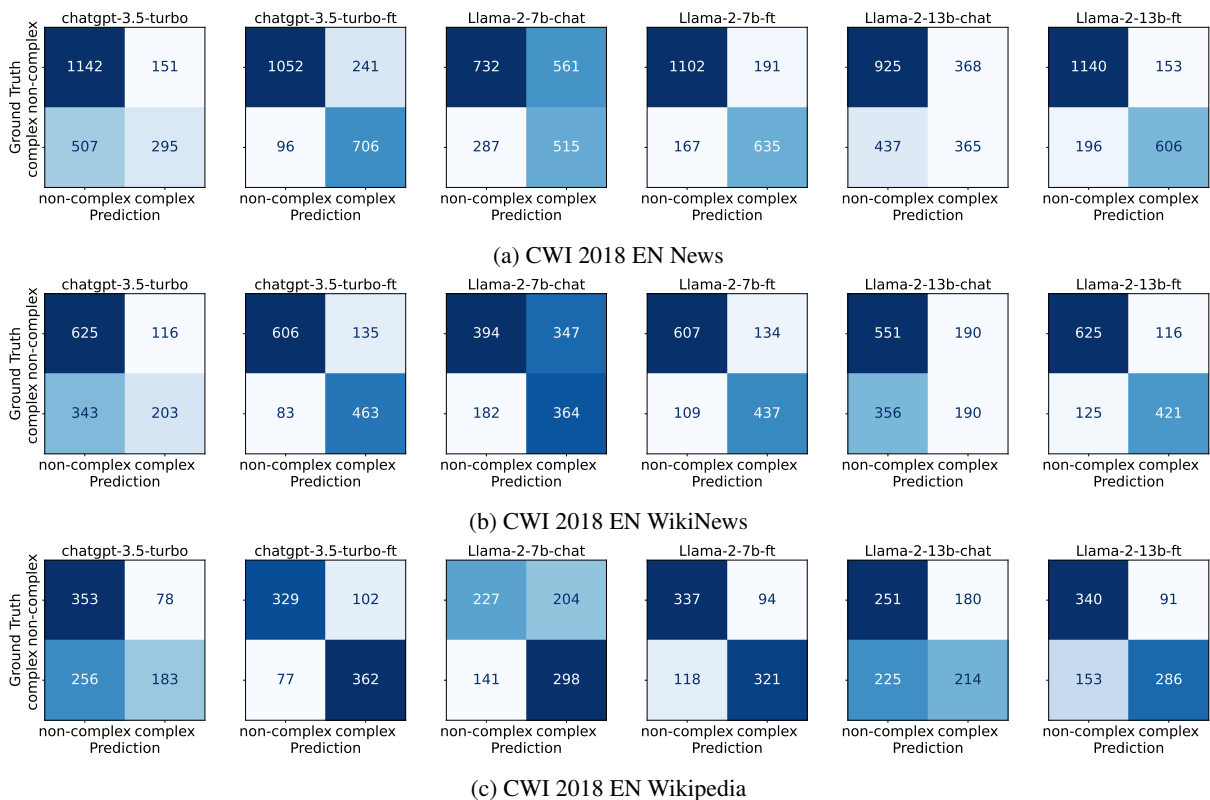


Figure 3: Confusion matrices computed on the English CWI datasets for News, WikiNews, and Wikipedia domains.

<b>ID</b>	<b>Sentence</b>	<b>Token</b>	<b>Complexity</b>	<b>Proof</b>
3595	Once the series had received the backing of the FIA , a management structure including new executive directors Brian Menell and Tony Teixeira were appointed to oversee the sale of franchises for the operation of international teams .	Brian	false	The word 'Brian' is a common proper noun and a typical English name, which is familiar to both native and non-native speakers. Its presence in the sentence is straightforward and does not add complexity.
3400	The first recorded case of an actor performing took place in 534 BC ( though the changes in calendar over the years make it hard to determine exactly ) when the Greek performer Thespis stepped on to the stage at the Theatre Dionysus and became the first known person to speak words as a character in a play or story .	play	false	The word 'play' is a basic English word frequently used in both its noun and verb forms. It is easily understood by both native and non-native speakers, especially in the context of theater.
1048	Also , if the reviewing administrator concludes that the block was justified , you will not be unblocked unless the reviewing administrator is convinced that you understand what you are blocked for , and that you will not do it again .	administrator	true	The word 'administrator' is long and contains multiple syllables, which can make it challenging to pronounce and remember. Additionally, its specific meaning in the context of authority or management may not be immediately clear to non-native speakers.
3670	Two is the base of the simplest numeral system in which natural numbers can be written concisely , being the length of the number a logarithm of the value of the number ( whereas in base 1 the length of the number is the value of the number itself ) ; the binary system is used in computers .	numeral	true	The word 'numeral' is less commonly used in everyday language and pertains to a specific field (mathematics). This specialization can make it less familiar and harder to understand for some readers.
2767	The Angara rocket family is a family of space-launch vehicles being developed by the Moscow-based Khrunichev State Research and Production Space Center .	space-launch	true	The term 'space-launch' is a compound word that refers to a specific and technical concept related to aerospace. Its specialized nature and the combination of two words can make it more difficult to understand.
1155	Early references from the Vadstena Abbey show how the Swedish nuns were baking gingerbread to ease indigestion in 1444 .	indigestion	true	The word 'indigestion' is relatively long and describes a specific medical condition related to digestion, which might not be commonly known or used in daily conversation, making it harder for some readers.
919	The roof of the nave is composed of a pair and knuckle frame , coated inside by pieces of tracery .	tracery	true	The word 'tracery' is an architectural term that may not be widely recognized outside of specialized contexts. Its specific meaning and less frequent use contribute to its complexity.

