

Supplementary information

Detecting hallucinations in large language models using semantic entropy

In the format provided by the authors and unedited

1887 Supplementary Material

1888

1889

1890 Note 1: Worked Example of Semantic Entropy Calculation

1891

1892 This note provides a worked example of the calculation of semantic entropy. As a

1893

1894 worked example, suppose that we have asked “Where is the Eiffel Tower?”. The

1895 model generates five answers with the length-normalised sequence log-probabilities

1896

1897 $\frac{1}{N} \prod_i^N p(\mathbf{s}_i | \mathbf{s}_i)$ (given in the first column). In this hypothetical example, we happened

1898

1899 to sample the literal string “Paris.” twice, because we are just randomly sampling

1900 from the language model which assigns the string high probability. But we also found

1901

1902 a different string that was equivalent to it, as well as some wrong answers. Note

1903

1904 that in the case of model APIs which do not report the log-probabilities (such as

1905 GPT-4 at time of writing) we will not have the final number, just the text. To

1906

1907 compute the semantic entropy, we cluster these generations into clusters that can

1908

1909 be considered to mean the same thing. We also add up the probabilities associated

1910

1911 Generation	1912 p	1911 Naive entropy			1911 Semantic entropy			
		1912 $p(\mathbf{s}_i)$	1912 $\log[p(\mathbf{s}_i)]$	1912 $p(\mathbf{s}_i) \log[p(\mathbf{s}_i)]$	1912 $\sum_{\mathbf{s}_i \in C_j} p(\mathbf{s}_i)$	1912 $p(C_j)$	1912 $\log[p(C_j)]$	1912 $p(C_j) \log[p(C_j)]$
1913 “Paris.”	0.20	0.33	-0.48	-0.16	0.55	0.90	-0.04	-0.04
1914 “Paris.”	0.20	0.33	-0.48	-0.16	-	-	-	-
1915 “It’s Paris.”	0.15	0.25	-0.61	-0.15	-	-	-	-
1916 “Rome.”	0.05	0.08	-1.09	-0.09	0.05	0.08	-1.09	-0.09
1917 “New York.”	0.01	0.02	-1.79	-0.03	0.01	0.02	-1.79	-0.03
Sum	0.61	1.00	-4.45	-0.59	0.61	1.00	-2.92	-0.16

1918 Supplementary Table 1: Worked example of Semantic Entropy Calculation.

1919 The raw token sequence probabilities for each generation, p , are in the first column.

1920 Note that they do not sum to one because they are the probabilities associated with

1921 each actually sampled outcome, and if we sample many generations their sum will

1922 exceed one. To calculate the naive entropy of the output distribution, in the second

1923 column we compute an estimator of the normalised probability for each generated

1924 sequence, $p(\mathbf{s}_i)$, by dividing each probability by the sum of the first column. These now

1925 do sum to one (up to a rounding error). One way to estimate the naive entropy would

1926 then be to multiply each log-probability by the probability, sum them (and multiply by

1927 -1, not shown). For the semantic entropy, we instead look at the probabilities summed

1928 up within each meaning-cluster. We compute log-probabilities in the same way, and

1929 then compute the entropy of the resulting distribution. The resulting entropy is much

1930 lower, because several generations meant the same thing as each other (final column).

1931 All values to two decimal places.

1932

with each one. This value is reported in the fifth column (the first under the heading “Semantic entropy”). To compute the semantic entropy we compute the negative sum of the expectation of the log-probabilities (Eq. (5)). That is, we get the result $0.16 = -0.9 \log(0.9) - 0.08 \log(0.08) - 0.02 \log(0.02)$.

For the “discrete” variant of semantic entropy we effectively treat the probability of sampling each of the generations as uniform, by using it as an empirical distribution that approximates the underlying distribution. This means that for the “Paris.” cluster we get a weight of $0.6=0.2+0.2+0.2$ and a weight of 0.2 for the other two clusters. That is, the discrete semantic entropy here is $0.41 = -0.6 \log(0.6) - 0.2 \log(0.2) - 0.2 \log(0.2)$.

Although these two methods produce different absolute results, we find that in practice they tend to agree fairly well on relative ordering, which is what is used in practice to classify confabulations. As a result, the cluster approximation of semantic entropy is a fairly good alternative in cases where the log-probabilities are not disclosed.

Note 2: Choosing an Entailment Estimator

Sentence-length Generations

We confirm that the bi-directional entailment classifier works as expected. Prior work has show that in some settings NLI methods can systematically fail⁷⁸, so we seek to ascertain whether these failures substantially affect typical question-answering. Two raters manually labeled 100 pairs of sentence-length generations from LLaMA 2 Chat 70B for three of our datasets for entailment, recording whether they believed that sentence A entailed sentence B. They rated each entailment as: **entailment**, **neutral**, **contradiction**. For the purpose of measuring agreement we combine the neutral and contradiction ratings, because our method is searching for positive entailment. We found that the human raters agreed with each other (87%) at roughly the same rate that they on average agreed with GPT-4 (87%) while they agreed with GPT-3.5 only

1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978

	DeBERTa	LLaMA 2 Chat 70B	GPT-3.5	GPT-4	Human A	Human B
Human A	0.81	0.78	0.83	0.89	-	0.87
Human B	0.78	0.80	0.84	0.85	0.87	-
Human Average	0.80	0.79	0.83	0.87	-	-

Supplementary Table 2: Manual entailment evaluation. Inter-rater agreement on entailment classification for pairs of sentence-length answers produced by LLaMA 2 Chat 70B to 100 questions from each of SQuAD, TriviaQA, and BioASQ (600 answers in total). On average, the human raters agreed with each other to approximately the same extent that they agreed with GPT-4, while GPT-3.5 was only slightly less predictive of human-assessed entailment.

slightly less on average (83%). As a result, because GPT-3.5 is more than an order of magnitude cheaper, we use GPT-3.5 for all entailment calculations for sentence-length generations on SQuAD, TriviaQA, BioASQ, SVAMP, and NQ Open. Presumably our method would perform better with a more expensive entailment estimator. Supplementary table 2 shows the detailed agreement results between the human raters and the entailment estimation.

