

Holmes

Benchmark the Linguistic Competence of Language Models

Anonymous ACL submission

Abstract

We introduce *Holmes*, a benchmark to assess the *linguistic competence* of language models (LMs) – their ability to grasp linguistic phenomena. Unlike prior prompting-based evaluations, *Holmes* assesses the linguistic competence of LMs via their internal representations using classifier-based probing. In doing so, we disentangle specific phenomena (e.g., part-of-speech of words) from other cognitive abilities, like following textual instructions, and meet recent calls to assess LMs’ linguistic competence in isolation. Composing *Holmes*, we review over 250 probing studies and feature more than 200 datasets to assess *syntax*, *morphology*, *semantics*, *reasoning*, and *discourse* phenomena. Analyzing over 50 LMs reveals that, aligned with known trends, their linguistic competence correlates with model size. However, surprisingly, model architecture and instruction tuning also significantly influence performance, particularly in *morphology* and *syntax*. Finally, we propose *FlashHolmes*, a streamlined version of *Holmes* designed to lower the high computation load while maintaining high-ranking precision.



[holmes-benchmark.github.io](https://github.com/holmes-benchmark)

1 Introduction

Linguistic competence is the unconscious understanding of language, like grasping grammatical rules (Chomsky, 1965). As language models (LMs) are trained on simple tasks like next word prediction (Brown et al., 2020), one might naturally wonder: *What is the linguistic competence of LMs, and how do they differ?* To answer such questions, benchmarks estimate cognitive abilities by providing textual instructions and evalu-

Holmes Rankings

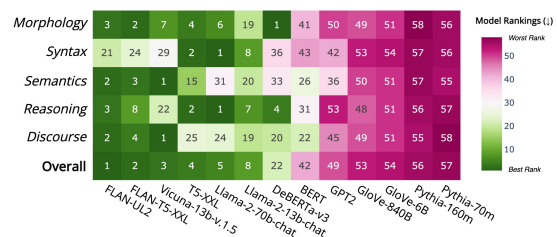


Figure 1: A subset of *Holmes* rankings (\downarrow) for various evaluated LMs. FLAN-UL2 outperforms the others *overall*, while different LMs prevail for the five distinct types of linguistic phenomena.

ate LMs’ responses, as done for mathematical reasoning (Cobbe et al., 2021) or factual knowledge (Petroni et al., 2019, 2020). However, they conflate latent abilities (like following provided instructions) with those under test, such as understanding specific linguistic phenomena, e.g., syntactic structures (Liang et al., 2023). As this entanglement makes it infeasible to draw definitive conclusions about distinct abilities (Hu and Levy, 2023), recent studies call to assess the linguistic competence of LMs comprehensively and in isolation (Lu et al., 2023; Mahowald et al., 2024).

In this work, we introduce the *Holmes* (Figure 2). A benchmark to assess the linguistic competence of LMs (Figure 1) regarding numerous linguistic phenomena. To fully disentangle the understanding of these phenomena and other abilities of LMs, we use classifier-based probing (Tenney et al., 2019a; Hewitt and Manning, 2019; Belinkov, 2022). A method that uses the LMs’ internal representations of text inputs to train linear models (probes) to predict specific aspects of phenomena, such as words’ part-of-speech (POS). We then ap-

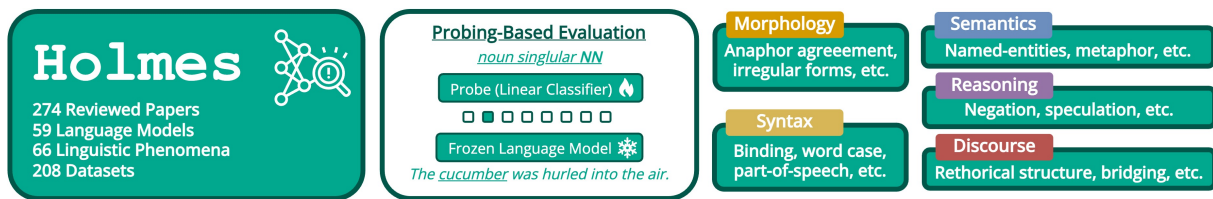


Figure 2: Overview of `Holmes` (left) with the five phenomena types (right) and an example of probing-based evaluations for part-of-speech: encoding the input tokens and predicting the POS tag for *cucumber*, here *NN*.

proximate the LMs’ grasp of these phenomena using the probes’ performance, rigorously verified using control tasks (Hewitt and Liang, 2019) and from an information theory perspective (Voita and Titov, 2020). With this particular and comprehensive scope, we thoroughly address the initially raised questions as follows:

Meta-Study (§ 3) The review of over 270 probing studies reveals a gap in comprehensively evaluating linguistic competence. Despite covering over 200 probing tasks and 150 LMs, individual studies focus on particular tasks and LMs. As a result, only three LMs were probed on over 20% of the tasks, and one single task was evaluated for more than 20% of the reviewed LMs. Notably, recent large LMs are significantly underrepresented.

Benchmark (§ 4) To address this identified deficiency, `Holmes` offers a structured framework to assess the English linguistic competence of LMs comprehensively. It features 208 distinct datasets covering *morphology*, *syntax*, *semantics*, *reasoning*, and *discourse* phenomena, including previously underrepresented ones like negation or rhetoric in text (Liang et al., 2023).

Results and Analysis (§ 5) From assessing 59 LMs (Figure 1), we find that no single one consistently excels the others and that their linguistic competence is more pronounced for *morphology* and *syntax* than the other phenomena types. Instead, we find **model size**, **model architecture**, and **instruction tuning** fundamentally affect their linguistic competence.

First, LMs’ linguistic competence, particularly for *morphology* and *syntax*, scales with their **model size**. This generalizes previous findings (Tenney et al., 2019b; Zhang et al., 2021) beyond LMs with 350 million parameters. Second, contrary to prompting evaluations (Lu et al., 2023) and aligned with other work (Waldis et al., 2024a; Gautam et al., 2024), **model architecture** is critical. The linguistic competence of decoder-only LMs is less pro-

nounced, and even 70 billion does not allow them to encode linguistic phenomena of words with comparable strength to encoder-only LMs of a similar size. Third, while previous studies focused on aligning LMs with human interactions through **instruction tuning** (Ouyang et al., 2022; Touvron et al., 2023; Zhou et al., 2023), we show for the first time its effect on their linguistic competence. It improves *morphology* and *syntax* but has mixed effects for the other types of phenomena. Lastly, we contrast the results of `Holmes` with OpenLLM (Beeching et al., 2023), an extensive LM benchmark focusing on user-centered applications like mathematical reasoning. We find that `Holmes` provides a unique but supplementary perspective, as rankings partly align, especially for reasoning-related phenomena.

Efficiency (§ 6) Finally, to mitigate the heavy computational burden of evaluating a new LM on `Holmes`, we form the streamlined version `FlashHolmes` by selectively excluding samples not significantly influencing overall rankings (Perlititz et al., 2023). Specifically, `FlashHolmes` approximates `Holmes` rankings with high precision while requiring only ~3% of the computation.

We summarize our contributions as follows:

- **Benchmark.** `Holmes` comprehensively and thoroughly assesses the linguistic competence of LMs in isolation, providing substantial ground for advancements in NLP.
- **Empirical insights.** Extensive experiments reveal that LMs’ linguistic competence is more pronounced for *morphology* and *syntax*, and size, architecture, and instruction tuning are crucial for LM differences.
- **Ease of use.** We provide tools to interactively explore `Holmes` results and straightforward code to evaluate upcoming LMs with efficiency in mind (`FlashHolmes`)¹.

¹Find resources at <https://holmes-benchmark.github.io>

2 Preliminaries

Language Models (LMs) Language Models compute probabilities for word sequences i , enabling tasks such as classifying i , textual comparisons between i and another sequence i' , and text generation based on i . We consider LMs as any model producing representations of input i , regardless of their specific type: **sparse** like bag-of-words (Harris, 1954); **static** such as GloVe (Pennington et al., 2014); or **contextualized** transformer-based LMs (Devlin et al., 2019; Raffel et al., 2020).

Linguistic Competence Following Chomsky (1965), linguistic competence is defined as the unconscious knowledge of language, encompassing the understanding of specific linguistic phenomena, including word dependencies and their distinct parts of speech (POS).

Linguistic Phenomena We define the linguistic competence of LMs as their ability to understand a diversity of linguistic phenomena. Specifically, we focus on five phenomena types: *morphology*, the structure of words; *syntax*, the structure of sentences; *semantics*, the meaning of words; *reasoning*, the use of words in logical deduction and other related phenomena like negation or speculation; *discourse*, the context in text like rhetorical structure. Following Mahowald et al. (2024), we categorize these phenomena types into two groups: *morphology* and *syntax* are **formal** phenomena, which include understanding grammatical rules and statistical patterns, while **functional** ones (*semantics*, *reasoning*, and *discourse*) focus on practical abilities like interpreting text sentiment or detecting the existence of speculation.

Datasets We define a dataset as text examples and labels covering a specific aspect of a linguistic phenomenon, like words and their POS tag. Typically, these labels are highly unambiguous to assess the specific aspect under test in isolation.

Probes Using probes, we empirically assess the linguistic competence of LMs regarding the featured linguistic phenomena in HOLMES. To this end, we employ probing tasks using the widely recognized classifier-based probing method (Tenney et al., 2019a; Hewitt and Manning, 2019; Belinkov, 2022), or known as diagnostic classifiers (Veldhoen et al., 2016; Giulianelli et al., 2018). Running such a probing task involves training a probe (linear model) using the specific dataset to

test a distinct aspect of a linguistic phenomenon in isolation. Therefore, we feed the text examples, encoded with a given LM, as training inputs. Subsequently, we use the probe’s performance to approximate how an LM understands the specific linguistic phenomenon under test. With a higher score, we assume the embeddings embody patterns relevant to this phenomenon, which enhances the accuracy (Tenney et al., 2019b).

3 Meta-Study

In this section, we survey 274 studies (§ 3.1), probing LMs’ linguistic competence. We analyze these studies regarding their evolution, covered probing tasks and LMs (§ 3.2), and identify the apparent need for consolidating existing resources (§ 3.3).

3.1 Scope

We analyze 28k papers (P) from 2015 to August 2023 of major NLP conferences (TACL, ACL, AACL, COLING, EACL, EMNLP, NAACL, and corresponding workshops) expanded with selected work from other venues such as ICLR. To identify relevant work, we follow a semiautomatic approach. First, we automatically select papers based on their meta-data and full text.² We select a total of 493 candidate papers matching at least one of the following three criteria ($P' = \{\forall p \in P | p \in P_1 \cup p \in P_2 \cup p \in P_3\}$):

P₁: papers contain *probing* or *probe* in the title.

P₂: papers contain *probing* or *probe* in the abstract and at least five times in the main content.

P₃: papers contain *probing* or *probe* at least ten times in the main content.

We manually verified these automatically curated candidates (P') and found 274 relevant papers (P_r). We selected them as they either evaluate LMs regarding one or more linguistic phenomena as part of the analysis or as a main contribution. This involves filtering papers using the term *probing* in other senses, such as *probing hash tables* (Bogoychev and Lopez, 2016).

3.2 Analysis

Next, we analyze these 274 relevant studies (P_r).

i) Scattered evolution calls for consolidation.

First, we analyze the evolution of the relevant studies. Figure 3 relates how these studies cite each other (**probing citations** C_p) compared to other

²We use PyPDF2 v3.0.0, DBLP and semanticscholar API.