Table 10: Few shot examples for CWI 2018 English - Wikipedia domain

ID	Sentence	Token	Complexity	Proof	Proof (En.)
4890	Unmittelbar nach den Anschlügen vom 11.	Unmittelbar	false	Das Wort 'Unmittelbar' ist nicht komplex, da es ein häufig verwendetes deutsches Adjektiv ist und weder selten noch schwierig zu verstehen ist.	The word 'Immediate' is not complex as it is a commonly used German adjective and is neither rare nor difficult to understand.
713	Janukowytsch findet dort die größte Unterstützung , während Juschtschenko das größte Wählerpotenzial sieht .	größte	false	Das Wort 'größte' ist ein Basisadjektiv in der deutschen Sprache und stellt keine besondere Schwierigkeit dar.	The word 'largest' is a basic adjective in the German language and does not pose any particular difficulty.
4106	Sie berichtete unter anderem über ihre derzeitige Tournee mit dem Thema Hitler-Tagebücher .	Tournee	true	Das Wort 'Tournee' stammt aus dem Französischen und wird in der deutschen Sprache seltener verwendet, was es für Nicht-Muttersprachler schwieriger macht.	The word 'tournee' comes from French and is used less frequently in the German language, making it more difficult for non-native speakers.
2738	Die Anwälte Berlusconi kündigten an , gegen die Verjährung einen Einspruch einzureichen , um einen Freispruch erster Klasse zu erreichen .	Freispruch	true	Das Wort 'Freispruch' kann komplex sein, da es ein spezifischer juristischer Begriff ist, der in alltäglichen Gesprächen selten vorkommt.	The word 'acquittal' can be complex because it is a specific legal term that rarely appears in everyday conversations.
3535	Der eineinhalbstündige feierliche Trauergottesdienst fand in der zu zwei Drittel gefüllten Friedenskirche im Nürnberger Stadtteil St. -Johannis statt .	Trauer-gottesdienst	true	Das Wort 'Trauergottesdienst' ist komplex, da es ein zusammengesetztes Substantiv ist und selten verwendet wird.	The word 'funeral service' is complex because it is a compound noun and is rarely used.
185	Konvergenz als Ursache der Fehleinordnung : Nach ihrer Analyse des Fibrinogen-Gens stellen etwa die äußerlich sehr ähnlichen Flamingos und Löffler zwei weit auseinanderliegende Gruppen auf den beiden Evolutionsästen dar .	Fibrinogen-Gens	true	Das Wort 'Fibrinogen-Gens' ist komplex, da es ein wissenschaftlicher Begriff ist, der in der allgemeinen Sprache nicht häufig vorkommt.	The word 'fibrinogen gene' is complex because it is a scientific term that is not commonly used in common language.
5726	Hauptgrund für die Verschlechterung des Zustandes sei der heiße und trockene Sommer 2003 mit hohen Ozonwerten .	Ozonwerten	true	Das Wort 'Ozonwerten' kann für Nicht-Muttersprachler schwierig sein, da es ein wissenschaftlicher Begriff ist und spezifisches Wissen über Luftqualität erfordert.	The word 'ozone levels' can be difficult for non-native speakers as it is a scientific term and requires specific knowledge of air quality.

Table 11: Few shot examples for CWI 2018 German. For proofs we also provide the translation in English.