In addition to validating human agreement with the entailment models, we also investigate the performance of semantic entropy with different entailment strategies. In supplementary table 3, we report AUROC values of semantic entropy with various entailment models for predictions from a LLaMA 2 Chat 70B model on TriviaQA, SQuAD, and BioASQ. The LLaMA 2 Chat 70B model performs worst, followed by the purpose-built DeBERTa model and GPT models, where this time version 3.5 slightly outperforms version 4 on average. For these experiments, unlike our main results, we used only 8 generations (normally 10) to estimate the entropy and measured accuracy relative to the reference answer using LLaMA 2 Chat 70B (normally GPT-4, see section on ‘Assessing Accuracy’).

	DeBERTa	LLaMA 2 Chat 70B	GPT-3.5	GPT-4
TriviaQA	0.83	0.70	0.85	0.83
SQuAD	0.76	0.71	0.77	0.80
BioASQ	0.75	0.73	0.87	0.79
Average	0.78	0.71	0.83	0.80

Supplementary Table 3: Entailment Method Ablation. AUROC values for semantic entropy when using different models to compute entailment for sentence-length generations from LLaMA 2 Chat 70B. Semantic entropy performs better when prompted GPT models predict entailment rather than a purpose-built DeBERTa model.

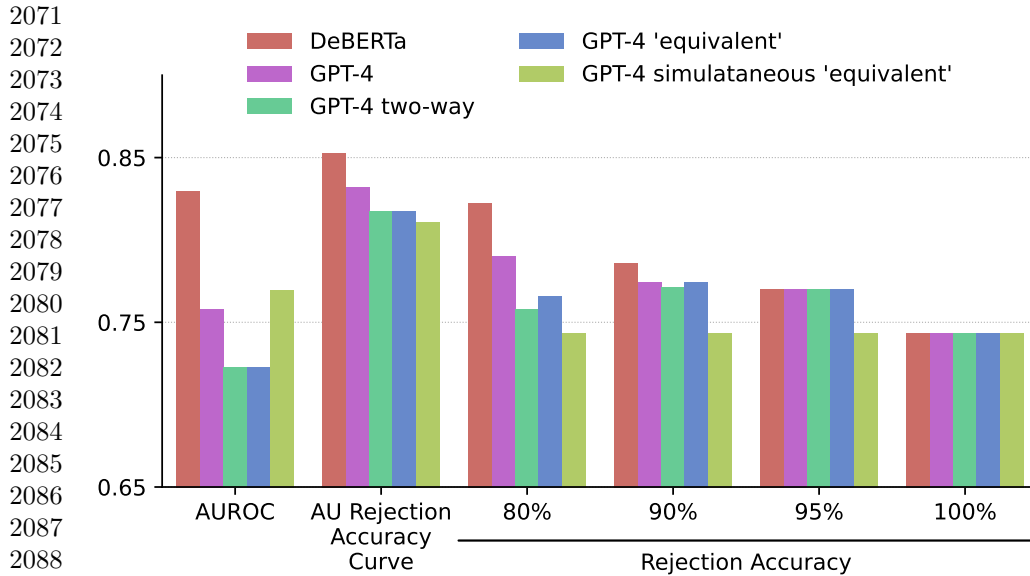
Paragraph-length Generations

In supplementary figure 1, we report our experiments for several entailment prediction variants for paragraph-length generations, in addition to our default non-defating bi-directional DeBERTa method. We also experimented with several entailment variants: “GPT-4 two-way” asks GPT-4 to evaluate whether the sentences mean the same thing directly (“Do the following two possible answers to the subquestion mean the same thing?” instead of “Does Possible Answer 1 semantically entail...”); “GPT-4 ‘equivalent’” instead asks “Are the following two possible answers to the subquestion semantically equivalent?”; while “GPT-4 simultaneous ‘equivalent’ ” provides all of the possible answers and asks “Are the following answers equivalent?”. All of these methods were substantially worse.

Note 3: Limitations to Clustering by Entailment

In idealized examples, it is clear when two sentences do or do not mean the same thing as each other. In practice, it can sometimes be that sentence A seems to mean the same as B, and B the same as C, but A and C don’t seem to mean the same thing. That is, because semantic equivalence is fuzzy it does not always intuitively behave transitively, meaning that the assumptions behind our equivalence classes do not hold in practice.

2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070



2090 **Supplementary Figure 1: Entailment method choice for paragraph biogra-**
 2091 **phies.** Implementing the non-defeating bi-directional entailment with DeBERTa
 2092 provided the best empirical results for paragraph biography confabulation detection.
 2093