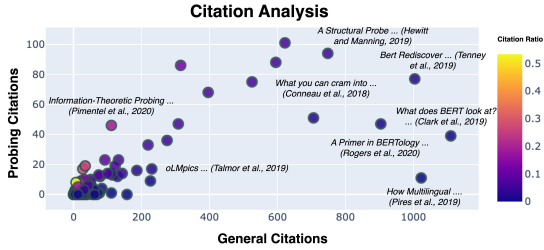


Figure 3: Citation analysis considering *probing citations* originating from the set of relevant work and every other citation (*general citations*). The color scale indicates the ratio (α) between them.

gathered citations (**general citations** C_g). Colored, we show the ratio α between these two measures $\alpha = \frac{|C_p|+1}{|C_g|+1}$. First, only a fraction of the works gained general attention, as 16 papers exceeded 200 general citations. Further, probing works cite each other rather sparsely, with an average probing citation ratio of $\alpha = 0.1$. Therefore, we see other fields are paying little attention to the linguistic competence of LMs. Paired with scattered citation patterns, we identify the need to consolidate existing resources to solidly ground research in this field.

ii) Probing work prioritizes tasks and analytics over methods. We categorize the selected work according to their probing focus: **methodological**, new methods, like control tasks (Hewitt and Liang, 2019) or minimum description length (Voita and Titov, 2020); **task-focused** assessing specific linguistic phenomena as main contributions, such as discourse relations in text (Koto et al., 2021); and **analytical** using probing tasks to analyze LMs, such as the impact of pre-training data (Zhang et al., 2021). Figure 4 shows: the majority (51.8%) of studies focus on specific probing tasks like numeric scales (Zhang et al., 2020), or morphosyntactic (Shapiro et al., 2021); 35.7% use probing as a supplementary analytical tool, for example, analyzing the effect of fine-tuning (Mosbach et al., 2020a; Zhu et al., 2022a); 12.5% address methodological problems related to probing (Wu et al., 2020; Immer et al., 2022; Zhu et al., 2022b).

iii) The dominance of classifier-based probing. Next, we analyze the specific employed probing method: **classifier**, using linear or shallow models to probe internal representations of LMs, as demonstrated in Tenney et al. (2019a); **mask**, letting LMs fill gaps to verify linguistic phenomena, as shown in Talmor et al. (2020) or Warstadt et al. (2020); **at-**

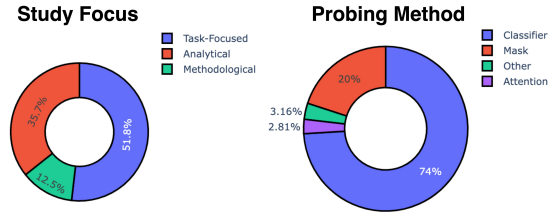


Figure 4: Categorization of the selected studies by their focus and their conducted probing method.

tention, which relies on attention patterns, as used in Pandit and Hou (2021) for bridging; and **other**, methods not belonging to the previous three categories, such as dimension selection (Torroba Hennigen et al., 2020). Most studies utilize the classifier-based probing method (74%), 20% conduct mask-based probing, and only a minority of work ($\sim 3\%$) considers attention patterns or other approaches.

iv) Tasks and LMs are barely broadly evaluated. Finally, we analyze which tasks and LMs the relevant probing studies consider. For example, Tenney et al. (2019b) considers BERT and probes POS tagging, semantic-role labeling (SRL), and other ones. Aggregated over all studies, we found a broad coverage of 289 unique tasks and 161 distinct LMs. Below, we delve into the details and highlight noteworthy findings.

We analyze how LMs and tasks are considered jointly in Figure 5. Despite the broad coverage, single studies, including fundamental ones, maintain a particular focus and consider only a fraction of LMs and tasks. For example, while most tasks (72%) were assessed on BERT, RoBERTa’s coverage has already declined to 42%. Conversely, part-of-speech tagging (POS), the most probed task, was only evaluated on 23% of the LMs, for example, not covering prominent examples like BART (Lewis et al., 2020). Notably, more recently released larger and powerful LMs, like PYTHIA (Biderman et al., 2023), UL2 (Tay et al., 2023), or LLAMA-2 (Touvron et al., 2023), and instruction-tuned LMs (FLAN-T5 (Chung et al., 2022), LLAMA-2-Chat (Touvron et al., 2023), or TK-Instruct (Wang et al., 2022) are missing almost entirely, with single more recent exceptions (Hu and Levy, 2023; Waldis et al., 2024a). Again, these insights underscore the need to consolidate existing resources for more dense coverage. This is further evident when considering Figure 5, where we sort LMs and tasks according to how often they were mentioned in the relevant works. Then, we plot

Coverage of Language Models and Probing Tasks

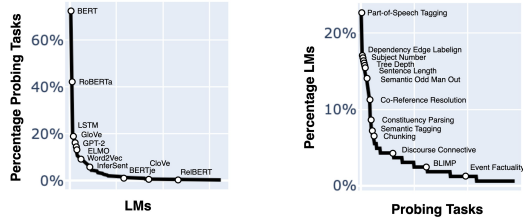


Figure 5: Overview of how many tasks single LMs cover and vice versa - single examples are highlighted.

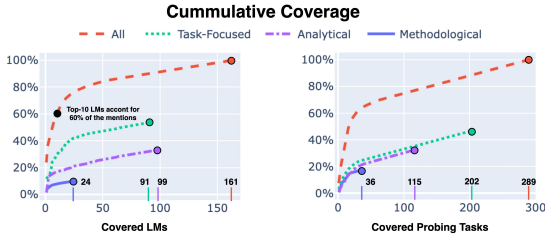


Figure 6: Cumulative coverage of LMs and tasks, considering all relevant studies and their focus.

their cumulative coverage concerning all mentions. For example, considering all studies (red line), the top-10 most mentioned LMs account for 80% of all LMs mentions (black dot). In contrast, the other 151 unique LMs account for only 40%. Comparing the paper focus, we see that methodological studies rely only on 24 LMs and 36 tasks. In contrast, task-focused and analytical work covers a similar number of LMs (91 and 99, respectively). However, due to their distinct focus, task-focused studies cover significantly more tasks (202) than analytical ones (115).

3.3 Summary

This meta-study emphasizes the need to consolidate existing resources for a comprehensive assessment of the linguistic competence of LMs — a manifold but rather blind spot in evaluation research. Apart from more thorough evaluations, such a stimulus can significantly boost future research, as happened in computer vision with ImageNet (Deng et al., 2009) or in NLP with GLUE and SuperGLUE (Wang et al., 2019a,b).

4 Holmes Benchmark

With *Holmes*, we provide an extensive ground to tackle these identified deficiencies in the existing literature and comprehensively investigate the English linguistic competence of LMs. Specifically, *Holmes* features 208 datasets addressing distinct

aspects of 66 phenomena covering *morphology*, *syntax*, *semantic*, *reasoning*, and *discourse*.

4.1 Datasets

To feature a total of 208 unique datasets, we leverage existing and established resources like OntoNotes (Weischedel et al., 2013), English Web Treebank (Silveira et al., 2014), or BLIMP (Warstadt et al., 2020) and create datasets addressing phenomena like the POS of words, their dependencies or determine the linguistic acceptability of sentences. Further, we include a range of less employed data, addressing contextualization of words (Klafka and Ettinger, 2020), reasoning (Talmor et al., 2020), semantic decomposition (White et al., 2016; Rudinger et al., 2018a,b; Govindarajan et al., 2019; Vashishtha et al., 2019), grammatical knowledge (Huebner et al., 2021), bridging (Pandit and Hou, 2021), and rhetorical (Carlson et al., 2001) and discourse (Webber et al., 2019) structure in text. Finally, we cover rarely probed phenomena like negation (Szarvas et al., 2008; Konstantinova et al., 2012; Vahtola et al., 2022), or word complexity (Paetzold and Specia, 2016).

4.2 Structure

Apart from the comprehensive scope, *Holmes* provides a clear structure for specific evaluations on different levels of aggregation. We first group the datasets according to the linguistic phenomena addressed. Then, we categorize these phenomena into their previously introduced type (see § 2) - *morphology*, *syntax*, *semantics*, *reasoning* and *discourse*. We rely on the categorization provided by the specific studies whenever given. The detailed categorization is given in § A.3.

4.3 Experimental Setup

Holmes evaluation follows the primarily used classifier-based probing paradigm, as described in § 2. Considering the internal representations allows us to maximally disentangle the understanding of distinct linguistic phenomena from each other and from other cognitive abilities (like following textual instructions). Further, this method allows us to assess any type of LMs, including sparse, static, or contextualized ones. Based on the specific dataset, we either select the embeddings of the specific input tokens (like single words for POS tagging) or average embeddings across a span or the whole sentence. We define a probing task as training a probe f_p (linear model without intermediate layers) using

these embeddings as inputs and the dataset labels as training signals. If not defined in the original data, we divide the dataset samples into train/dev/test split following a ratio of 70/10/20. We repeat this procedure five times using different random seeds and aggregate the results afterward.

4.4 Evaluations

We approximate how well an LM encodes specific linguistic phenomena using the absolute prediction performance of the probes. In addition, we rigorously evaluate the reliability of probing results using control tasks and from an information theory perspective (Voita and Titov, 2020; Hewitt and Liang, 2019). Different from commonly used prompting assessments, this particular evaluation protocol refrains from known fallacies in which the results and conclusions are sensible with specific instructions (Mizrahi et al., 2024; Min et al., 2022) or few-shot examples (Lu et al., 2023).

Task Score Metric Based on a dataset’s specific task type, we use a corresponding performance measure, macro F_1 for classification or Pearson correlation for regression. In addition, we calculate the standard deviation σ of the probe across multiple seeds. A lower σ indicates a better encoding of a given linguistic phenomenon since the measurement is robust to noise. Further, we use the task score for ranking-based evaluation of all evaluated LMs $L = \{l_1, \dots, l_m\}$ within `Holmes`. We calculate the mean winning rate mwr (in percentage), telling us how many times one LM l_1 wins against others (Liang et al., 2023). With a higher mwr , we assume an LM encodes tested linguistic phenomena better than others.

Compression Next, we evaluate the probes’ reliability from an information-theoretic perspective. Following Voita and Titov (2020), we use the compression I to measure how well a probe compresses input data. A higher I means fewer bits are needed, indicating that the given linguistic phenomenon is more clearly encoded in the embeddings.

Selectivity A reliable probe should grasp patterns relevant to the tested phenomena in the internal representations of LMs but should not be able to learn anything else. Therefore, we expect high performance when evaluating the specific dataset but low performance when we randomize training signals. We check this using control tasks introduced in Hewitt and Liang (2019). Specifically, we calcu-

late the selectivity S as the difference between the probe trained with the original labels y and the control task where we train the probe with randomly assigned labels y' . With a higher S , we assume the detected patterns are relevant for the specific phenomena under test, as random patterns do not lead to similar performance.

5 Holmes Results

Using `Holmes`, we evaluate a diverse collection of 59 LMs.³ Using the results of these extensive experiments, we first answer the research question: *what is the linguistic competence of LMs?* In doing so, we discuss the reliability of results (i), the linguistic competence of LMs concerning the unique structure of `Holmes` (ii), and how these results relate to other downstream abilities (iii). Subsequently, we examine *how linguistic competence varies among LMs*, as we find LMs prevailing for different types of linguistic phenomena (Figure 1) and delve into the effects of model architecture (iv), size (v), and instruction tuning (vi).

i) The reliability of Holmes. First, we show the reliability of probing-based evaluation, using *deviation* σ , *compression* I , and *selectivity* S results in Figure 7. Single outliers are datasets that are too hard for all LMs, as the sample size is too small, or the linguistic phenomena under test are too complex, as the ability to detect spans causes speculations in a text. We average these metrics for every dataset across all LMs. Note, for *selectivity*, we consider only base-sized model (10m-200m parameters) for computational efficiency.

First, we found a low average deviation ($\sigma = 0.02$), indicating the high reliability of probes across random seeds. These results also highlight the stability of probing results, compared to prompting-based ones where results across many paraphrased prompts lead to a deviation of $\sigma = 0.07$ reported in Mizrahi et al. (2024). Next, substantial compression (average $I = 1.9$) and selectivity (average $S = 0.31$) further confirm the probes’ reliability. Interestingly, one identifies two parallel trends for selectivity. Harder datasets with many labels, like POS tagging, are arranged around a selectivity of 0.1 to 0.4 and a task metric of 0.3. In contrast, for easier binary classification tasks (such as linguistic applicability), we observe selectivity around 0.2 to 0.5 and a task metric of 0.6 to 0.9.

³Find a complete list in Appendix § A.2.