ID	Sentence	Token	Complexity	Proof	Proof (En.)
11798	En 1911, escapó de su casa y se alistó en una expedición militar, organizada por Ricciotti Garibaldi, para liberar a Albania del control turco.	Garibaldi	false	El apellido 'Garibaldi' no es difícil porque es un nombre propio conocido, especialmente en el contexto de la historia y la cultura italiana.	The surname 'Garibaldi' is not difficult because it is a well-known proper name, especially in the context of Italian history and culture.
10963	Estos magos fueron, según la tradición, adorar al Mesías que acababa de nacer en Belén de Judea, el que posteriormente se llamaría Jesús de Nazaret.	adorar	true	La palabra 'adorar' puede considerarse difícil debido a su uso menos común y su connotación religiosa específica.	The word 'worship' may be considered difficult due to its less common use and its specific religious connotation.
8294	En marzo de 2011 firma con el BK Jimki dónde sustituirá a Meleschenko, entrenador interino desde la renuncia de Sergio Scariolo tras no conseguir el pase para el Top-16 de la Euroliga.	interino	true	La palabra 'interino' puede ser difícil debido a su uso en un contexto específico y profesional, lo que requiere un conocimiento preciso del término.	The word 'interim' can be difficult due to its use in a specific and professional context, which requires precise knowledge of the term.
6171	Linda con las poblaciones de Yepes, Huerta de Valdecarábanos y el término segregado de La Guardia, todas de Toledo.	Linda	true	La palabra 'Linda' es difícil porque se trata de un término geográfico específico que puede no ser conocido por todos los hablantes.	The word 'Linda' is difficult because it is a specific geographical term that may not be known to all speakers.
5911	Estuvieron presentes el presidente de Estados Unidos Bill Clinton y el presidente de la República de Corea Kim Young Sam, y se dedicó a los hombres y mujeres que sirvieron en la guerra.	Bill	false	El nombre 'Bill' no es difícil porque es un nombre propio común y fácil de reconocer, especialmente en el contexto de figuras públicas como Bill Clinton.	The name 'Bill' is not difficult because it is a common and easy to recognize proper name, especially in the context of public figures like Bill Clinton.
2673	Cada uno de los vectores columna de la matriz "A" se llama modo propio de vibración, y los "Ci" son las amplitudes relativas de cada modo propio.	amplitudes	true	La palabra 'amplitudes' es técnica y específica del campo de las matemáticas y la física, lo que puede hacerla difícil para quienes no están familiarizados con estos temas.	The word 'amplitudes' is technical and specific to the field of mathematics and physics, which can make it difficult for those unfamiliar with these topics.
1945	El Ducado de Prusia o Prusia Ducal (en alemán: "Herzogtum Preußen"; en polaco: "Prusy Książęce") fue un ducado entre 1525-1701 en la región más oriental de Prusia heredero del Estado monástico de los Caballeros Teutónicos.	monástico	true	La palabra 'monástico' es difícil porque es un término especializado que se refiere a la vida y organización de los monasterios, lo que puede no ser familiar para todos.	The word 'monastic' is difficult because it is a specialized term referring to the life and organization of monasteries, which may not be familiar to everyone.

Table 12: Few shot examples for CWI 2018 Spanish. For proofs we also provide the translation in English.

ID	Sentence	Token	Complexity	Proof
6043	Containers lost at sea and compensation (debate)	Containers	Very Easy	The word 'Containers' is a common and easily understood term in English, referring to objects used for holding or transporting items.
4290	We have also shown that chondrogenesis can be initiated and chondrogenic differentiation will take place even in the absence of both BMP2 and BMP4 or BMP2 and BMP7.	differentiation	Easy	The word 'differentiation' is slightly technical and commonly used in biological contexts, making it easy but not very easy.
2143	Their scribes and the Pharisees murmured against his disciples, saying, "Why do you eat and drink with the tax collectors and sinners?"	scribes	Neutral	The term 'scribes' is not commonly used in everyday language and refers to a specific historical role, requiring some background knowledge to understand.
5144	Our data suggest that while recombination events destined to be resolved as COs can proceed normally in Trip13 mutants, DSBs that enter the NCO repair pathway are incompletely resolved or processed inefficiently.	COs	Difficult	The acronym 'COs' is specialized and requires specific knowledge in genetics to understand that it refers to 'crossovers' in the context of recombination events.
4873	In the mouse model of RA, small genetic contributions are also often observed.	RA	Very Difficult	The acronym 'RA' stands for 'rheumatoid arthritis,' a term that is highly specialized and not immediately clear without specific medical knowledge.