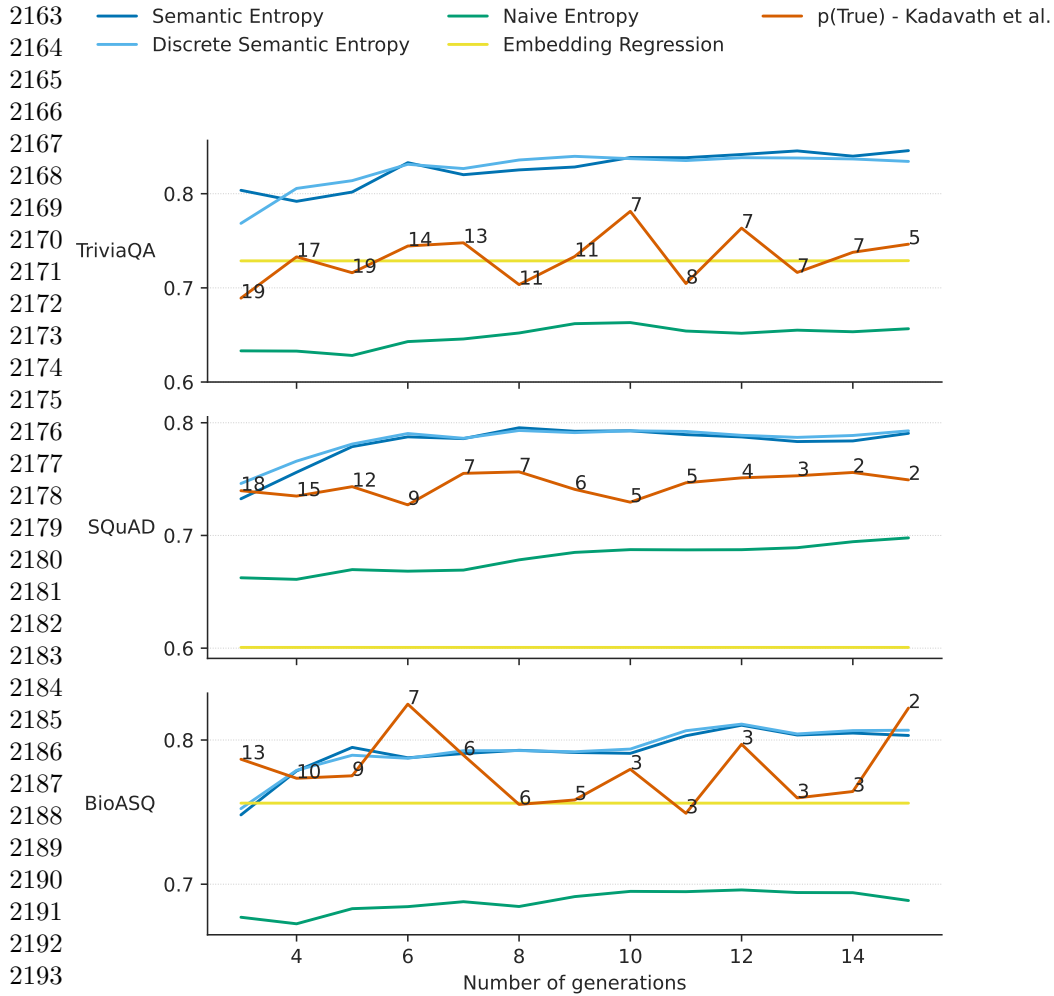
2094
 2095 Similarly, there are cases where bi-directional entailment does *not* mean that two
 2096 sentences mean the same thing. For example, “John drove his car to the store.” and
 2097 “John went to the store in his car.” generally imply each other and would be marked
 2098 as “entailment” by most classifiers, and this reflects the fact that they mean more-
 2099 or-less the same thing. But, for example, this depends somewhat on the context and
 2100 various aspects of implicature⁷⁶. For example, if we have other reasons to think that
 2101 John might have been the owner of the car but a passenger, rather than the driver,
 2102 then we might judge the two sentences to not mean the same thing as each other.
 2103 As an alternative example of a failure of bi-directional entailment to correspond with
 2104 semantic equivalence, sentences with scalar adverbs such as “Paris might be in France”
 2105 and “Paris might not be in France” can entail each other while meaning something
 2106 quite different⁷⁷.

For questions whose answers are straightforward, relatively objective, factual, and not vague these problems may not be significant. In particular, we did not observe any of these problems arising in manual inspection of outputs during any of our experiments. Nevertheless, for subtle situations and applications we would encourage practitioners to check these assumptions.

Note 4: Computational Cost and Choosing the Number of Generations

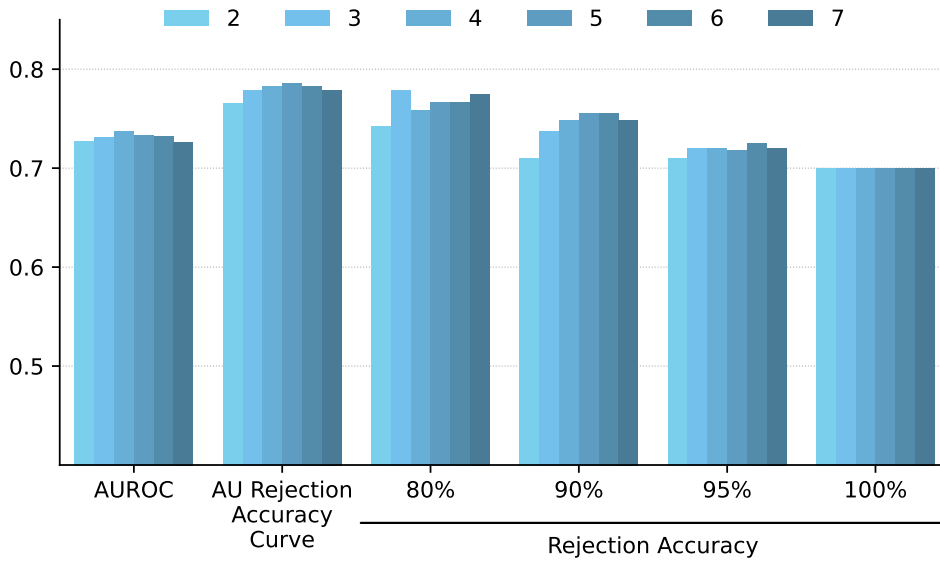
The bi-directional equivalence algorithm is combinatorially complex in M , the number of samples generated, as it requires $\binom{M}{2}$ -many comparisons in the worst-case. In practice, however, the computational cost is small compared to the cost of generating sequences.

