

Contributive Attribution for Question Answering via Tree-based Context Pruning

Anonymous ACL submission

Abstract

The development of large language models for question answering has benefited from understanding which context sentences are responsible for their answer. These sentences are commonly called contributive attribution. Recent works use the probability drop of the answer for a modified context to estimate how well sentences in the context match the attribution. Unfortunately, this metric does not convey the necessity and sufficiency qualities that the natural language processing community has defined in previous works. We propose a metric composed of a necessary and a sufficiency score based on probability drops to fill this gap. Then, to illustrate the soundness of the metric in practice, we develop a hierarchical method, TreeFinder, which progressively selects finer parts of the context through tree-based pruning using the metric. It begins with a few coarse-grained chunks and iteratively narrows the top k chunks according to our metric down to sentence-level granularity. At each iteration, we calculate our metric using ablation-based log-probability differences and filter out irrelevant chunks. Experimental results on HotpotQA demonstrate that TreeFinder outperforms ContextCite and TracLLM in contributive attribution quality when it is composed of a few sentences. Further experiments on Loogle and LongBench-v2 show that TreeFinder ranks sentences for attribution score better than ContextCite in long contexts.

1 Introduction

Revealing which sentences are determinant to the answer of Large Language Models (LLMs) in Question-Answering (QA) systems is an important step to building trusted agents. In an

ideal QA system, users should be able to trace the answers effortlessly to eliminate doubts, and methods should be built to automatically spot uncertainty.

Currently, LLMs are still challenged by language comprehension benchmarks, especially when contexts are long (Lee et al., 2024), due to the sparsity of rewards (Su et al., 2025) and limited data (Villalobos et al., 2022). This may be, in part, because current models are prone to missing key information in their context (Yang et al., 2025), especially as the number of tokens increases (Hsieh et al., 2024). Another part of the problem is that LLMs are prone to distractions, such as irrelevant information (Yoon et al., 2024) or adversarial attacks (Zou et al., 2023).

In this work, we focus on contributive attributions (Worledge et al., 2024a). Those are the key sentences that the model identifies as important to create its answer. Identifying these attributions enables, for example, the creation of robust fact-checking solutions by verifying that the answer is entailed by the attributions. They can also be used to generate a second answer that should match the first or enhance it (Huang et al., 2024). As a last example, they can show the biases of the model when presented with contradictory statements. Indeed, not having all variations equally weighed could indicate a preference.

These attributions should not be confused with corroborative attributions, which are the sentences supporting the correct answer (Worledge et al., 2024a,b). By contrast, contributive attributions (or rationales, citations, context traceback (Wang et al., 2025b)) fundamentally attempt to explain the answer provided by the LLM. They are thus very valuable, since

001

002

003

004

005

006

007

008

009

010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081 their mismatch with corroborative ones may
082 help to understand the failures of LLM when
083 their answers are not satisfactory.

084 Methods for finding contributive attributions
085 often leverage the innards of LLM, such as us-
086 ing attention scores or gradients similar to Grad-
087 CAM (Selvaraju et al., 2019). However, these
088 depend on many choices, such as the particular
089 head or layer computing these scores, that vary
090 with architecture (Pirene et al., 2025). Be-
091 yond making users lost in choices, these meth-
092 ods could become unusable as new LLM archi-
093 tectures, possibly incompatible, are developed.

094 In this work, we first propose a metric that
095 defines the quality of a set of sentences as the
096 contributive attribution using probability drops.
097 This metric consists of a sufficiency score and
098 a necessity score. Afterwards, we detail an al-
099 gorithm that segments a context, measures the
100 scores of the segments, selects the best ones
101 with a Top-k filter, and repeats greedily until
102 the desired sentence-level coarseness. Our al-
103 gorithm avoids enforcing a linear model, such
104 as with LIME (Ribeiro et al., 2016) or SHAP
105 (Lundberg and Lee, 2017) and thus does not
106 make assumptions on the linearity of the met-
107 ric under context ablations.