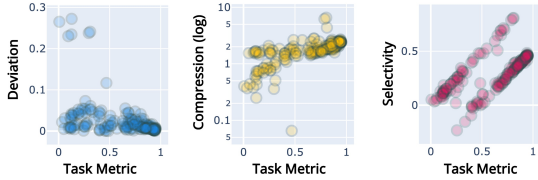


Figure 7: Reliability evaluation using *deviation*, *compression* (log), and *selectivity* on the y-axis for all 208 probing datasets. The x-axis represents the task metrics (either person correlation or macro F_1).

Further, we measure a significant ($p < 0.05$) positive correlation between the task metrics and the compression ($\tau = 0.64$) and selectivity ($\tau = 0.65$). This further confirms our reliability assumption and allows us to trust the task metric as the primary evaluation measure.

ii) The story of Holmes. We focus on what `Holmes` tells us in general and regarding formal and functional phenomena, as defined in § 2. We report in Figure 8 the *task metric*, *discriminability*, and *selectivity*, averaged for every phenomena type. Note, discriminability (Rodriguez et al., 2021) quantifies the alignment of LMs ranking of one specific dataset compared to the overall rankings using the Kendall Tau correlation. Considering these three metrics, all tested LMs strongly encode formal phenomena (*morphology* and *syntax*), which often depend on the local neighborhood of words. Therefore, we assume that LMs approximate these co-occurrences during pre-training with high precision. For example, the specific POS tag of a word, like *man* (*noun*), primarily depends on its surroundings, such as the frequent predecessor *the*. In contrast, LMs encode less information about functional phenomena (*semantics*, *reasoning*, and *discourse*) since they show a relatively low performance regarding the task metric. For these functional phenomena, we assume more complex co-occurrences are required to capture the broad context in language, such as the rhetorical relation of two distant text spans. Despite these differences between formal and functional phenomena types, they contribute to the benchmark in a balanced way. A low to medium discriminability indicates that none of these types of linguistic phenomena dominates the overall LM rankings.

This balanced influence of the five phenomena types is further visible when considering their ranking correlations (Figure 9, left). A high average correlation of 67.8 ± 6.6 with the overall results

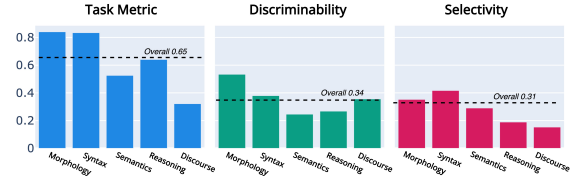


Figure 8: Average *task metric*, *difficulty*, and *discriminability* for each phenomena type. The dashed lines show the average measure over all datasets.

(last column/row) hints that they are facets of a broader occurrence but share common characteristics. Still, breaking into categories is meaningful, as the phenomena types (first five columns/rows) are medium correlated (average of 53.9 ± 14.5). Analyzing the results of phenomena types further highlights the value of this distinction. While results of *morphology* and *syntax* are similarly correlated with the overall results (68.2 and 70.2), their direct correlation (69.1) indicates their supplementary nature. Further, *discourse* results show the lowest correlation with others (44.8 ± 16.1), indicating the particular scope.

iii) The companions of Holmes. We analyze how the results of `Holmes` and those from other evaluations focusing on downstream applications align (Figure 9, right). We select the OpenLLM benchmark (Beeching et al., 2023), as it covers a wide range of open LMs, in contrast to others like HELM (Liang et al., 2023). First, `Holmes` and OpenLLM results of jointly evaluated LMs are medium correlated, hinting that the linguistic competence of LMs is partly aligned with their downstream abilities. The nature of this alignment is further evident when focusing on *morphology*, *reasoning*, and *discourse*. Interestingly, and in contrast to *syntax* and *semantics*, their correlation to the OpenLLM and `Holmes` overall results is similar. Therefore, these three phenomena presumably represent skills that are more tested in the general benchmarks. These correlation patterns are consistent across the three most meaningful OpenLLM datasets (*MMLU*, *TruthfulQA*, and *GSM8K*). As *TruthfulQA* shows lower correlations with the linguistic phenomena and other datasets within OpenLLM, we presume this dataset captures distinctly different skills (possibly knowledge).⁴ These insights show how different benchmarks provide a different scope and supplement themselves simulta-

⁴Further, it's also known that we need to expect this dataset to be fully leaked (Balloccu et al., 2024).

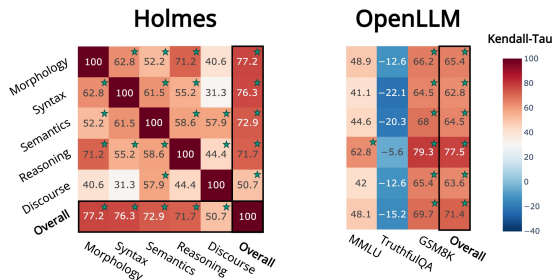


Figure 9: Kendall-tau correlation within `Holmes` (left) and compared to the `OpenLLM` benchmark (right). Green stars indicate significant correlations ($p < 0.05$).

neously. Further, the above analysis shows, again, the value of assessing the linguistic competence of LMs across different phenomena types, for fine-grained analyses.

iv) The effect of language model architecture.

Next, we discuss the impact of model architecture on the linguistic competence of LMs. In Figure 11 (left), we compare encoder and decoder LMs. Due to the absence of big encoder LMs, we consider five *encoder* and six *decoder* LMs with up to 220m parameters. Encoder LMs show a higher *mwr* of 52% than decoder LMs (21%). This observation is the most saturated for *morphology* or *syntax*, encompassing a variety of token-level phenomena, like part-of-speech. We assume that the missing bi-directional encoding of decoder LMs causes this lower performance because the available context of one token heavily depends on its position. Thus, even common tokens, like *the*, have different potential representations - at the beginning or in the middle of a sentence. These instabilities are further evident when considering Figure 11 (right) which reports the accuracy for the top-20 most common POS tokens (such as *the*) based on the *pos*, *xpos*, *upos* dataset. Given their high frequency, one expects stable prediction performance. Surprisingly, encoder LMs (BERT and RoBERTa) show higher median accuracy and clearly lower deviations compared to the same-size decoder counterpart (GPT2). While scaling model size to 12B (Pythia) and 70B (Llama-2) allows for improved accuracy and lower deviations, decoder LMs do not match the encoder performance, even up to **700 times bigger**.

v) The effect of scaling parameters. We discuss how the number of parameters influences the linguistic competence of LMs. Given the variety of LMs of different sizes, we focus on the Pythia

(decoder-only) and T5 (encoder-decoder) families. From Figure 10, we observe for both Pythia and T5 that the linguistic competence scales with model size, and it is particularly pronounced after exceeding 0.5B (Pythia) and 1.0B (T5) parameters. Again, model architecture is crucial, as T5 LMs (encoder-decoder) exhibit a clearly higher mean winning rate of 40 – 70% than Pythia (decoder-only) ones with *mwr* of 20 – 60%. Further, we found formal phenomena evolving differently with increased model size than functional ones. Specifically, *morphology* and *syntax* start at a lower level, with an apparent performance jump after 0.5B (Pythia) and 1.0B (T5) parameters, followed by slow but steady growth. Differently, *semantics*, *reasoning*, and *discourse* start at a higher *mwr*, followed by a continuous improvement as the model size grows. From these results, we assume more parameters allow LMs to better approximate simpler co-occurrences in the near neighborhood of words to understand formal phenomena like word dependencies. In contrast, more parameters do not have the same pronounced effect on functional phenomena, like rhetorical relations, which require an LM to acquire more distant and complex word co-occurrences.

Model	Morphology	Syntax	Semantics	Reasoning	Discourse	Overall
Comparison against <code>Llama-2</code> with 7 billion parameters						
Llama-2-Chat	-8%	+3%	-5%	-9%	-3%	-2%
Comparison against <code>T5</code> with 11 billion parameters						
FLAN-T5	+10%	+2%	-2%	+6%	-2%	+1%
Comparison against <code>Pythia</code> with 12 billion parameters						
Dolly-v2	+4%	-3%	-9%	-3%	+4%	-4%
Comparison against <code>Llama-2</code> with 13 billion parameters						
Tulu-2	+5%	+2%	-15%	0%	-30%	-8%
Orca-2	-1%	-3%	-4%	+4%	-5%	-2%
Llama-2-chat	+3%	+1%	-6%	+3%	-1%	-1%
Vicuna-v1.5	+23%	+7%	-3%	+6%	-6%	+4%
Comparison against <code>UL2</code> with 20 billion parameters						
FLAN-UL2	+40%	+16%	+7%	+13%	+1%	+13%
Comparison against <code>Mixtral</code> with ~47 billion parameters						
Mixtral-Instruct	+4%	+3%	0%	+6%	-2%	+2%
Comparison against <code>Llama-2</code> with 70 billion parameters						
Tulu-2	+15%	0%	-11%	-3%	0%	-2%
Llama-2-Chat	+23%	+14%	+2%	+4%	+17%	+10%
Average	+10%	+4%	-3%	+4%	-2%	+1%

Table 1: Effect of instruction tuning on the mean winning rate compared to the pre-trained LMs.

vi) The effect of instruction tuning. Finally, we focus on how instruction tuning affects LMs' linguistic competence and compare the tuned LMs with their base models—for example, FLAN-UL2 vs. UL2. From results in Table 1, we note less saturated effects for the overall scope while being more pronounced for the five phenomenon types - again emphasizing the structured and comprehensive evaluation of linguistic competence. On average, we found instruction tuning has the highest

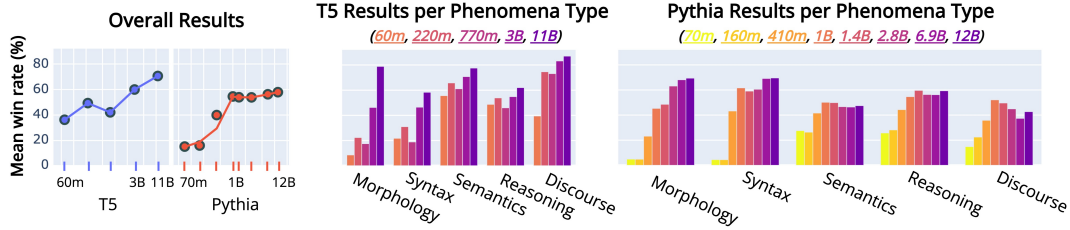


Figure 10: Effect of scaling LM parameters considering the T5 and Pythia model families providing eight and five different sizes. We address the overall scope (left) and the different types of linguistic phenomena (right).

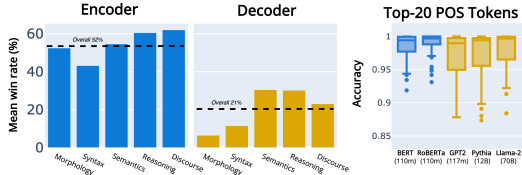


Figure 11: Comparison of the phenomenon types for encoder and decoder LMs (left) and on the right, the accuracy of the top-20 most common tokens of the three part-of-speech probing datasets for BERT, RoBERTa, GPT2, Pythia, and Llama-2.

effect on *morphology* (+10%) followed by *syntax* (+4%), *reasoning* (+4%), and a negative effect for *semantics* -3% and *discourse* -2%. These results confirm previous assumptions that instruction tuning updates are often superficial (Yadav et al., 2023; Hershcovitch et al., 2024; Sharma et al., 2023) and that LMs are better at mimicking language (formal phenomena) than understanding it, measured with functional phenomena (Mahowald et al., 2024). Further, larger models benefit more from instruction tuning. Llama-2-70b-Chat and FLAN-UL2 gain up to +23% and +40% for *morphology* and +10% and +13% on average. In addition, decoder-only LMs (Llama-2 and Pythia) tend to show less pronounced positive effects than encoder-decoder LMs (FLAN-T5-XXL and FLAN-UL2). However, they better understand *reasoning* phenomena. When comparing LMs based on Llama-2-13b, we see that specific fine-tuning methods shape the LMs differently. The top-ranked 13b LM for Holmes and OpenLLM, Vicuna, was trained on 125k instructions, less than other models. Thus, high quality is more important than the number of instructions for LMs’ linguistic competence. Tulu loses performance while being trained on a large mixture of data (approx. 330k instructions), the same for its 70b version. Finally, the focus of Orca-2 on reasoning is also reflected in its embedding space. These insights show again that while provid-

ing a particular perspective, Holmes shows clear differences between LMs and allows us to map them to methodological decisions.