First, M does not necessarily need to be very large. We show how the confabulation-detection performance (measured by AUROC) changes with M for sentence-length generations in supplementary figure 2 and for FactualBio paragraphs in supplementary figure 3. For sentence-length generations, after roughly $M = 5$ there are diminishing returns, although going up to $M = 10$ can still help. In this paper, we use $M = 10$ for sentence-length generations as well as short-phrase generations. For this ablation, we produce generations using LLaMA 2 Chat 70B but several experimental characteristics differ from those of our main results. We check entailment using GPT-4 (standardly GPT-3.5), measure accuracy using LLaMA 2 Chat 70B (standardly GPT-4), and use 8 generations to estimate entropy (standardly 10). For paragraph-length biographies, we find that four total factoids (three new generations plus the original factoid) seems optimal (see supplementary figure 3). Unlike the standard setting, more generations is not strictly better, because it decreases the relative weight on the original factoid which increases the risk of a badly posed question that generates irrelevant answers.



Supplementary Figure 2: Number of sentence-length generations used for entropy. We find diminishing returns to increasing the number of generations sampled for the semantic entropy estimation, but select 10 as a reasonable number for results in this paper. The numbers annotating $p(\text{True})$ illustrate the number of few-shot examples we were able to include without exceeding the maximal input size for each dataset and number of generations.

Second, when using the DeBERTa-large model, it is so much smaller than the main language model, each pair comparison is much faster than generating even one token from the main model. Using GPT-3.5 to do clustering is considerably more expensive than DeBERTa.

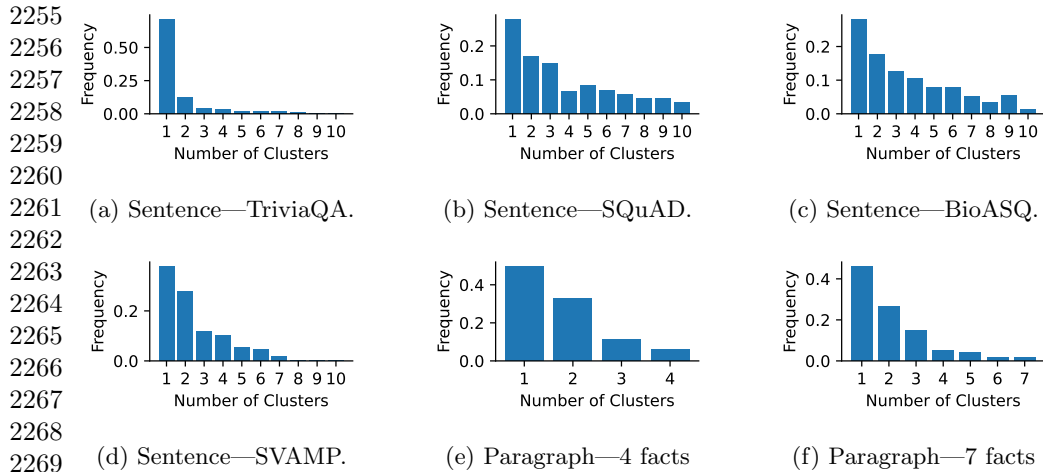


Supplementary Figure 3: Number of factoid generations used in paragraph-length biographies. We find that confabulation-detection performance is not very sensitive to the number of generations, but that four total factoids per question (including the original one) results in competitive performance.

Third, because semantic equivalence is transitive we only need to compare one member of each equivalence class to the remaining sequences (see algorithm in Extended Data Figure 1). The number of semantic clusters in our tasks is empirically often quite low which means that far fewer than the worst-case number of comparisons are actually needed in practice. In supplementary figure 4, we show some empirical numbers of clusters for several key datasets.

Fourth, because the LLM often generates identical sequences in practice, we can cache entailments. For example, if the LLM’s three generations in response to a question are “Paris.”, “It’s in Paris.”, and “Paris.” we can do a (very computationally cheap) *string level* comparison of the final “Paris.” to the previous generations and, on finding that it is the identical text to the earlier string, we can use the previously calculated entailments. We find that in practice this reduces the computational costs

2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254



2270 **Supplementary Figure 4:** In many cases, relatively few clusters are found, which
 2271 can improve computational efficiency. The easiest dataset (TriviaQA) generally has
 2272 fewest clusters because the answers are confident. In our results for the paragraph-
 2273 length task, we use 4-factoids per question as shown in (supplementary figure 4e), but
 2274 increasing the number of generations does not greatly increase the number of clusters
 2275 (supplementary figure 4f). All sentence-length-generation plots are for LLaMA 2 Chat
 2276 70B while the paragraph-generation plots are GPT-4.
 2277