Table 13: Few shot examples for LCP 2021 single-word expressions

ID	Sentence	Token	Complexity	Proof
526	Therefore, $TGF\beta$ and BMP signaling are playing distinct but necessary roles to maintain articular cartilage.	necessary roles	Very Easy	The phrase 'necessary roles' is straightforward, commonly used in English, and easily understood within the context of the sentence.
212	In this confidence, I was determined to come first to you, that you might have a second benefit;	second benefit	Easy	The phrase 'second benefit' is relatively simple, but the context may slightly challenge the reader, making it less immediate to understand.
1376	We will be very strict on enforcing this fundamental principle in this case as well.	fundamental principle	Neutral	The term 'fundamental principle' requires a moderate understanding of abstract concepts and formal language, making it neutral in difficulty.
503	neither to pay attention to myths and endless genealogies, which cause disputes, rather than God's stewardship, which is in faith-	endless genealogies	Difficult	The phrase 'endless genealogies' is less common and refers to complex and potentially obscure biblical or historical references, adding to its difficulty.
1008	Such polymorphisms should yield biomarkers suitable for more readily accessible samples, such as peripheral blood or buccal smears.	buccal smears	Very Difficult	The term 'buccal smears' is highly specialized and technical, typically known only to those with specific biomedical knowledge, making it very difficult.

Table 14: Few shot examples for LCP 2021 multi-word expressions

Tasks		
abstract_narrative_understanding	fantasy_reasoning	nonsense_words_grammar
analytic_entailment	figure_of_speech_detection	odd_one_out
bbq_lite_json	formal_fallacies_syllogisms_negation	penguins_in_a_table
causal_judgment	general_knowledge	phrase_relatedness
cause_and_effect	human_organs_senses	play_dialog_same_or_different
codenames	hyperbaton	presuppositions_as_nli
contextual_parametric_knowledge_conflicts	implicatures	question_selection
crash_blossom	implicit_relations	reasoning_about_colored_objects
crass_ai	intent_recognition	riddle_sense
dark_humor_detection	irony_identification	ruin_names
disambiguation_qa	logical_deduction	strange_stories
empirical_judgments	logical_fallacy_detection	temporal_sequences
entailed_polarity	metaphor_boolean	timedial
epistemic_reasoning	metaphor_understanding	tracking_shuffled_objects
evaluating_information_essentiality	movie_dialog_same_or_different	winowhy

Table 15: All tasks selected from the BIG-bench benchmark that were used during the meta-learning process.

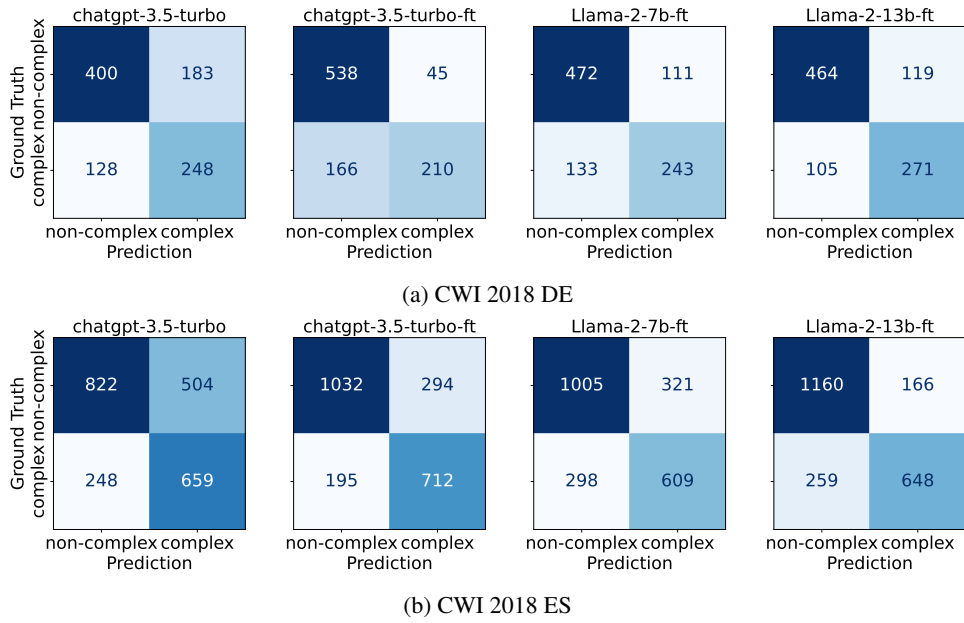


Figure 4: Confusion matrices computed on the German and Spanish CWI datasets.