6 Efficiency

Seamless, easy, cost-effective integration of new LMs is crucial for widely adopting a benchmark. As Holmes covers many datasets and examples, it is computationally heavy in encoding text and training the probes. It takes approx. 6 GPU days to encode the 70 million tokens (~230k pages of text) and 2 days to run the 208 probes for a 70b model. To account for this issue, we introduce FlashHolmes, a streamlined version of Holmes. It allows the evaluation of new LMs with a fraction of the compute while maintaining evaluation integrity.

Besides excluding licensed data (18 probing datasets), we analyze the effect of discarding training instances. As a result, we reduce the computation for encoding and the actual probing simultaneously. We follow Perlitz et al. (2023) and calculate the *rank resolution*, 95% CI of model rank difference. This measure indicates the maximum expected rank deviation from evaluating an LM on FlashHolmes compared to Holmes. For example, a rank resolution of one means that an LM evaluated on FlashHolmes and Holmes has the same rank or switch place with its neighbors with a probability of 95%. Figure 12 shows the resulting rank resolution when training only on a fraction of the instances, from 1/2 to 1/512. Solely focusing on efficiency (1/512) still provides a decent rank resolution of ~2.7. In contrast, considering 1/2 of the training data results in the best reliability of ~1.0. To balance benchmark reliability and efficiency, we compose FlashHolmes using 1/32 of the training instances. Precisely, it reduces the computation expenses of evaluating LMs to ~3% of what Holmes would have required while pre-

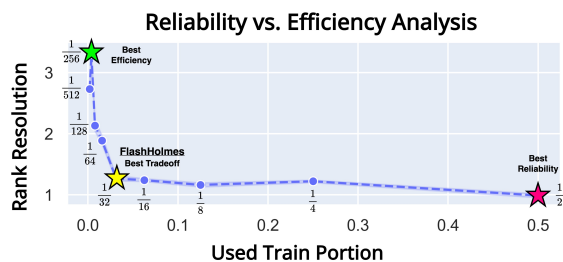


Figure 12: Analysis of the reliability vs. efficiency trade-off when reducing the number of training data.

706 serving a high rank-correlation of ~ 1.3 .

707 7 Related Work

708 **Benchmarking LMs** Benchmarks approximate
 709 LMs abilities like general language understanding
 710 (Wang et al., 2019b,a), out-of-distribution gener-
 711 alization (Yang et al., 2023; Waldis et al., 2024b),
 712 adversarial scenarios (Nie et al., 2020; Wang et al.,
 713 2021), or retrieval like *BEIR* (Thakur et al., 2021)
 714 or *MTEB* (Muennighoff et al., 2023). With the
 715 advent of larger LMs, the methodological focus
 716 shifted to prompting-based evaluations which eval-
 717 uate the LMs’ response to provided instructions
 718 (Brown et al., 2020; Hendrycks et al., 2021; Sri-
 719 vastava et al., 2022) covering application-oriented
 720 tasks (Liang et al., 2023), or mathematical reason-
 721 ing (e.g., *GSM8K* (Cobbe et al., 2021)).

722 Assessing the Linguistic Competence of LMs

723 The analysis of LMs’ linguistic competence ranges
 724 from analyzing static word vectors (Köhn, 2015),
 725 sentence embeddings (Conneau et al., 2018; Adi
 726 et al., 2017), the internals of translation models
 727 (Shi et al., 2016; Bau et al., 2019), or contextual-
 728 ized LMs (Tenney et al., 2019b,a; Hewitt and Man-
 729 ning, 2019). Other work addressed methodolog-
 730 ical aspects, such as using control tasks (Hewitt
 731 and Liang, 2019), assessing LMs from an infor-
 732 mation theory perspective (Voita and Titov, 2020;
 733 Pimentel et al., 2020), or evaluating causal effects
 734 in LMs (Elazar et al., 2021). Finally, another line
 735 of work focuses on whether LMs follow human
 736 understanding of linguistic competence when solv-
 737 ing downstream tasks (Belinkov, 2022; Aw et al.,
 738 2023; Mahowald et al., 2024). However, Mosbach
 739 et al. (2020b) and Waldis et al. (2024a) found fine-
 740 tuning for downstream tasks actually hurting the
 741 understanding of linguistic phenomena.

742 While prior studies assessing the linguistic com-
 743 petence of LMs tend to focus on a limited set of
 744 linguistic phenomena or models, *Holmes* provides

745 extensive coverage of both phenomena and eval-
 746 uated LMs. Unlike recent evaluations based on
 747 prompting methods (Blevins et al., 2023; Liang
 748 et al., 2023; Amouyal et al., 2024), *Holmes* as-
 749 sesses the internal representations of LMs directly.
 750 This approach allows for detailed analysis of spe-
 751 cific model characteristics, such as architecture,
 752 and helps separate the linguistic competence from
 753 other cognitive abilities. Thereby, we respond to
 754 recent calls for a thorough and explicit evaluation
 755 of linguistic phenomena (Hu and Levy, 2023; Lu
 756 et al., 2023; Mahowald et al., 2024).

757 8 Conclusion

758 *Holmes* marks the most up-to-date and extensive
 759 consolidation of existing resources addressing the
 760 need to assess the linguistic competence of LMs in
 761 isolation. Our experiments demonstrate that LMs’
 762 linguistic competence is pronounced regarding for-
 763 mal phenomena but lacks functional ones when
 764 information about broader textual contexts, such as
 765 rhetorical structure, is required. Further, size, ar-
 766 chitecture, and instruction tuning crucially account
 767 for differences among LMs. As LM and resources
 768 in the landscape of linguistics continue to grow,
 769 we will actively extend *Holmes* with further prob-
 770 ing datasets, evaluate upcoming LMs, and plan to
 771 incorporate multilingualism.

772 Ethical Considerations and Limitations

773 **Language** *Holmes* as well as *FlashHolmes*
 774 solely assess linguistic phenomena for the English
 775 language. As we plan to expand the benchmark and
 776 scope of multilingual data, we focus momentarily
 777 on English because of the widespread availability
 778 of resources, including curated corpora and the
 779 diversity of available LMs.

780 **Phenomena and LM Coverage** We agree with
 781 Liang et al. (2023) and see one fundamental aspect
 782 in composing a benchmark in acknowledging its
 783 incompleteness. Both linguistic phenomena and
 784 LMs are a subset of the variety of available re-
 785 sources. We consolidated them carefully to provide
 786 a comprehensive scope of the linguistic compe-
 787 tence and various LMs. However, as benchmarks
 788 evolve as tools to assess LMs, we will further ex-
 789 pand *Holmes* both with the existing and upcoming
 790 LMs and data resources.

791 **Data Availability** Linguistic annotations, in par-
 792 ticular more complex ones targeting phenomena

like *discourse*, are money and time-wise expensive. Out of 208 datasets included in `Holmes`, 18 probing datasets are based on licensed resources and are not freely available. However, with `FlashHolmes`, we provide an effective and efficient alternative based on open-access resources. Furthermore, upon confirming the granted access, we are happy to share our probing datasets, including those based on the licensed resources.

Bias As `Holmes` relies on existing resources, it inherits the bias embodied in this data. Examples of such bias are gender equality or gender fairness, like the use of neo pronouns such as *em* in [Lauscher et al. \(2023\)](#).

References

Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. 2017. [Fine-grained analysis of sentence embeddings using auxiliary prediction tasks](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.

Ehsan Aghazadeh, Mohsen Fayyaz, and Yadollah Yaghoobzadeh. 2022. [Metaphors in pre-trained language models: Probing and generalization across datasets and languages](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2037–2050, Dublin, Ireland. Association for Computational Linguistics.

Samuel Amouyal, Aya Meltzer-Asscher, and Jonathan Berant. 2024. [Large language models for psycholinguistic plausibility pretesting](#). In *Findings of the Association for Computational Linguistics: EAACL 2024*, pages 166–181, St. Julian’s, Malta. Association for Computational Linguistics.

Khai Loong Aw, Syrielle Montariol, Badr AlKhamissi, Martin Schrimpf, and Antoine Bosselut. 2023. [Instruction-tuning aligns llms to the human brain](#). *CoRR*, abs/2312.00575.

Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondrej Dusek. 2024. [Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93, St. Julian’s, Malta. Association for Computational Linguistics.

Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James R. Glass. 2019. [Identifying and controlling important neurons in neural machine translation](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Edward Beeching, Clémentine Fourier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. [Open llm leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard). https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard.

Yonatan Belinkov. 2022. [Probing classifiers: Promises, shortcomings, and advances](#). *Computational Linguistics*, 48(1):207–219.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. [What do neural machine translation models learn about morphology?](#) In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. 2023. [Pythia: A suite for analyzing large language models across training and scaling](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 2397–2430. PMLR.

Julia Birke and Anoop Sarkar. 2006. [A clustering approach for nearly unsupervised recognition of nonliteral language](#). In *11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 329–336, Trento, Italy. Association for Computational Linguistics.

Terra Blevins, Hila Gonen, and Luke Zettlemoyer. 2023. [Prompting language models for linguistic structure](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6649–6663, Toronto, Canada. Association for Computational Linguistics.

Nikolay Bogoychev and Adam Lopez. 2016. [N-gram language models for massively parallel devices](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1944–1953, Berlin, Germany. Association for Computational Linguistics.

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. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33*:

904					
905					
906					
907	Lynn Carlson, Daniel Marcu, and Mary Ellen Okurovsky. 2001. Building a discourse-tagged corpus in the framework of rhetorical structure theory . In <i>Proceedings of the SIGDIAL 2001 Workshop, The 2nd Annual Meeting of the Special Interest Group on Discourse and Dialogue, Saturday, September 1, 2001 to Sunday, September 2, 2001, Aalborg, Denmark</i> . The Association for Computer Linguistics.				
915	Noam Chomsky. 1965. <i>Aspects of the Theory of Syntax</i> . The MIT Press, Cambridge.				
917	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models . <i>CoRR</i> , abs/2210.11416.				
928	Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: pre-training text encoders as discriminators rather than generators . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.				
934	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems . <i>CoRR</i> , abs/2110.14168.				
940	Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single $\&\!#\ast$ vector: Probing sentence embeddings for linguistic properties . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.				
948	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm .				
953	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database . <i>2009 IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 248–255.				
958	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of				
	deep bidirectional transformers for language understanding . 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.				
	Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2021. Amnesic probing: Behavioral explanation with amnesic counterfactuals . <i>Transactions of the Association for Computational Linguistics</i> , 9:160–175.				
	Rudolf Franz Flesch. 1948. A new readability yardstick . <i>The Journal of applied psychology</i> , 32(3):221–233.				
	William Gantt, Lelia Glass, and Aaron Steven White. 2022. Decomposing and recomposing event structure . <i>Transactions of the Association for Computational Linguistics</i> , 10:17–34.				
	Vagrant Gautam, Eileen Bingert, D. Zhu, Anne Lauscher, and Dietrich Klakow. 2024. Robust pronoun use fidelity with english llms: Are they reasoning, repeating, or just biased? <i>CoRR</i> , abs/2404.03134.				
	Mario Giulianelli, Jack Harding, Florian Mohnert, Dieuwke Hupkes, and Willem Zuidema. 2018. Under the hood: Using diagnostic classifiers to investigate and improve how language models track agreement information . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 240–248, Brussels, Belgium. Association for Computational Linguistics.				
	Venkata Govindarajan, Benjamin Van Durme, and Aaron Steven White. 2019. Decomposing generalization: Models of generic, habitual, and episodic statements . <i>Transactions of the Association for Computational Linguistics</i> , 7:501–517.				
	Zellig S Harris. 1954. Distributional structure . <i>Word</i> , 10(2-3):146–162.				
	Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.				
	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: decoding-enhanced bert with disentangled attention . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.				
	Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals . In <i>Proceedings of the 5th International</i>				