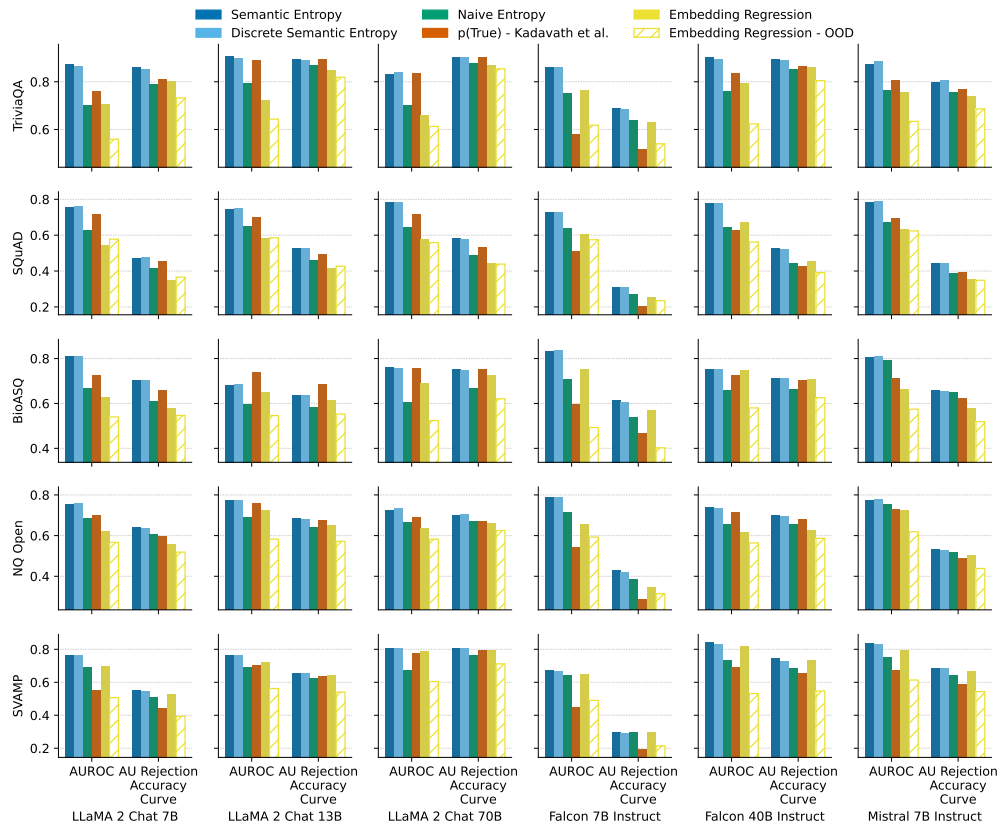
2278
 2279 by 51.4% for TriviaQA, 12.3% for BioASQ, and 18.0% for SQuAD (with the size of
 2280 the improvement caused by the proportion of identical answers produced by the model
 2281 for those datasets).
 2282

2283
 2284
 2285 **Note 5: Further Details for Sentence-Length Generations.**
 2286

2287 Here, we provide an unaggregated view of the sentence-length AUROCs which form
 2288 Figure 2. Individual datasets and models follow a very similar pattern to the average,
 2289 as shown in supplementary figure 5
 2290

2291 We also provide a more detailed view of the rejection accuracies at proportions of
 2292 answers retained for sentence-length generations in supplementary figure 6.
 2293

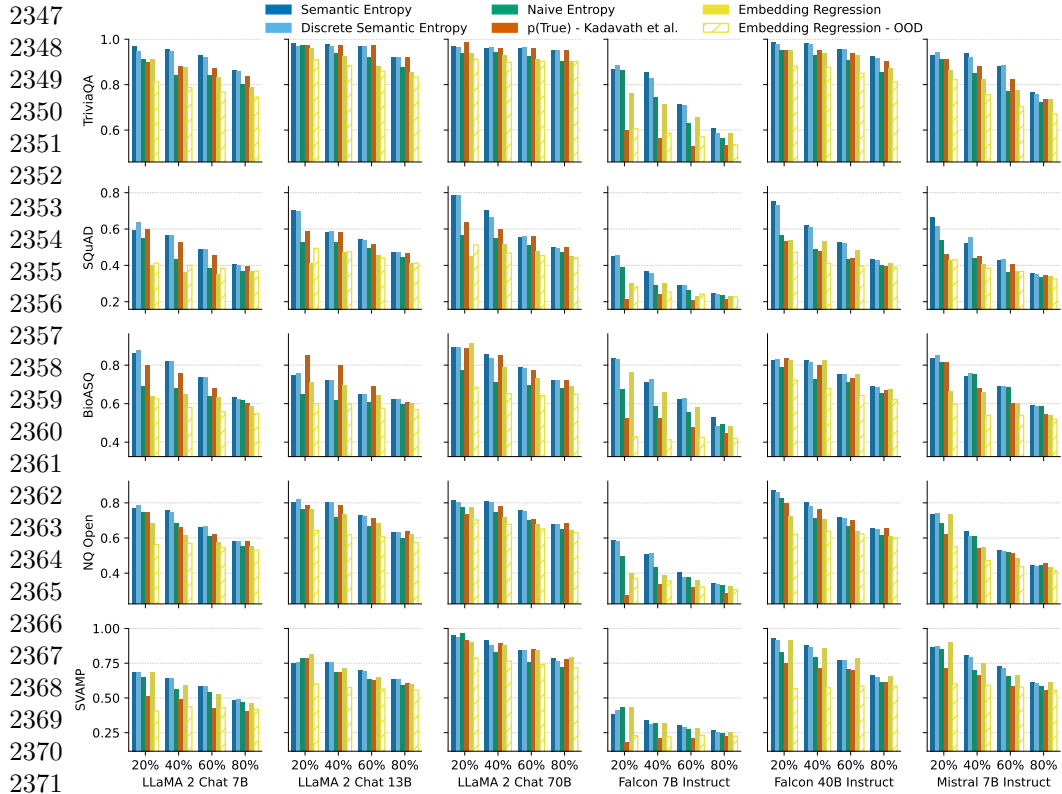
2294
 2295
 2296
 2297
 2298
 2299
 2300



Supplementary Figure 5: Sentence-length confabulation detection—full AUROC. An unaggregated view of the AUROCs shown in Figure 2.