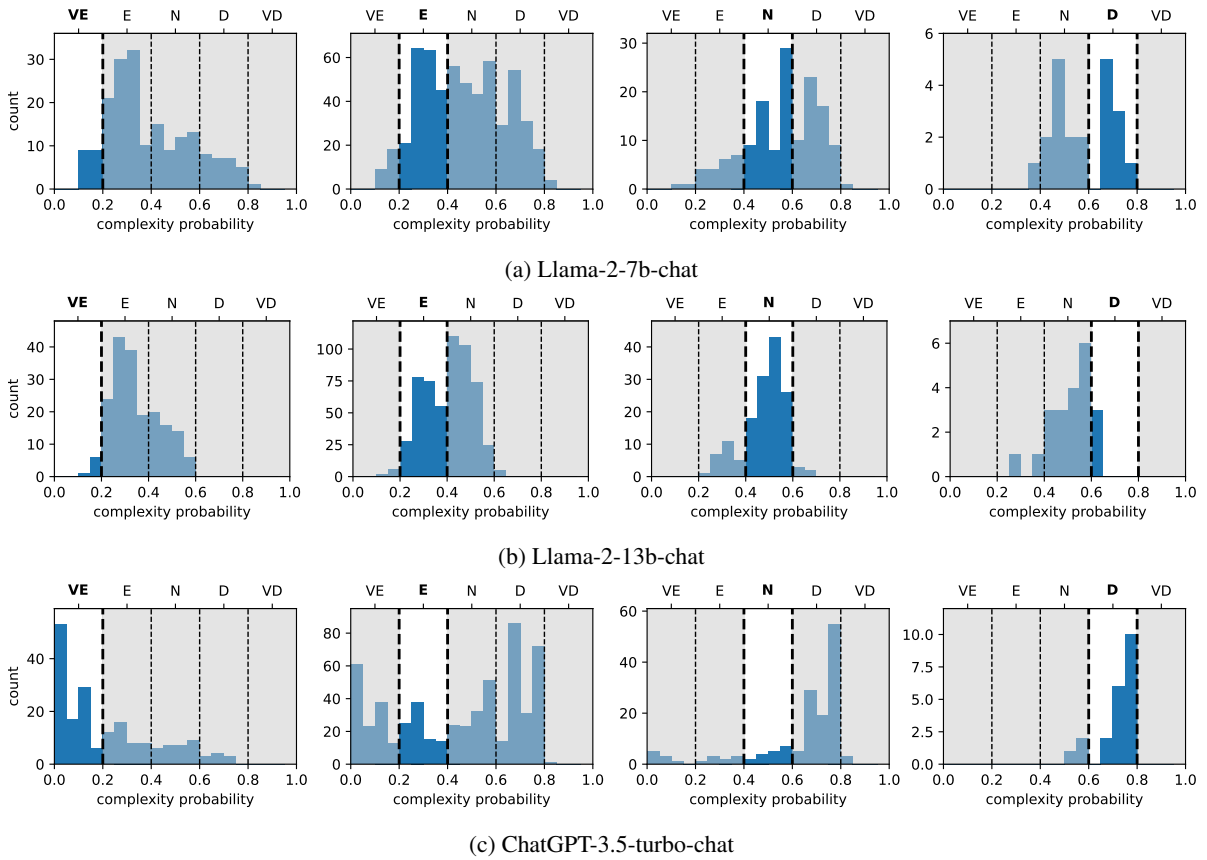


Figure 5: Predictive probability distribution of zero-shot LLMs on LCP single-word test set. Highlighted is the ground truth interval. Neither model predicts in the VD interval. VE – very easy, E – easy, N – neutral, D – difficult, VD – very difficult.