1016			
1017			
1018			
1019	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.		
1020			
1021			
1022			
1023			
1024			
1025	Moshik Hershcovitch, Leshem Choshen, Andrew Wood, Ilias Enmouri, Peter Chin, Swaminathan Sundararaman, and Danny Harnik. 2024. Lossless and near-lossless compression for foundation models . <i>CoRR</i> , abs/2404.15198.		
1026			
1027			
1028			
1029			
1030	John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.		
1031			
1032			
1033			
1034			
1035			
1036			
1037	John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations . 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 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.		
1038			
1039			
1040			
1041			
1042			
1043			
1044			
1045	Yufang Hou. 2018. Enhanced word representations for bridging anaphora resolution . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 1–7, New Orleans, Louisiana. Association for Computational Linguistics.		
1046			
1047			
1048			
1049			
1050			
1051			
1052	Yufang Hou. 2020. Bridging anaphora resolution as question answering . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1428–1438, Online. Association for Computational Linguistics.		
1053			
1054			
1055			
1056			
1057	Jennifer Hu and Roger Levy. 2023. Prompting is not a substitute for probability measurements in large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5040–5060, Singapore. Association for Computational Linguistics.		
1058			
1059			
1060			
1061			
1062			
1063	Philip A. Huebner, Elinor Sulem, Fisher Cynthia, and Dan Roth. 2021. BabyBERTa: Learning more grammar with small-scale child-directed language . In <i>Proceedings of the 25th Conference on Computational Natural Language Learning</i> , pages 624–646, Online. Association for Computational Linguistics.		
1064			
1065			
1066			
1067			
1068			
1069	Alexander Immer, Lucas Torroba Hennigen, Vincent Fortuin, and Ryan Cotterell. 2022. Probing as quantifying inductive bias . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1839–1851, Dublin, Ireland. Association for Computational Linguistics.		
1070			
1071			
1072			
			1073
			1074
			1075
			1076
			1077
			1078
			1079
			1080
			1081
			1082
			1083
			1084
			1085
			1086
			1087
			1088
			1089
			1090
			1091
			1092
			1093
			1094
			1095
			1096
			1097
			1098
			1099
			1100
			1101
			1102
			1103
			1104
			1105
			1106
			1107
			1108
			1109
			1110
			1111
			1112
			1113
			1114
			1115
			1116
			1117
			1118
			1119
			1120
			1121
			1122
			1123
			1124
			1125
			1126
			1127
			1128
			1129
			1130

1131 *Proceedings of the 6th Workshop on Representation*
1132 *Learning for NLP (RepL4NLP-2021)*, pages 8–19,
1133 Online. Association for Computational Linguistics. 1189 1190

1134 Zhenzhong Lan, Mingda Chen, Sebastian Goodman,
1135 Kevin Gimpel, Piyush Sharma, and Radu Soricut.
1136 2020. [ALBERT: A lite BERT for self-supervised](#)
1137 [learning of language representations](#). In *8th Inter-*
1138 *national Conference on Learning Representations,*
1139 *ICLR 2020, Addis Ababa, Ethiopia, April 26-30,*
1140 *2020*. OpenReview.net. 1191 1192 1193 1194

1141 Anne Lauscher, Debora Nozza, Ehm Miltersen, Archie
1142 Crowley, and Dirk Hovy. 2023. [What about “em”?](#)
1143 [how commercial machine translation fails to handle](#)
1144 [\(neo-\)pronouns](#). In *Proceedings of the 61st Annual*
1145 *Meeting of the Association for Computational Lin-*
1146 *guistics (Volume 1: Long Papers)*, pages 377–392,
1147 Toronto, Canada. Association for Computational Lin-
1148 guistics. 1195 1196 1197 1198 1199

1149 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan
1150 Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
1151 Veselin Stoyanov, and Luke Zettlemoyer. 2020.
1152 [BART: Denoising sequence-to-sequence pre-training](#)
1153 [for natural language generation, translation, and com-](#)
1154 [prehension](#). In *Proceedings of the 58th Annual Meet-*
1155 *ing of the Association for Computational Linguistics,*
1156 *pages 7871–7880*, Online. Association for Computa-
1157 tional Linguistics. 1200 1201 1202

1158 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris
1159 Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian
1160 Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar,
1161 Benjamin Newman, Binhang Yuan, Bobby Yan,
1162 Ce Zhang, Christian Alexander Cosgrove, Christo-
1163 pher D Manning, Christopher Re, Diana Acosta-
1164 Navas, Drew Arad Hudson, Eric Zelikman, Esin
1165 Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren,
1166 Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel
1167 Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun,
1168 Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar
1169 Khattab, Peter Henderson, Qian Huang, Ryan An-
1170 drew Chi, Sang Michael Xie, Shibani Santurkar,
1171 Surya Ganguli, Tatsunori Hashimoto, Thomas Icard,
1172 Tianyi Zhang, Vishrav Chaudhary, William Wang,
1173 Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Ko-
1174 reeda. 2023. [Holistic evaluation of language models](#).
1175 *Transactions on Machine Learning Research*. Fea-
1176 tured Certification, Expert Certification. 1203 1204 1205 1206 1207 1208 1209 1210

1177 Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg.
1178 2016. [Assessing the ability of LSTMs to learn syntax-](#)
1179 [sensitive dependencies](#). *Transactions of the Associa-*
1180 *tion for Computational Linguistics*, 4:521–535. 1211 1212 1213 1214 1215 1216 1217

1181 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
1182 dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
1183 Luke Zettlemoyer, and Veselin Stoyanov. 2019.
1184 [Roberta: A robustly optimized BERT pretraining](#)
1185 [approach](#). *CoRR*, abs/1907.11692. 1218 1219 1220 1221

1186 Ilya Loshchilov and Frank Hutter. 2019. [Decoupled](#)
1187 [weight decay regularization](#). In *7th International*
1188 *Conference on Learning Representations, ICLR 2019,*
New Orleans, LA, USA, May 6-9, 2019. OpenRe-
1189 view.net. 1222 1223 1224 1225 1226 1227

Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Har-
1189 ish Tayyar Madabushi, and Iryna Gurevych. 2023.
1192 [Are emergent abilities in large language models just](#)
1193 [in-context learning?](#) *CoRR*, abs/2309.01809. 1228 1229 1230 1231 1232 1233 1234 1235 1236

Kyle Mahowald, Anna A. Ivanova, Idan A. Blank,
1189 Nancy Kanwisher, Joshua B. Tenenbaum, and
1196 Evelina Fedorenko. 2024. [Dissociating language](#)
1197 [and thought in large language models](#). *Trends in*
1198 *Cognitive Sciences*. 1237 1238 1239 1240 1241 1242 1243 1244

George A. Miller. 1995. [Wordnet: A lexical database](#)
1195 [for english](#). *Communications of the ACM*, 38(11):39–
1200 41. 1201 1202

Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe,
1195 Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-
1196 moyer. 2022. [Rethinking the role of demonstrations:](#)
1197 [What makes in-context learning work?](#) In *Proceed-*
1198 *ings of the 2022 Conference on Empirical Methods in*
1203 *Natural Language Processing*, pages 11048–11064,
1204 Abu Dhabi, United Arab Emirates. Association for
1205 Computational Linguistics. 1206 1207 1208 1209 1210

Arindam Mitra, Luciano Del Corro, Shweti Mahajan,
1211 Andrés Cudas, Clarisse Simões, Sahaj Agrawal, Xuxi
1212 Chen, Anastasia Razdaibiedina, Erik Jones, Kriti Ag-
1213 garwal, Hamid Palangi, Guoqing Zheng, Corby Ros-
1214 set, Hamed Khanpour, and Ahmed Awadallah. 2023.
1215 [Orca 2: Teaching small language models how to rea-](#)
1216 [son](#). *CoRR*, abs/2311.11045. 1217

Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror,
1218 Dafna Shahaf, and Gabriel Stanovsky. 2024. [State of](#)
1219 [what art? A call for multi-prompt LLM evaluation](#).
1220 *CoRR*, abs/2401.00595. 1221

Michael Mohler, Mary Brunson, Bryan Rink, and Marc
1222 Tomlinson. 2016. [Introducing the LCC metaphor](#)
1223 [datasets](#). In *Proceedings of the Tenth International*
1224 *Conference on Language Resources and Evaluation*
1225 *(LREC’16)*, pages 4221–4227, Portorož, Slovenia.
1226 European Language Resources Association (ELRA). 1227

Roser Morante and Eduardo Blanco. 2012. [*SEM 2012](#)
1228 [shared task: Resolving the scope and focus of nega-](#)
1229 [tion](#). In **SEM 2012: The First Joint Conference on*
1230 *Lexical and Computational Semantics – Volume 1:*
1231 *Proceedings of the main conference and the shared*
1232 *task, and Volume 2: Proceedings of the Sixth Interna-*
1233 *tional Workshop on Semantic Evaluation (SemEval*
1234 *2012)*, pages 265–274, Montréal, Canada. Associa-
1235 tion for Computational Linguistics. 1236

Marius Mosbach, Anna Khokhlova, Michael A. Hed-
1237 derich, and Dietrich Klakow. 2020a. [On the interplay](#)
1238 [between fine-tuning and sentence-level probing for](#)
1239 [linguistic knowledge in pre-trained transformers](#). In
1240 *Proceedings of the Third BlackboxNLP Workshop on*
1241 *Analyzing and Interpreting Neural Networks for NLP,*
1242 *pages 68–82*, Online. Association for Computational
1243 Linguistics. 1244

1245	Marius Mosbach, Anna Khokhlova, Michael A. Hedderich, and Dietrich Klakow. 2020b. On the Interplay Between Fine-tuning and Sentence-level Probing for Linguistic Knowledge in Pre-trained Transformers . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2502–2516, Online. Association for Computational Linguistics.	1303
1246		1304
1247		1305
1248		1306
1249		1307
1250		1308
1251		
1252	Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive text embedding benchmark . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.	1309
1253		1310
1254		1311
1255		1312
1256		1313
1257		
1258	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	1314
1259		1315
1260		1316
1261		1317
1262		1318
1263		1319
1264		
1265	Allen Nie, Erin Bennett, and Noah Goodman. 2019. DisSent: Learning sentence representations from explicit discourse relations . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4497–4510, Florence, Italy. Association for Computational Linguistics.	1320
1266		1321
1267		1322
1268		1323
1269		1324
1270		1325
1271		1326
1272		1327
1273		1328
1274		1329
1275		
1276		1330
1277		1331
1278		1332
1279		1333
1280		1334
1281		1335
1282		1336
1283		
1284		1337
1285		1338
1286		1339
1287		1340
1288		
1289		1341
1290		1342
1291		1343
1292		1344
1293		1345
1294		1346
1295		
1296		1347
1297		1348
1298		1349
1299		1350
1300		1351
1301		1352
1302		1353
		1354
		1355
		1356
		1357
		1358
		1359