Note 6: Assessing Model Accuracy

We check the quality of our automated ground-truth evaluation (using GPT-4 to compare the model generation with the reference answer) against human judgement by hand on sentence-length answers produced by LLaMA 2 Chat 70B responding to 100 questions from TriviaQA, SQuAD, and BioASQ. In each case, we check whether a generated answer matches the reference answer. Even if doing this reveals that the reference answer is wrong, which sometimes happens, we are interested in knowing whether humans and the automatic methods agree on the match, not on whether they know the actual correct answer. Supplementary table 4 shows the two human raters



2372 **Supplementary Figure 6: Sentence-length confabulation detection—**
 2373 **rejection accuracy.** A more detailed view of the rejection accuracies of the figures
 2374 provides further elaboration of the findings in Figure 2.
 2375

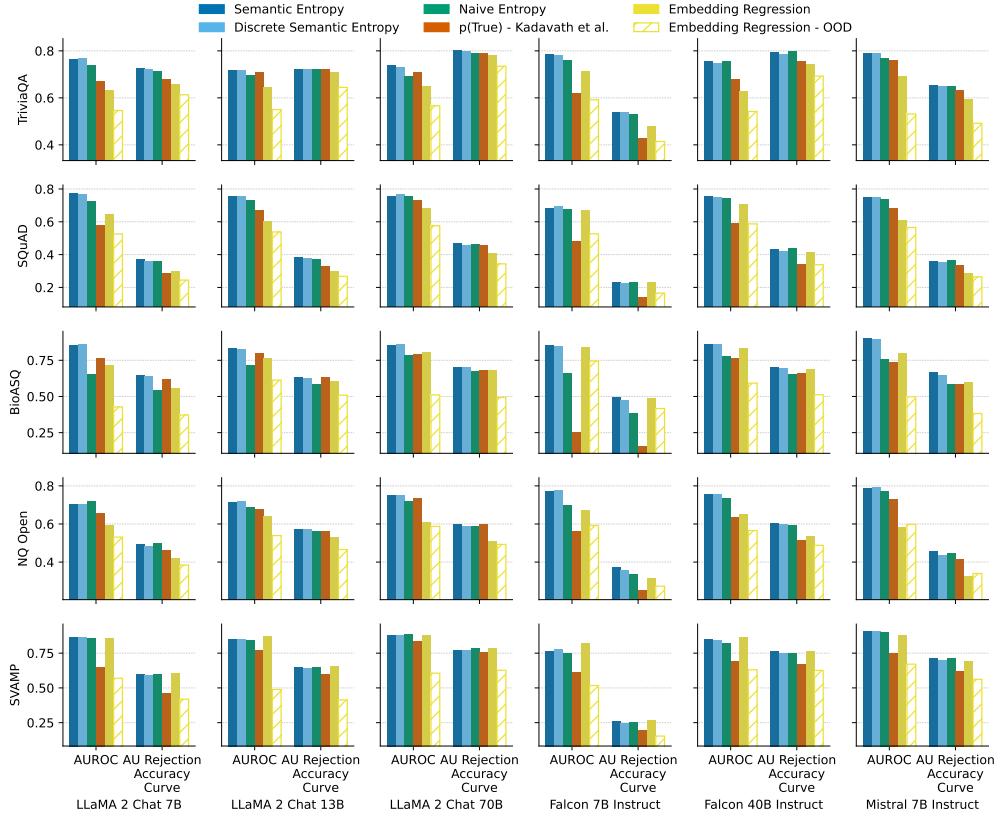
2376
 2377 agreed with each other at roughly the same rate (92%) as they agreed with GPT-4 on
 2378 average (93%). While GPT-3.5 is only slightly worse, in order to get the best ground-
 2379 truth estimation feasible, we use GPT-4 to compare the generated answer with the
 2380 reference answers provided in the dataset for results in this paper.
 2381
 2382
 2383

2384
 2385 **Note 7: Results for Short-Phrase Generation**
 2386

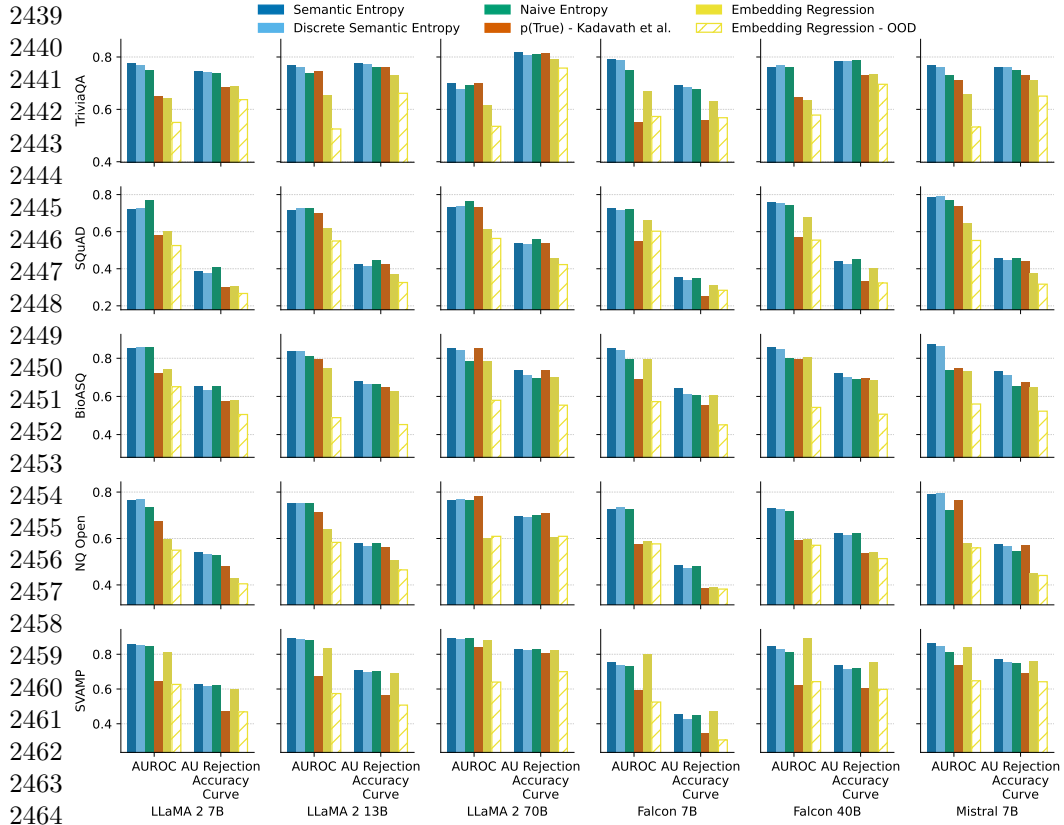
2387 In the main text, we provide results for sentence-length and paragraph-length gen-
 2388 eration. Here, we show results for a “short-phrase” scenario where we prompt the
 2389 model to “Answer the following question as briefly as possible”. Additionally, we here
 2390
 2391
 2392