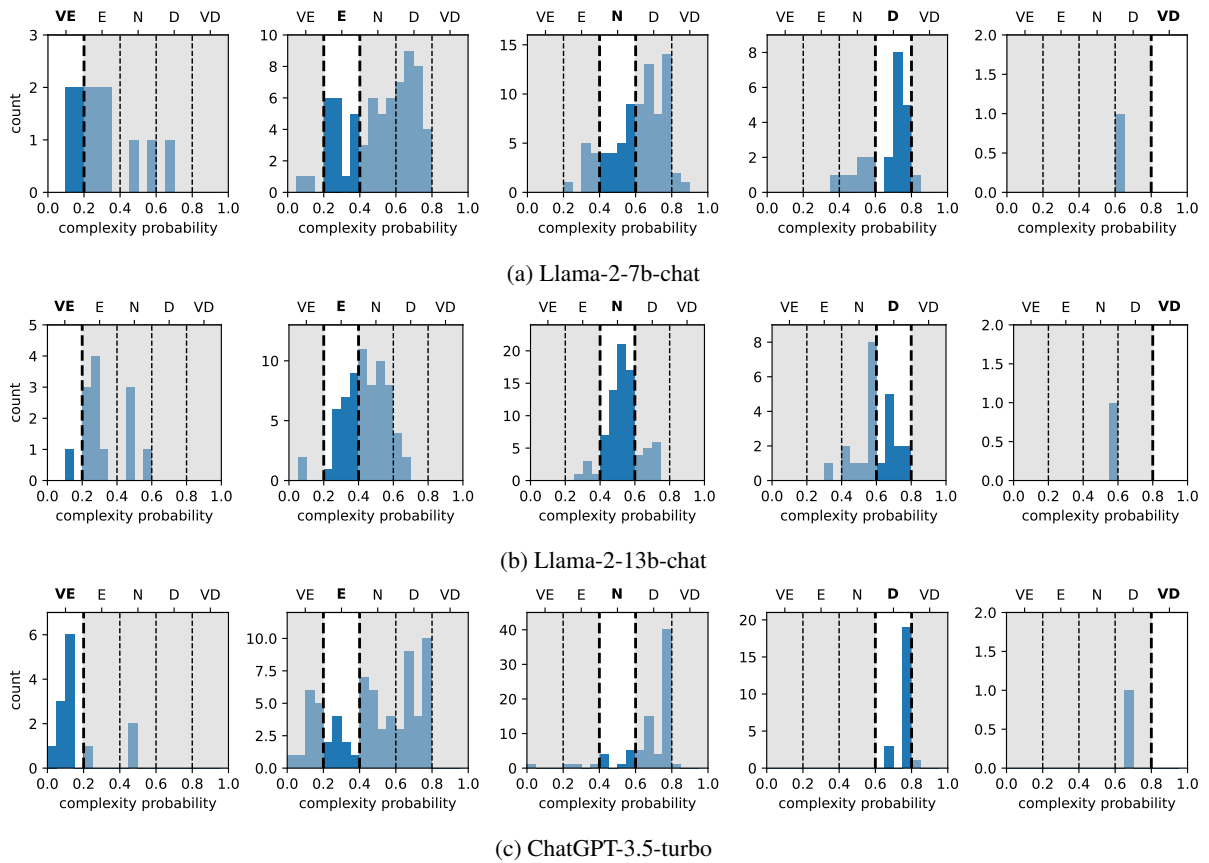


Figure 6: Predictive probability distribution of zero-shot LLMs on LCP multi-word test set. Highlighted is the ground truth interval. Neither model predicts in the VD interval. VE – very easy, E – easy, N – neutral, D – difficult, VD – very difficult.

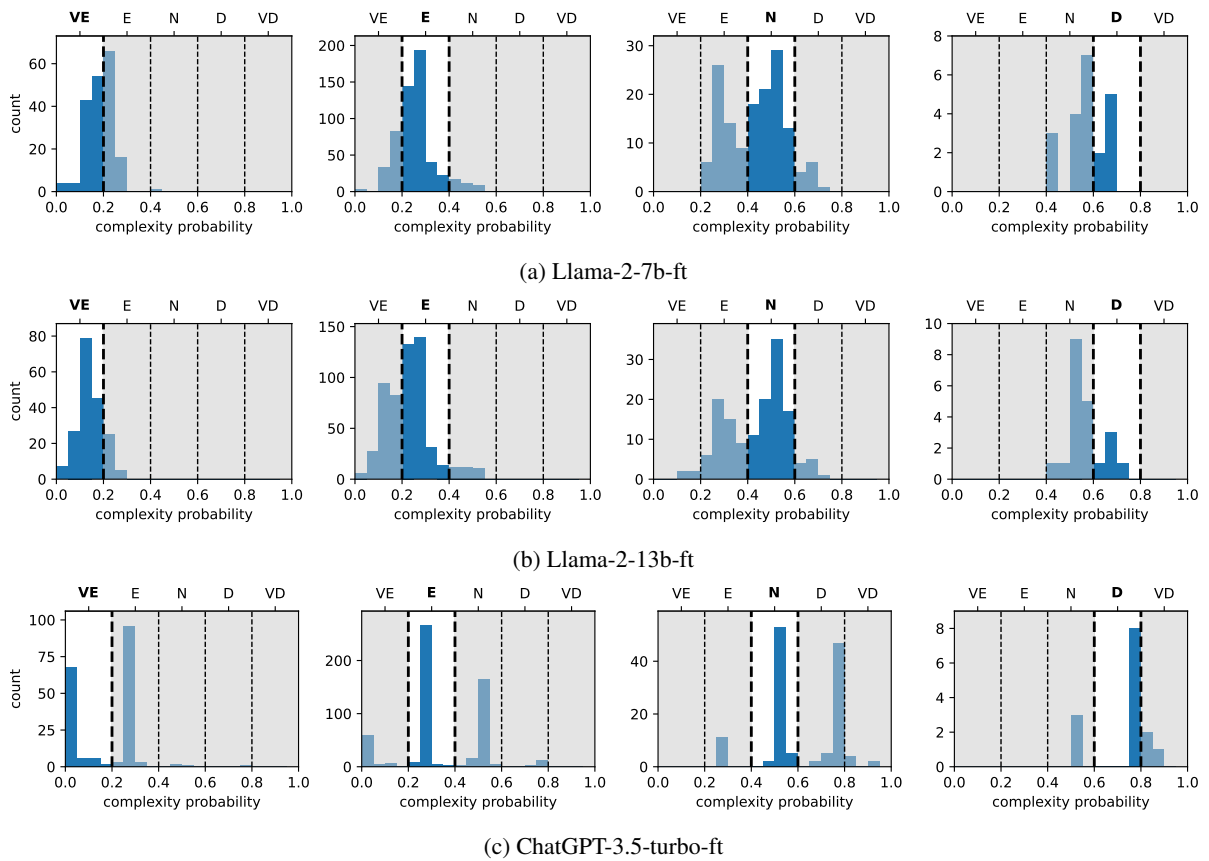


Figure 7: Predictive probability distribution of fine-tuned LLMs on LCP single-word test set. Highlighted is the ground truth interval. Neither model predicts in the VD interval. VE - very easy, E - easy, N - neutral, D - difficult, VD - very difficult.