1360					
1361					
1362					
1363					
1364	Rachel Rudinger, Aaron Steven White, and Benjamin Van Durme. 2018b.	Neural models of factuality .	In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 731–744, New Orleans, Louisiana. Association for Computational Linguistics.		
1365					
1366					
1367					
1368					
1369					
1370					
1371					
1372	Naomi Shapiro, Amandalynne Paullada, and Shane Steinert-Threlkeld. 2021.	A multilabel approach to morphosyntactic probing .	In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 4486–4524, Punta Cana, Dominican Republic. Association for Computational Linguistics.		
1373					
1374					
1375					
1376					
1377					
1378	Pratyusha Sharma, Jordan T. Ash, and Dipendra Misra. 2023.	The truth is in there: Improving reasoning in language models with layer-selective rank reduction .	<i>CoRR</i> , abs/2312.13558.		
1379					
1380					
1381					
1382	Xing Shi, Inkit Padhi, and Kevin Knight. 2016.	Does string-based neural MT learn source syntax?	In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1526–1534, Austin, Texas. Association for Computational Linguistics.		
1383					
1384					
1385					
1386					
1387					
1388	Natalia Silveira, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Chris Manning. 2014.	A gold standard dependency corpus for English .	In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)</i> , pages 2897–2904, Reykjavik, Iceland. European Language Resources Association (ELRA).		
1389					
1390					
1391					
1392					
1393					
1394					
1395					
1396	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013.	Recursive deep models for semantic compositionality over a sentiment treebank .	In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.		
1397					
1398					
1399					
1400					
1401					
1402					
1403					
1404	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameeet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi,				
1405					
1406					
1407					
1408					
1409					
1410					
1411					
1412					
1413					
1414					
1415					
1416					
1417					
	Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. 2022.	Beyond the imitation game: Quantifying and extrapolating the capabilities of language models .	<i>CoRR</i> , abs/2206.04615.		1418 1419 1420 1421 1422 1423
	Gerard J Steen, Aletta G Dorst, J Berenike Herrmann, Anna A Kaal, Tina Krennmayr, Tryntje Pasma, et al. 2010.	A method for linguistic metaphor identification .	Converging evidence in language and communication research. John Benjamins Publishing Company Amsterdam.		1424 1425 1426 1427 1428 1429
	Shivchander Sudalairaj, Abhishek Bhandwaldar, Aldo Pareja, Kai Xu, David D. Cox, and Akash Srivastava. 2024.	LAB: large-scale alignment for chatbots .	<i>CoRR</i> , abs/2403.01081.		1430 1431 1432 1433
	György Szarvas, Veronika Vincze, Richárd Farkas, and János Csirik. 2008.	The BioScope corpus: annotation for negation, uncertainty and their scope in biomedical texts .	In <i>Proceedings of the Workshop on Current Trends in Biomedical Natural Language Processing</i> , pages 38–45, Columbus, Ohio. Association for Computational Linguistics.		1434 1435 1436 1437 1438 1439 1440
	Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020.	oLMpics-On What Language Model Pre-training Captures .	<i>Transactions of the Association for Computational Linguistics</i> , 8:743–758.		1441 1442 1443 1444
	Yi Tay, Mostafa Dehghani, Vinh Q. Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Dara Bahri, Tal Schuster, Huaixiu Steven Zheng, Denny Zhou, Neil Houlsby, and Donald Metzler. 2023.	UL2: unifying language learning paradigms .	In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.		1445 1446 1447 1448 1449 1450 1451 1452
	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019a.	BERT rediscovered the classical NLP pipeline .	In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4593–4601, Florence, Italy. Association for Computational Linguistics.		1453 1454 1455 1456 1457 1458
	Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019b.	What do you learn from context? probing for sentence structure in contextualized word representations .	In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.		1459 1460 1461 1462 1463 1464 1465 1466 1467
	Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021.	BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models .	In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .		1468 1469 1470 1471 1472 1473 1474

1475	Lucas Torroba Hennigen, Adina Williams, and Ryan Cotterell. 2020. Intrinsic probing through dimension selection . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 197–216, Online. Association for Computational Linguistics.	1533
1476		1534
1477		1535
1478		1536
1479		1537
1480		1538
1481	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models . <i>CoRR</i> , abs/2307.09288.	1539
1482		1540
1483		1541
1484		1542
1485		1543
1486		1544
1487		1545
1488		1546
1489		1547
1490		1548
1491		1549
1492		1550
1493		1551
1494		1552
1495		1553
1496		1554
1497		1555
1498		1556
1499		1557
1500		1558
1501		1559
1502		1560
1503		1561
1504	Teemu Vahtola, Mathias Creutz, and Jörg Tiedemann. 2022. It is not easy to detect paraphrases: Analysing semantic similarity with antonyms and negation using the new SemAntoNeg benchmark . In <i>Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 249–262, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	1562
1505		1563
1506		1564
1507		1565
1508		1566
1509		1567
1510		1568
1511		1569
1512	Siddharth Vashishtha, Benjamin Van Durme, and Aaron Steven White. 2019. Fine-grained temporal relation extraction . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2906–2919, Florence, Italy. Association for Computational Linguistics.	1570
1513		1571
1514		1572
1515		1573
1516		1574
1517		1575
1518	Sara Veldhoen, Dieuwke Hupkes, and Willem H. Zuidema. 2016. Diagnostic classifiers revealing how neural networks process hierarchical structure . In <i>Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016</i> , volume 1773 of <i>CEUR Workshop Proceedings</i> . CEUR-WS.org.	1576
1519		1577
1520		1578
1521		1579
1522		1580
1523		1581
1524		1582
1525		1583
1526		1584
1527	Elena Voita and Ivan Titov. 2020. Information-theoretic probing with minimum description length . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 183–196, Online. Association for Computational Linguistics.	1585
1528		1586
1529		1587
1530		1588
1531		1589
1532		1590
	Andreas Waldis, Yufang Hou, and Iryna Gurevych. 2024a. Dive into the chasm: Probing the gap between in- and cross-topic generalization . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 2197–2214, St. Julian’s, Malta. Association for Computational Linguistics.	1533
		1534
		1535
		1536
		1537
		1538
	Andreas Waldis, Yufang Hou, and Iryna Gurevych. 2024b. How to handle different types of out-of-distribution scenarios in computational argumentation? a comprehensive and fine-grained field study . <i>CoRR</i> , abs/2309.08316.	1539
		1540
		1541
		1542
		1543
	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. Superglue: A stickier benchmark for general-purpose language understanding systems . In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	1544
		1545
		1546
		1547
		1548
		1549
		1550
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	1551
		1552
		1553
		1554
		1555
		1556
		1557
	Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2021. Adversarial GLUE: A multi-task benchmark for robustness evaluation of language models . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	1558
		1559
		1560
		1561
		1562
		1563
		1564
		1565
	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. How far can camels go? exploring the state of instruction tuning on open resources . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1566
		1567
		1568
		1569
		1570
		1571
		1572
		1573
		1574
	Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krma Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages	1575
		1576
		1577
		1578
		1579
		1580
		1581
		1582
		1583
		1584
		1585
		1586
		1587
		1588
		1589
		1590

1591	5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1648
1592		1649
1593	Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohanney, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The benchmark of linguistic minimal pairs for English . <i>Transactions of the Association for Computational Linguistics</i> , 8:377–392.	1650
1594		1651
1595		1652
1596		1653
1597		1654
1598		1655
1599	Bonnie Webber, Rashmi Prasad, Alan Lee, and Aravind Joshi. 2019. The penn discourse treebank 3.0 annotation manual . <i>Philadelphia, University of Pennsylvania</i> .	1656
1600		1657
1601		1658
1602		1659
1603	Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Ninanwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, et al. 2013. Ontonotes release 5.0 . <i>Linguistic Data Consortium, Philadelphia, PA</i> , 23:170.	1660
1604		1661
1605		1662
1606		1663
1607		1664
1608	Aaron Steven White, Drew Reisinger, Keisuke Sakaguchi, Tim Vieira, Sheng Zhang, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. 2016. Universal decompositional semantics on Universal Dependencies . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1713–1723, Austin, Texas. Association for Computational Linguistics.	1665
1609		1666
1610		1667
1611		1668
1612		1669
1613		1670
1614		1671
1615		1672
1616	Zhiyong Wu, Yun Chen, Ben Kao, and Qun Liu. 2020. Perturbed masking: Parameter-free probing for analyzing and interpreting BERT . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4166–4176, Online. Association for Computational Linguistics.	1673
1617		1674
1618		1675
1619		1676
1620		1677
1621		1678
1622	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions . <i>CoRR</i> , abs/2304.12244.	1679
1623		1680
1624		1681
1625		1682
1626		1683
1627	Prateek Yadav, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023. Compeft: Compression for communicating parameter efficient updates via sparsification and quantization . <i>CoRR</i> , abs/2311.13171.	1684
1628		1685
1629		1686
1630		1687
1631	Linyi Yang, Shuibai Zhang, Libo Qin, Yafu Li, Yidong Wang, Hanmeng Liu, Jindong Wang, Xing Xie, and Yue Zhang. 2023. GLUE-X: Evaluating natural language understanding models from an out-of-distribution generalization perspective . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 12731–12750, Toronto, Canada. Association for Computational Linguistics.	1688
1632		1689
1633		1690
1634		1691
1635		1692
1636		1693
1637		1694
1638		1695
1639	Amir Zeldes. 2017. The GUM corpus: Creating multilayer resources in the classroom . <i>Language Resources and Evaluation</i> , 51(3):581–612.	1696
1640		1697
1641		1698
1642	Xikun Zhang, Deepak Ramachandran, Ian Tenney, Yanai Elazar, and Dan Roth. 2020. Do language embeddings capture scales? In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP</i> , pages 292–299, Online. Association for Computational Linguistics.	1699
1643		1700
1644		1701
1645		1702
1646		1703
1647		1704
	Yian Zhang, Alex Warstadt, Xiaocheng Li, and Samuel R. Bowman. 2021. When do you need billions of words of pretraining data? In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 1112–1125, Online. Association for Computational Linguistics.	1705
		1706
		1707
		1708
		1709
		1710
		1711
		1712
		1713
		1714
		1715
		1716
		1717
		1718
		1719
		1720
		1721
		1722
		1723
		1724
		1725
		1726
		1727
		1728
		1729
		1730
		1731
		1732
		1733
		1734
		1735
		1736
		1737
		1738
		1739
		1740
		1741
		1742
		1743
		1744
		1745
		1746
		1747
		1748
		1749
		1750
		1751
		1752
		1753
		1754
		1755
		1756
		1757
		1758
		1759
		1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
		1769
		1770
		1771
		1772
		1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780
		1781
		1782
		1783
		1784
		1785
		1786
		1787
		1788
		1789
		1790
		1791
		1792
		1793
		1794
		1795
		1796
		1797
		1798
		1799
		1800

A Additional Details of HOLMES

A.1 Additional Details on the Evolution of Probing Literature

We analyze publication trends by year and venue as shown in Table 2. Less work was published between 2015-2018 (*earlier*) focusing on LSTM-based (Linzen et al., 2016; Conneau et al., 2018) and static LMs (Köhn, 2015; Linzen et al., 2016; Belinkov et al., 2017; Conneau et al., 2018). With the release of BERT (Devlin et al., 2019) in 2019, we note increasing attention to analyzing linguistic abilities within LMs, with a peak of 90 papers in 2022.⁵ Considering the venue, more than half of the relevant work (149 papers) was published at major conferences (ACL and EMNLP), and 68 papers were published at ACL, EACL, NAACL, and COLING.⁶ In addition, we observe a constant contribution of TACL, various workshops, such as Analyzing and Interpreting Neural Networks for NLP or Representation Learning for NLP.

A.2 Experimental Details

Probing Hyperparameters Following previous work (Hewitt and Liang, 2019; Voita and Titov, 2020), we use fixed hyperparameters for training the probes: 20 epochs, where we find the best one using dev instances; AdamW (Loshchilov and Hutter, 2019) as optimizer; a batch size of 64; a learning rate of 0.0005; a dropout rate of 0.2; a warmup rate of 10% of the steps; random seeds: [0, 1, 2, 3, 4]

Hardware We run all of our experiments using 12 Nvidia RTX A6000 GPUs. Every GPU provides 48GB of memory and 10752 CUDA Cores.

Considered LMs Table 8 outlines the details of the LMs we evaluate on HOLMES in this work.

A.3 Linguistic Task Categorization

We show in Table 3, Table 4, Table 7, Table 5, and Table 6 which resources HOLMES use to cover *morphology*, *syntax*, *semantics*, *reasoning*, and *discourse* phenomena. This includes 33 works providing the data, the specific linguistic phenomena, or both. For example, for *readability* we use the data of Weischedel et al. (2013) and calculated the flesch score (Flesch, 1948).

⁵Note that EMNLP-23 and ACL-23 proceedings were not published when conducting this meta-analysis.

⁶Note that EMNLP-23 and ACL-23 proceedings were not published when conducting this meta-study.

	<i>earlier</i>	2019	2020	2021	2022	2023	Total
ACL	2	10	12	9	34	25	92
AAACL	-	-	-	-	1	-	1
COLING	-	-	10	-	9	-	19
EACL	-	-	-	7	-	15	22
EMNLP	2	4	13	17	21	-	57
NAACL	-	3	-	9	14	-	26
TACL	1	1	2	3	3	1	11
Workshops	4	4	10	10	7	1	36
Other	1	2	1	1	1	4	10
Probing	10	24	48	56	90	46	274
All Papers	8,056	3,111	3,822	4,294	5,133	3,647	28,063

Table 2: Evolution of probing studies. Note that EMNLP-23 and ACL-23 proceedings were not published when conducting this meta-study.