	F1	LLaMA 2 Chat 70B	GPT-3.5	GPT-4	Human A	Human B
Human A	0.48	0.92	0.88	0.92	-	0.92
Human B	0.47	0.93	0.90	0.93	0.92	-
Human Average	0.47	0.92	0.89	0.93	-	-

Supplementary Table 4: Ground truth evaluation. We evaluate different automatic accuracy measures against human evaluation on sentence-length answers produced by LLaMA 2 Chat 70B on 100 randomly chosen questions from TriviaQA, SQuAD, and BioASQ. We find that GPT-4 agrees with both human raters at roughly the same level as they agree with each other.



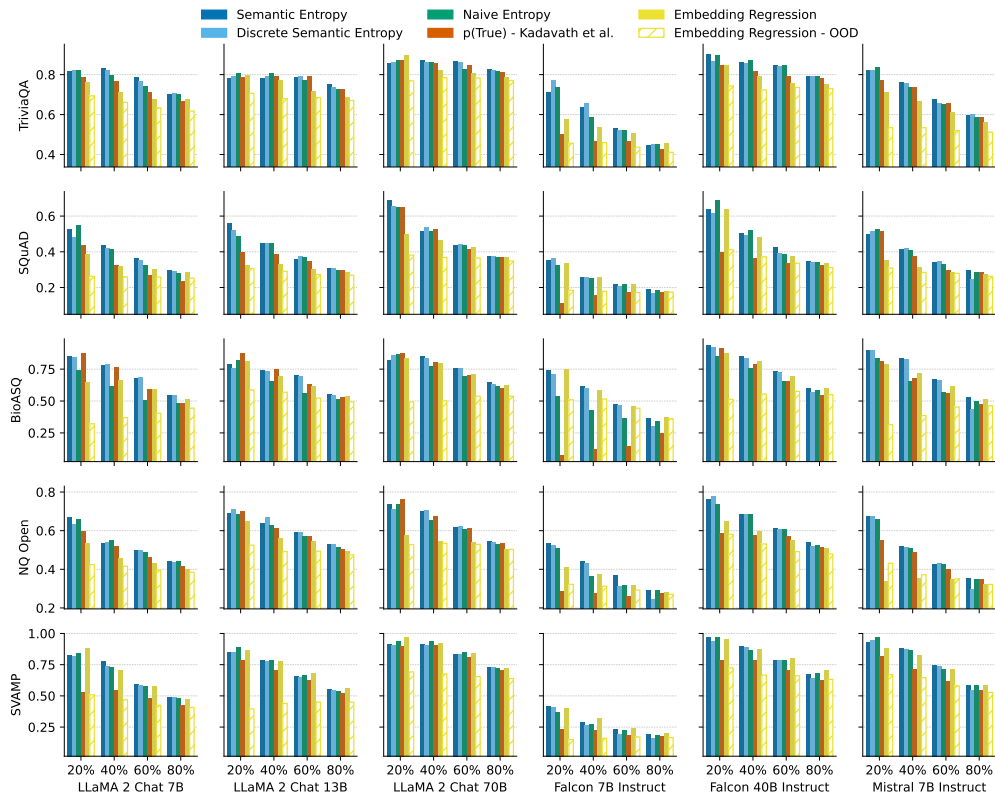
Supplementary Figure 7: Short-Phrase confabulation detection. With a prompt that encourages short generations, we achieve a lower average answer length (15.9 ± 21.8 characters compared to 95.6 ± 69.6 for sentence-length generations). Semantic entropy still works well compared to $p(\text{True})$ and embedding regression, although its advantage over naive entropy is smaller to due to the lower syntactic variability of the shorter generations.