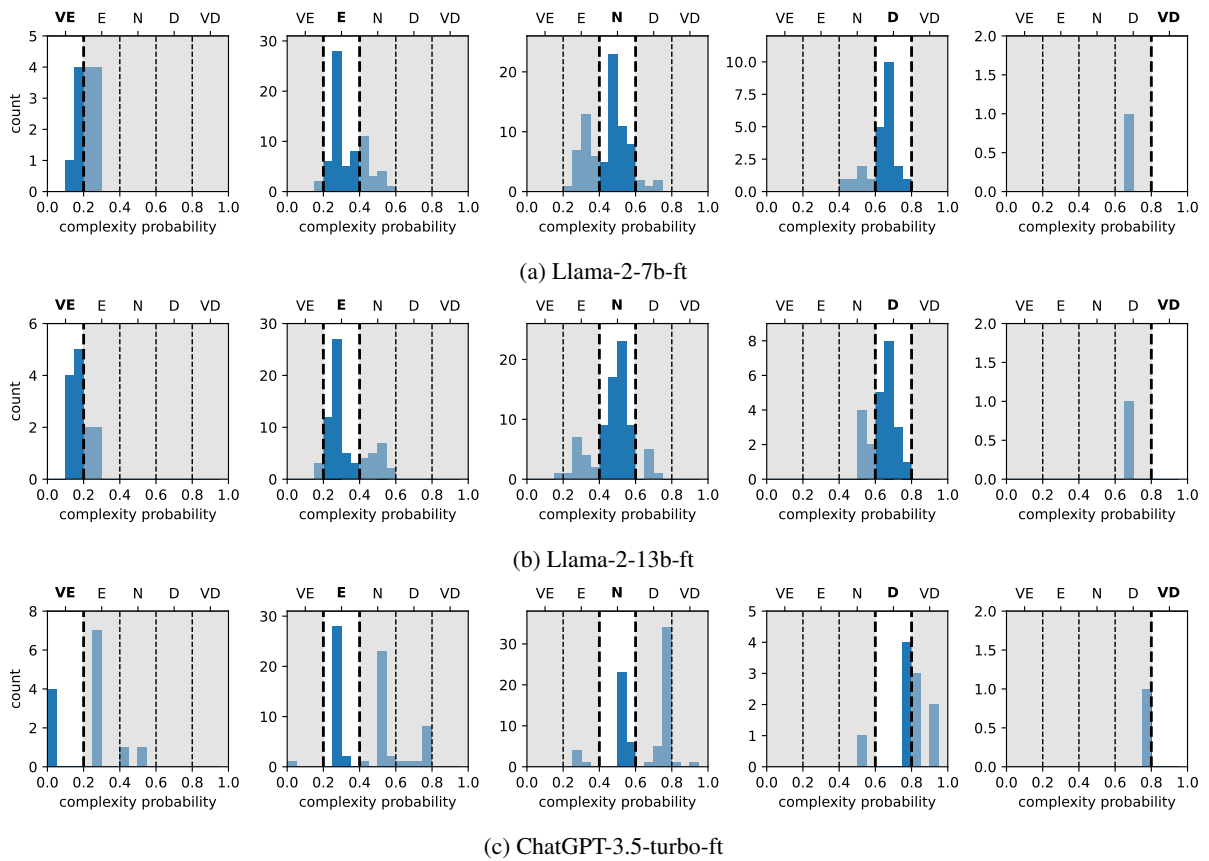


Figure 8: Predictive probability distribution of fine-tuned LLMs on LCP multi-word test set. Highlighted is the ground truth interval. Neither model predicts in the VD interval. VE - very easy, E - easy, N - neutral, D - difficult, VD - very difficult.