Morphology First, we feature 19 tasks verifying *morphology* phenomena: *Anaphor agreement*, *determiner noun agreement*, *subject-verb agreement* and *irregular forms* (Warstadt et al., 2020; Huebner et al., 2021).

Syntax The second group of 75 tasks verifies the following *syntax* phenomena: *Part-of-speech tagging* and *constituent labeling* (Weischedel et al., 2013); *dependency labeling* (Silveira et al., 2014); *bigram-shift*, *tree-depth*, *top-constituent-task*, and *sentence-length* (Conneau et al., 2018); *subject- & object-number*, and *deoncausative-inchoative alternation* based on Klafka and Ettinger (2020); *binding*, *control/raising*, *negative polarity item licensing*, *island-effects*, *argument-structure*, *ellipsis*, and *filler-gap* (Warstadt et al., 2020; Huebner et al., 2021).

Semantics Third, consider 67 tasks covering *semantics* phenomena: *Named-entity labeling* and *semantic-role labeling* (Weischedel et al., 2013); *subject- and object-number*, *tense*, *semantic odd man out*, *word content*, and *coordination inversion* (Conneau et al., 2018); *semantic relation classification* (Hendrickx et al., 2010); *semantic proto-roles* (Rudinger et al., 2018a); *factuality* (Rudinger et al., 2018b); *genericity* (Govindarajan et al., 2019); *event structure* (Gantt et al., 2022); *time* (Vashishtha et al., 2019); *word sense* (White et al., 2016); *sentiment analysis* (Socher et al., 2013); *object- and subject-animacy*, *object- and subject-gender*, *verb-tense*, and *verb-dynamic* Klafka and Ettinger (2020); *metaphor* (Mohler et al., 2016; Birke and Sarkar, 2006; Steen et al., 2010); *complex word identification* (Paetzold and Specia, 2016); and *passive* (Krasnowska-Kieraś and Wróblewska, 2019). In addition, we derive

1765 *synonym-/antonym-detection* task using WordNet
1766 (Miller, 1995) and the texts from OntoNotes v5
1767 Weischedel et al. (2013).

1768 **Reasoning** Forth, 19 tasks cover *reasoning*
1769 phenomena: *Paraphrasticity* with negation and
1770 antonyms (Vahtola et al., 2022); *negation detec-*
1771 *tion* (Szarvas et al., 2008; Konstantinova et al.,
1772 2012; Morante and Blanco, 2012); *negation-span*
1773 *classification* (Szarvas et al., 2008; Konstantinova
1774 et al., 2012); *negation-correspondence* (Szarvas
1775 et al., 2008; Konstantinova et al., 2012); *specula-*
1776 *tion detection*, *speculation-span classification*, and
1777 *speculation-correspondence* (Szarvas et al., 2008);
1778 and *always-never*, *age comparison*, *objects com-*
1779 *parison*, *antonym negation*, *property conjunction*,
1780 *taxonomy connection*, *encyclopedic composition*,
1781 and *multi-hop composition* (Talmor et al., 2020).

1782 **Discourse** Finally, Holmes embodies 28 task
1783 addressing *discourse* phenomena: *Co-reference*
1784 *resolution* Weischedel et al. (2013); *bridging*
1785 (Hou, 2018, 2020; Pandit and Hou, 2021); *dis-*
1786 *course connective* (Nie et al., 2019); *sentence or-*
1787 *der* and *next-sentence prediction* (Narayan et al.,
1788 2018); *discourse correspondence*, *discourse or-*
1789 *der*, *discourse relation*, *discourse distance*, *dis-*
1790 *course explicit classes*, *discourse implicit classes*
1791 (Webber et al., 2019; Kurfali and Östling, 2021);
1792 and *rst-count/-depth/-distance/-relation/-relation-*
1793 *group/-successively/-type* (Carlson et al., 2001;
1794 Koto et al., 2021; Kurfali and Östling, 2021; Zeldes,
1795 2017).

1796 A.4 Details of Probing Dataset Composition

1797 Whenever possible, we rely on established prob-
1798 ing datasets and transform instances into a unified
1799 format: **1**) an input x which is either one or a pair
1800 of span(s) or sentence(s), including the string and
1801 an optional starting and ending index in the con-
1802 text c when task type is either a span or span-pair
1803 classification; **2**) an optional textual context c to
1804 encode x , for example the sentence in which a span
1805 occurs; and **3**) a corresponding label y . If given,
1806 we use the original train/dev/test splits. However,
1807 if this division does not exist, we use a 70/10/20
1808 ratio to form these splits. Furthermore, we adapt
1809 the design of some tasks to map to our task format.
1810 Exemplary, for the oLMmpics (Talmor et al., 2020)
1811 dataset, we transform the mask-filling tasks into a
1812 binary classification where the *correct* label corre-
1813 sponds to a sentence with a correctly filled mask

and *incorrect* to a sentence where the mask was
1814 filled wrongly. 1815

OnToNotes Following Tenney et al. (2019b,a),
1816 we use the *OntoNotes* (Weischedel et al., 2013)
1817 dataset to derive *part-of-speech tagging*, *con-*
1818 *stituent labeling*, *named-entity labeling*, *semantic*
1819 *role*, and *co-reference resolution* probing datasets.
1820 Further, we consider with *constituent maximum*
1821 *depth* and *constituent node length* further proper-
1822 ties of the constituent tree this dataset *OntoNotes*. 1823

Dependency Corpus As in Tenney et al.
1824 (2019b,a), we use Universal Dependencies anno-
1825 tations of the English Web Treebank to form a
1826 *dependency labeling* datasets. 1827

Context Probes Presented in Klafka and Ettinger
1828 (2020), we compose nine datasets to verify infor-
1829 mation about context words. 1830

BLiMP Dataset Using the data presented in the
1831 BLiMP benchmark (Warstadt et al., 2020), we de-
1832 rive 67 probing datasets verifying specific phenom-
1833 ena, like *island effect*, covering *morphology*, *syn-*
1834 *tax*, and *semantics*. Unlike the original version,
1835 we compose a binary classification task for every
1836 phenomenon. Precisely, whether to accept or reject
1837 a given sentence, where rejecting means that the
1838 given linguistic phenomena is violated. 1839

Zorro Dataset As for the BLiMP tasks, we con-
1840 vert the 21 distinct Zorro tasks into a binary classi-
1841 fication task on whether a sentence accepts or rejects
1842 the given linguistic phenomena is violated. 1843

SemEval-2010 Task 8 For *semantic relation*
1844 *classification* we rely on the dataset of Hendrickx
1845 et al. (2010). 1846

Decompositional Semantics Initiative The *De-*
1847 *compositional Semantics Initiative*⁷ provides a
1848 large number of datasets to verify semantic phe-
1849 nomena. Apart of the common use *semantic proto-*
1850 *roles* (Rudinger et al., 2018a), we use their collec-
1851 tion of works to compose probing datasets for *fac-*
1852 *tuality* (Rudinger et al., 2018b), genericity (Govin-
1853 darajan et al., 2019), event structure (Vashishtha
1854 et al., 2019), time (Vashishtha et al., 2019), and
1855 word sense (White et al., 2016). 1856

Sentiment Analysis We use the commonly used
1857 work of Socher et al. (2013) and form a probing
1858 dataset targeting sentiment. 1859

⁷<https://decomp.io/>

Metaphor As in Aghazadeh et al. (2022), we use the data from Mohler et al. (2016); Birke and Sarkar (2006); Steen et al. (2010) to form three metaphor datasets.

Complex Word Identification We consider word complexity for the first time and use the data presented in Paetzold and Specia (2016). It provides annotations for different complexity levels of words.

Passive We use data from Krasnowska-Kieraś and Wróblewska (2019) to form a probing dataset assessing knowledge about passive language.

Synonym / Antonym Replacement Using the text of the *OntoNotes* (Weischedel et al., 2013) and Wordnet (Miller, 1995), we form a probing dataset to detect synonym and antonym replacement. Specifically, the binary classification task is: given two texts (the original and an updated one), was the updated one changed by replacing a word with its synonym or antonym?

Negation With this work, we verify for the first time *negation* based on human annotated datasets (Vahtola et al., 2022; Szarvas et al., 2008; Konstantinova et al., 2012). Specifically, we form different probing datasets.

- Is a text negated or not?
- Given two text spans, does the negation within the first one correspond to the second one?
- Given a text span, is it the cue or the scope of the negation?

oLMpics We form probing datasets addressing different lexical reasoning using the data presented in Talmor et al. (2020). As they provide multiple choices, we form *correct* instances by filling the gap with the correct option and *wrong* ones by filling in the other options. Specifically, we form dataset for *always-never*, *age comparison*, *objects comparison*, *antonym-negation*, *multi-hop composition property conjunction*, *taxonomy conjunction*, and *encyclopedic composition*.

Bridging We rely on the data presented in Pandit and Hou (2021) and form two probing datasets. One is to verify whether a text is linguistically applicable, considering bridging (antecedent matches anaphora). And a second one to verify whether an antecedent and anaphora match.

Discourse Connective Using data from Nie et al. (2019), we form a probing dataset to assess whether a given connective marker matches the discourse of the given text.

Sentence Order and Next Sentence Prediction Following Narayan et al. (2018), we form two datasets to verify the order of good or badness of a given sentence and whether two sentences occur after each other.

Discourse Representation Theory We use data from Webber et al. (2019) to compose eight probing datasets addressing *discourse representation theory*:

- Four probing dataset predicting the class of a given span. We distinguish between *implicit*, *explicit*, *implicit-coarse*, and *explicit-coarse*.
- The absolute distance, number of words, between two spans in the text.
- Whether the order of two spans is correct or not.
- Whether two spans have discourse relation or not.
- The specific discourse relation of two spans.

Rhetorical Structure Theory Using annotations from Carlson et al. (2001); Zeldes (2017), we compose 14 probing datasets addressing *rhetorical theory*. Specifically, we compose the following seven types of datasets for both works:

- The rhetorical type of a text span, either nucleus or satellite.
- The number of children of a text span within the rhetorical tree of the text.
- The depth of a text span within the rhetorical tree of the text.
- The number of edges between two text spans within the rhetorical tree.
- The specific rhetorical relation between two text spans like *conclusion*.
- The relation group of a specific rhetorical relation between two text spans like *evaluation* for the relation *conclusion*.
- Whether two text spans occur after each other in the rhetorical tree.

Phenomena	Text	Text-Pair	Span	Span-Pair	Warstadt et al. (2020)	Huebner et al. (2021)
	<i>anaphor agreement</i>	3				✓
<i>determiner noun agreement</i>	10				✓	✓
<i>irregular forms</i>	3				✓	✓
<i>subject-verb agreement</i>	10				✓	✓

Table 3: Overview of resources and linguistic phenomena mapping for *morphology*. It shows the number of datasets for the phenomena by dataset type.

Phenomena	Text	Text-Pair	Span	Span-Pair	Vahtola et al. (2022)	Szarvas et al. (2008)	Konstantinova et al. (2012)	Morante and Blanco (2012)	Talmor et al. (2020)
	<i>age comparison</i>	1							
<i>always-never</i>	1								✓
<i>antonym negation</i>	1								✓
<i>encyclopedic composition</i>	1								✓
<i>multi-hop composition</i>	1								✓
<i>negation</i>	3	1	2	2	✓	✓	✓	✓	
<i>objects comparison</i>	1								✓
<i>property conjunction</i>	1								✓
<i>speculation</i>	1		1	1		✓			
<i>taxonomy connection</i>	1								✓

Table 5: Overview of resources and linguistic phenomena mapping for *reasoning*. It shows the number of datasets for the phenomena by dataset type.

Phenomena	Text	Text-Pair	Span	Span-Pair	Weischedel et al. (2013)	Silveira et al. (2014)	Comneau et al. (2018)	Fleisch (1948)	Klafka and Ertinger (2020)	Warstadt et al. (2020)	Huebner et al. (2021)
	<i>argument-structure</i>	20									✓
<i>bigram-shift</i>	1						✓				
<i>binding</i>	8									✓	✓
<i>case</i>	1										✓
<i>constituent parsing</i>	2		1		✓						
<i>control/raising</i>	5									✓	✓
<i>deoncausative-inchoative alternation</i>	1								✓		
<i>dependency parsing</i>			1		✓						
<i>ellipsis</i>	3									✓	✓
<i>filler-gap</i>	9									✓	✓
<i>island-effects</i>	10									✓	✓
<i>local attractor</i>	1										✓
<i>object-number</i>	2									✓	
<i>part-of-speech</i>		3			✓	✓	✓				
<i>readability</i>	1				✓			✓			
<i>sentence-length</i>	1										
<i>subject-number</i>	2						✓			✓	
<i>top-constituent-task</i>	1						✓				
<i>tree-depth</i>	1						✓				

Table 4: Overview of resources and linguistic phenomena mapping for *syntax*. It shows the number of datasets for the phenomena by dataset type.