2465 **Supplementary Figure 8: Short-Phrase confabulation detection—Non-**
 2466 **instruction-tuned models.** Models that have not been instruction-tuned have
 2467 slightly different output-distribution characteristics to instruction-tuned models. How-
 2468 ever, we find that both classes of models show broadly similar results and that semantic
 2469 entropy continues to outperform baselines in this setting.
 2470

2471
 2472 employ five few-shot QA demonstrations before the main question, all following the
 2473 template introduced above, which further encourages the LLM to predict with brevity
 2474 as the reference answers are usually very short for our selection of datasets. The short-
 2475 phrase scenario is less practically relevant, as users commonly interact with longer
 2476 LLM generations, although some settings value brevity and directness.
 2477
 2478

2480 For short-phrase generations we use the DeBERTa entailment classifier method
 2481 to check for strict bi-directional entailment, because this is cheaper and well-suited
 2482 to short phrases. In addition to the instruction-tuned LLaMA, Falcon, and Mistral
 2483
 2484

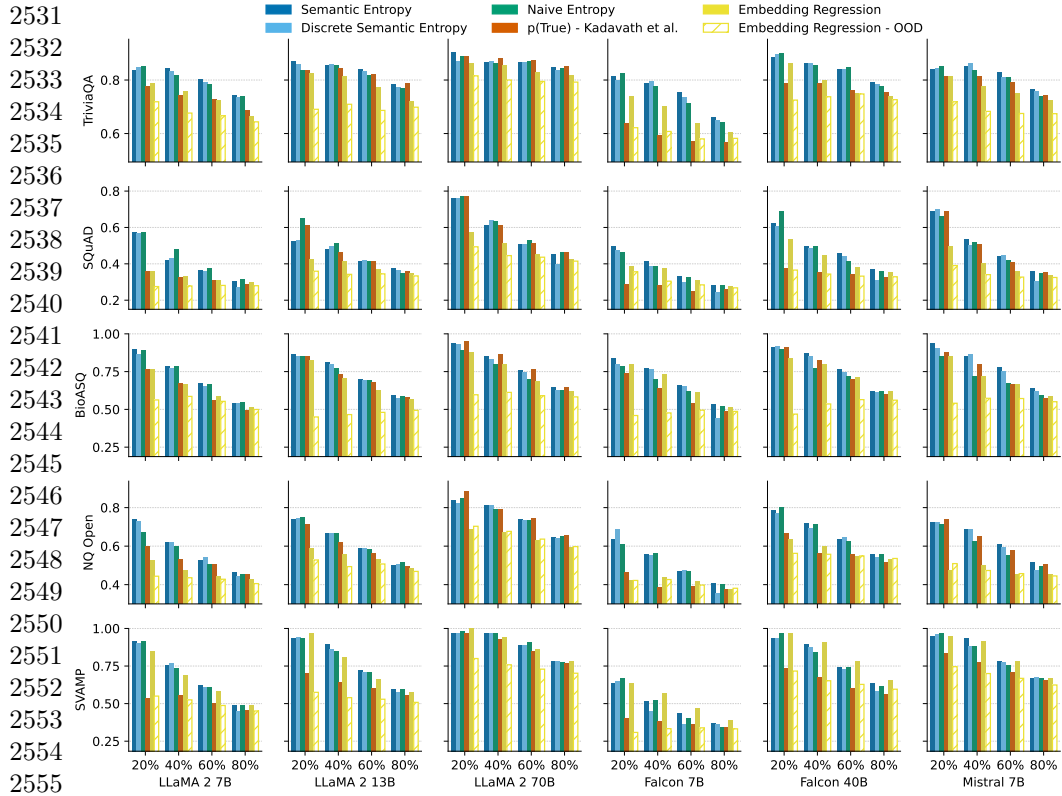


Supplementary Figure 9: Short-Phrase confabulation detection—rejection accuracy. A more detailed view of the rejection accuracies of the figures provides further elaboration of the findings in supplementary figure 7.

models we use in the sentence-length experiments, we additionally report results on the non-instruction-tuned LLaMA, Falcon, and Mistral models (which are not effective when applied to the sentence-length generation setting).

In supplementary figure 7, we show that results in this setting are broadly similar to those in the longer setting, with semantic entropy improving over $p(\text{True})$, embedding regression baseline, and naive entropy. In supplementary figure 8 we show this also holds for non-instruction-tuned models. Although non-instruction-tuned models have

2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530



2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556 **Supplementary Figure 10: Short-Phrase confabulation detection—rejection**
2557 **accuracy for non-instruction tuned models.** A more detailed view of the rejection
2558 accuracies of the figures provides further elaboration of the findings in supplementary
2559 figure 7.

2560
2561
2562 slightly different distributional characteristics—tending to have better-calibrated out-
2563 put probabilities—we find that semantic entropy continues to outperform the baselines
2564 in this setting.

2567 Averaged across the 60 combinations of tasks and models we study for the short-
2568 phrase setting, semantic entropy and discrete semantic entropy achieve the best mean
2569 AUROC (see below) values of 0.792 and 0.790 while naive entropy 0.760, $p(\text{True})$ 0.683,
2570 and the embedding regression baseline 0.708 lag behind it. Notably, semantic entropy
2571 continues to improve over $p(\text{True})$ and the embedding regression baseline by about 0.10
2572 AUROC. However, naive entropy performs better than for longer generations, although
2573
2574
2575
2576

semantic entropy does still improve over it by 0.03 on average (and by significantly more for individual datasets and models). This is because longer answers exhibit more of the syntactic variation that causes naive entropy to fail, and requiring answers to be as short as possible reduces the opportunity for variation.

In supplementary figure 9 we provide a more detailed view of the rejection accuracies at different proportions of answers retained. Lastly, supplementary figure 10 shows rejection accuracies for non-instruction tuned models.

References

- [76] Herbert Paul Grice. *Studies in the Way of Words*. Harvard University Press, Cambridge, 1989.
- [77] Isabelle Lorge and Janet B. Pierrehumbert. Not wacky vs. definitely wacky: A study of scalar adverbs in pretrained language models. *BlackBoxNLP Workshop at EMNLP*, 2023.
- [78] Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. Sources of hallucination by large language models on inference tasks. In *Findings of EMNLP*, 2023.

2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622