Model	En-W $\uparrow$	En-WN $\uparrow$	En-W $\uparrow$	De $\uparrow$	Es $\uparrow$
<i>Zero-shot</i>					
Llama-2-7b-chat	63.8	57.7	55.3	54.0	51.5
Llama-2-13b-chat	63.2	58.8	53.3	54.9	40.8
Vicuna-v1.5-7b	59.8	56.8	51.6	50.3	60.6
Vicuna-v1.5-13b	63.1	59.1	51.7	59.6	50.8
Llama-3-8b-chat	70.0	62.9	61.5	60.3	55.9
ChatGPT-3.5-turbo	69.5	64.6	62.0	60.1	63.3
GPT-4o	<b>76.8</b>	<b>73.2</b>	<b>71.0</b>	<b>73.0</b>	<b>75.8</b>
<i>Zero-shot CoT</i>					
Llama-2-7b-chat	56.2	59.1	57.2	45.7	46.0
Llama-2-13b-chat	61.6	56.5	54.7	53.2	49.7
Vicuna-v1.5-7b	60.3	58.7	56.9	50.0	58.3
Vicuna-v1.5-13b	66.8	61.8	55.9	58.9	58.7
Llama-3-8b-chat	66.9	63.0	64.4	59.9	55.5
ChatGPT-3.5-turbo	69.9	68.0	64.6	58.8	47.2
GPT-4o	<b>75.7</b>	<b>75.4</b>	<b>73.9</b>	<b>69.8</b>	<b>71.9</b>
<i>Few-shot</i>					
Llama-2-7b-chat	63.9	55.5	61.4	43.0	50.1
Llama-2-13b-chat	65.3	65.0	60.8	56.8	62.1
Vicuna-v1.5-7b	64.5	59.5	57.8	58.8	63.3
Vicuna-v1.5-13b	65.3	63.7	59.0	56.9	63.6
Llama-3-8b-chat	72.9	69.9	<b>67.6</b>	60.3	61.7
ChatGPT-3.5-turbo	72.3	66.4	65.1	65.5	72.7
GPT-4o	<b>76.6</b>	<b>71.8</b>	66.8	<b>75.9</b>	<b>75.7</b>
<i>Few-shot CoT</i>					
Llama-2-7b-chat	64.2	53.0	57.5	45.3	59.0
Llama-2-13b-chat	68.9	60.9	57.5	59.0	63.5
Vicuna-v1.5-7b	65.5	63.6	63.6	60.9	65.3
Vicuna-v1.5-13b	70.2	65.1	62.8	61.8	73.4
Llama-3-8b-chat	66.5	63.6	61.4	67.7	72.6
ChatGPT-3.5-turbo	72.4	66.0	65.2	72.9	<b>73.9</b>
GPT-4o	<b>77.2</b>	<b>68.9</b>	<b>68.5</b>	<b>73.8</b>	73.5
<i>Fine-tuned</i>					
Llama-2-7b-ft	82.9	81.1	76.7	<b>79.4</b>	75.7
Llama-2-13b-ft	83.3	81.3	74.6	76.6	<b>81.0</b>
Vicuna-v1.5-7b-ft	84.3	79.3	77.2	76.7	75.4
Vicuna-v1.5-13b-ft	85.2	80.3	<b>80.8</b>	79.0	74.3
Llama-3-8b-ft	<b>86.3</b>	82.9	77.5	78.1	76.9
ChatGPT-3.5-turbo-ft	83.9	<b>83.1</b>	79.4	78.0	74.4

Table 16: The accuracy on the test sets from CWI 2018 Shared Dataset. Notation: En - English dataset, De - German, and Es - Spanish; for English datasets, N - News domain, WN - WikiNews, W - Wikipedia. In bold, we denote the best score.

<b>Sentence</b>	<b>Word</b>	<b>Answer</b>	<b>Proof</b>	<b>Ground Truth</b>
Toronto traded Stewart to Chicago early in the 1947-48 season in a deal that brought Max Bentley to the Leafs.	brought	False	The verb 'brought' is not complex as it is a commonly used verb in everyday language with a straightforward meaning	True
His use of Russian language formed the basis of the style of novelists Ivan Turgenev, Ivan Goncharov, and Leo Tolstoy, as well as that of subsequent lyric poets such as Mikhail Lermontov.	basis	False	The word 'basis' is not particularly complex, but the concept it represents (i.e., the foundation or starting point of something) may be unfamiliar to some beginner English learners	True
Vertical distance measurements in the "down" direction are commonly referred to as depth.	depth	True	The word 'depth' has multiple meanings, including a vertical extent or height, making it potentially complex for beginning English learners to understand without proper context or explanation	True
The lack of oxygen above 2,400 meters (8,000 ft) can cause serious illnesses such as altitude sickness, high altitude pulmonary edema, and high altitude cerebral edema.	ft	True	The abbreviation 'ft' is commonly used in English to represent feet, which is a unit of measurement. However, in this context, it may be challenging for beginners to understand because they might not be familiar with the abbreviation.	False

Table 17: Examples of predictions and proofs for the Llama-2-13b-chat model on the CWI English Wikipedia dataset.