Phenomena	Text	Text-Pair	Span	Span-Pair	Weischedel et al. (2013)	Pandit and Hou (2021)	Nie et al. (2019)	Narayan et al. (2018)	Webber et al. (2019)	Carlson et al. (2001)	Zeldes (2017)
	<i>bridging</i>	1			1		✓				
<i>co-reference resolution</i>				1	✓						
<i>discourse connective</i>		1					✓				
<i>discourse representation theory</i>				8					✓		
<i>next-sentence prediction</i>		1					✓				
<i>rethorical structure theory</i>			6	8						✓	✓
<i>sentence order</i>		1					✓				

Table 6: Overview of resources and linguistic phenomena mapping for *discourse*. It shows the number of datasets for the phenomena by dataset type.

Phenomena	Text	Text-Pair	Span	Span-Pair	Weischedel et al. (2013)	Conneau et al. (2018)	Klafka and Eittinger (2020)	Wärstadt et al. (2020)	Huebner et al. (2021)	Hendrickx et al. (2010)	Rudinger et al. (2018a)	Rudinger et al. (2018b)	Govindarajan et al. (2019)	Gantt et al. (2022)	Vashishtha et al. (2019)	White et al. (2016)	Socher et al. (2013)	Mohler et al. (2016)	Birke and Sarkar (2006)	Steen et al. (2010)	Paetzold and Specia (2016)	Krasnowska-Kieras and Wroblewska (2019)	Miller (1995)
	Text	Text-Pair	Span	Span-Pair																			
<i>complex word identification</i>			1																			✓	
<i>coordination inversion</i>	1				✓																		
<i>event structure</i>		4	2												✓								
<i>factuality</i>				1								✓											
<i>genericity</i>		6											✓										
<i>metaphor</i>		4																	✓	✓	✓		
<i>named-entity labeling</i>		1			✓																		
<i>negative polarity item licensing</i>	4							✓	✓														
<i>object-animacy</i>	1						✓																
<i>object-gender</i>	1						✓																
<i>passive</i>	1																					✓	
<i>quantifiers</i>	6								✓														
<i>semantic relation classification</i>		1								✓													
<i>semantic proto-roles</i>			20								✓												
<i>semantic odd man out</i>	1				✓																		
<i>semantic-role labeling</i>			1		✓																		
<i>sentiment analysis</i>	1																✓						
<i>subject-animacy</i>	1						✓																
<i>subject-gender</i>	1						✓																
<i>synonym-/antonym-detection</i>	1																						✓
<i>tense</i>	2					✓	✓																
<i>time</i>		1													✓								
<i>verb-dynamic</i>	1						✓																
<i>word content</i>	1					✓																	
<i>word sense</i>			1													✓							

Table 7: Overview of resources and linguistic phenomena mapping for *semantics*. It shows the number of datasets for the phenomena by dataset type.

Model	Citation	Size	Pre-Training Objective	Pre-Training Data	Huggingface Tag
<i>Encoder-Only Language Models</i>					
ALBERT	Lan et al. (2020)	10 million	MLM+SOP	16GB	albert-base-v2
BERT	Tenney et al. (2019a)	110 million	MLM+NSP	16GB	bert-base-uncased
DeBERTa	He et al. (2021)	100 million	MLM	80GB	microsoft/deberta-base
DeBERTa-v3	He et al. (2023)	86 million	MLM+DISC	160GB	microsoft/deberta-v3-base
ELECTRA	Clark et al. (2020)	110 million	MLM	16GB	google/electra-base-discriminator
RoBERTa	Liu et al. (2019)	110 million	MLM+DISC	160GB	roberta-base
<i>Decoder-Only Language Models</i>					
GPT2	Radford et al. (2019)	117 million	LM	40GB	gpt2
Pythia-70m	Biderman et al. (2023)	70 million	LM	300 billion tokens	EleutherAI/pythia-70m
Pythia-160m	Biderman et al. (2023)	160 million	LM	300 billion tokens	EleutherAI/pythia-160m
Pythia-410m	Biderman et al. (2023)	410 million	LM	300 billion tokens	EleutherAI/pythia-410m
Pythia-1b	Biderman et al. (2023)	1 billion	LM	300 billion tokens	EleutherAI/pythia-1b
Pythia-1.4b	Biderman et al. (2023)	1.4 billion	LM	300 billion tokens	EleutherAI/pythia-1.4b
Pythia-2.8b	Biderman et al. (2023)	2.8 billion	LM	300 billion tokens	EleutherAI/pythia-2.8b
Pythia-6.9b	Biderman et al. (2023)	6.9 billion	LM	300 billion tokens	EleutherAI/pythia-6.9b
Pythia-12b	Biderman et al. (2023)	12 billion	LM	300 billion tokens	EleutherAI/pythia-12b
Pythia-70m-dedup	Biderman et al. (2023)	70 million	LM	207 billion tokens	EleutherAI/pythia-70m-dedup
Pythia-160m-dedup	Biderman et al. (2023)	160 million	LM	207 billion tokens	EleutherAI/pythia-160m-dedup
Pythia-410m-dedup	Biderman et al. (2023)	410 million	LM	207 billion tokens	EleutherAI/pythia-410m-dedup
Pythia-1b-dedup	Biderman et al. (2023)	1 billion	LM	207 billion tokens	EleutherAI/pythia-1b-dedup
Pythia-1.4b-dedup	Biderman et al. (2023)	1.4 billion	LM	207 billion tokens	EleutherAI/pythia-1.4b-dedup
Pythia-2.8b-dedup	Biderman et al. (2023)	2.8 billion	LM	207 billion tokens	EleutherAI/pythia-2.8b-dedup
Pythia-6.9b-dedup	Biderman et al. (2023)	6.9 billion	LM	207 billion tokens	EleutherAI/pythia-6.9b-dedup
Pythia-12b-dedup	Biderman et al. (2023)	12 billion	LM	207 billion tokens	EleutherAI/pythia-12b-dedup
Dolly-v2	Conover et al. (2023)	12 billion	LM+IT	300 billion token + 15K instructions	databricks/dolly-v2-12b
Llama-2-7b	Touvron et al. (2023)	7 billion	LM	2.4 trillion tokens	meta-llama/Llama-2-7b-hf
Llama-2-13b	Touvron et al. (2023)	13 billion	LM	2.4 trillion tokens	meta-llama/Llama-2-13b-hf
Llama-2-70b	Touvron et al. (2023)	70 billion	LM	2.4 trillion tokens	meta-llama/Llama-2-70b-hf
Llama-2-7b-chat	Touvron et al. (2023)	7 billion	LM+IT	2.4 trillion tokens + 27.5K instructions	meta-llama/Llama-2-7b-chat-hf
Llama-2-13b-chat	Touvron et al. (2023)	13 billion	LM+IT	2.4 trillion tokens + 27.5K instructions	meta-llama/Llama-2-13b-chat-hf
Llama-2-70b-chat	Touvron et al. (2023)	70 billion	LM+IT	2.4 trillion tokens + 27.5K instructions	meta-llama/Llama-2-70b-chat-hf
IBM-Merlinite	Sudalairaj et al. (2024)	7 billion	LM+IT	2.4 trillion tokens + 1400k instructions	ibm/merlinite-7b
IBM-Labradorite	Sudalairaj et al. (2024)	13 billion	LM+IT	2.4 trillion tokens + 1400k instructions	ibm/labradorite-13b
Vicuna-13b-v1.5	Zheng et al. (2023)	13 billion	LM+IT	2.4 trillion tokens + 125k instructions	lmsys/vicuna-13b-v1.5
Orca-2-13b	Mitra et al. (2023)	13 billion	LM+IT	2.4 trillion tokens + 817K instructions	microsoft/Orca-2-13b
Wizard-13B-v1.2	Xu et al. (2023)	13 billion	LM	unknown	WizardLM/WizardLM-13B-V1.2
Tulu-2-13b	Wang et al. (2023)	13 billion	LM+IT	2.4 trillion tokens + 330k instructions	allenai/tulu-2-13b
Tulu-2-dpo-13b	Wang et al. (2023)	13 billion	LM+IT	2.4 trillion tokens + 330k instructions	tulu-2-dpo-13b
Tulu-2-70b	Wang et al. (2023)	70 billion	LM+IT	2.4 trillion tokens + 330k instructions	allenai/tulu-2-70b
Tulu-2-dpo-70b	Wang et al. (2023)	70 billion	LM+IT	2.4 trillion tokens + 330k instructions	tulu-2-dpo-70b
Mistral-7b	Jiang et al. (2023)	7 billion	LM	unknown	mistralai/Mistral-7B-v0.1
Mistral-7b-Inst	Jiang et al. (2023)	7 billion	LM	unknown	mistralai/Mistral-7B-Instruct-v0.1
Mixtral-8x7b	Jiang et al. (2024)	47 billion	LM	unknown	mistralai/Mixtral-8x7B-v0.1
Mixtral-8x7b-Inst	Jiang et al. (2024)	47 billion	LM	unknown	mistralai/Mistral-7B-v0.1
<i>Encoder-Decoder Language Models</i>					
BART	Lewis et al. (2020)	121 million	DAE	160GB	google/facebook/bart-base
T5-small	Raffel et al. (2020)	60 million	DAE	800GB	google/t5-small-lm-adapt
T5-base	Raffel et al. (2020)	220 million	DAE	800GB	google/t5-base-lm-adapt
T5-large	Raffel et al. (2020)	770 million	DAE	800GB	google/t5-large-lm-adapt
T5-xl	Raffel et al. (2020)	3 billion	DAE	800GB	google/t5-xl-lm-adapt
T5-xxl	Raffel et al. (2020)	11 billion	DAE	800GB	google/t5-xxl-lm-adapt
FLAN-T5-small	Raffel et al. (2020)	60 million	DAE+IT	800GB + 1.8k tasks	google/t5-small-lm-adapt
FLAN-T5-base	Raffel et al. (2020)	220 million	DAE+IT	800GB + 1.8k tasks	google/t5-base-lm-adapt
FLAN-T5-large	Raffel et al. (2020)	770 million	DAE+IT	800GB + 1.8k tasks	google/t5-large-lm-adapt
FLAN-T5-xl	Raffel et al. (2020)	3 billion	DAE+IT	800GB + 1.8k tasks	google/t5-xl-lm-adapt
FLAN-T5-xxl	Raffel et al. (2020)	11 billion	DAE+IT	800GB + 1.8k tasks	google/t5-xxl-lm-adapt
TK-Instruct	Wang et al. (2022)	11 billion billion	DAE+IT	800GB + 1.6k tasks	allenai/tk-instruct-11b-def
UL2	Tay et al. (2023)	20 billion	DAE	800GB	google/ul2
FLAN-UL2	Tay et al. (2023)	20 billion	DAE+IT	800GB + 100k instructions	google/flan-ul2
<i>Static Language Models</i>					
Glove-6B	Pennington et al. (2014)	-	WP	6 billion tokens	glove.6B.300d
Glove-840B	Pennington et al. (2014)	-	WP	840 billion tokens	glove.840B.300d

Table 8: Overview of the evaluated LMS covering the corresponding citation, model size, model architecture, pre-training objective & data, and the Huggingface model tag. Regarding the pre-training objective, we distinguish between masked language modeling (MLM), sentence order prediction (SOP), next sentence prediction (NSP), next word prediction (LM), instruction fine-tuning (IT), word denoising (DAE), and word probabilities from word co-occurrences (WP). For pre-training data, we report known numbers, either as the size of the corpora in gigabytes (GB), the number of pre-training tokens, the number of instructions for fine-tuning, or the number of tasks for instruction fine-tuning.