108 The paper is structured as follows. In [Sec-](#)
109 [tion 1.3](#), we give our motivation to use our met-
110 ric and we formally introduce the problem state-
111 ment. Then, we go over related works. Af-
112 terwards, in [Section 2](#), we explain our method-
113 ology and write the pseudocode to implement
114 our algorithm. We continue in [Section 3](#) with
115 our experiments, followed by the results and in
116 [Section 4](#), we conclude. Finally, for complete-
117 ness, we provide four appendices. [Appendix A](#)
118 provides details of how we divide the context
119 in each iteration. In [Appendix B](#), we display
120 the plots for more datasets. In [Appendix C](#), we
121 identify the top five sentences given by each
122 compared algorithm for two specific samples.
123 And, in [Appendix D](#), we show how contributive
124 and corroborative contributions diverge in prac-
125 tical cases.

126 1.1 Contributive attributions

127 The literature often defines multiple qualities
128 that attributions should possess. Two qualities,
129 a compactness and a sufficiency condition, have

130 been described by [Lei et al. \(2016\)](#). Then, a
131 necessity condition has been used by [Yu et al.](#)
132 [\(2019\)](#) and [DeYoung et al. \(2020\)](#). These were
133 used in subsequent works such as ([Brinner and](#)
134 [Zarri , 2023](#)) and ([Zhang et al., 2023](#)) for find-
135 ing attributions in two label problems, but not
136 for full sentence answers.

137 In this section, we want to define contributive
138 attributions through probability drops to satisfy
139 the three qualities of necessity, sufficiency, and
140 compactness. Using this definition, we derive
141 a single metric to optimize for. By construc-
142 tion, this metric only finds contributive attribu-
143 tions, since the probabilities and the answer are
144 derived from the same model.

145 From this point on, every piece of text will
146 be considered a sequence $C = [s_1, \dots, s_n]$ of
147 atoms, which are sentences in our experiments.
148 We use $|C| = n$ to denote the length of such a
149 sequence. In addition, we use $R \subseteq C$ to denote
150 that R is a subsequence of C . Note that a sub-
151 sequence is not necessarily contiguous, but the
152 relative order of its elements is conserved.

153 We consider a closed QA setting, where a
154 question Q is associated with a context C made
155 up of a number of $|C|$ separable elements (e.g.,
156 sentences). Let M be a generative model from
157 which we can generate an answer A and de-
158 rive its probability knowing Q and C giving
159 $P_M(A|Q, C)$.

160 We want to find $R \subseteq C$ defined as

$$161 \begin{aligned} R &:= \arg \min_{r \subseteq C} |r| \\ &s.t. \begin{cases} \text{Equation (2)} \\ \text{Equation (3)} \end{cases} \end{aligned} \quad (1)$$

$$162 \begin{aligned} \log P_M(A | Q, r) \\ = \log P_M(A | Q, C) \end{aligned} \quad (2)$$

$$163 \begin{aligned} \log P_M(A | Q, C \setminus r) \\ = \log P_M(A | Q, \emptyset) \end{aligned} \quad (3)$$

164 In practice, the only solution of [Equation \(1\)](#)
165 with equality constraints is $r = C$ because
166 the probability of the answer changes even by
167 adding irrelevant sentences. We therefore relax
168 these requirements to add small tolerances. In
169 this context, we adopt the following formal def-
170 inition of R :

$$R := \arg \min_{r \subseteq C} |r| \quad (4)$$

$$s.t. \begin{cases} \text{Equation (5)} \\ \text{Equation (6)} \end{cases}$$

$$\log P_M(A | Q, C) - \log P_M(A | Q, r) \leq \varepsilon_{\text{suf}} \quad (5)$$

$$\log P_M(A | Q, C \setminus r) - \log P_M(A | Q, \emptyset) \leq \varepsilon_{\text{nec}} \quad (6)$$

where ε_{suf} and ε_{nec} are small constants that control tolerance for the sufficiency and necessity constraints, respectively.

We add the left-hand sides of the sufficiency and necessity inequalities, weighed respectively by $1 - \alpha$ and α , to produce a single metric. We get

$$a_M(Q, C, A, r) = \alpha(\log P_M(A | Q, C \setminus r) - \log P_M(A | Q, \emptyset)) + (1 - \alpha)(\log P_M(A | Q, C) - \log P_M(A | Q, r)). \quad (7)$$

1.2 Probability drop

Imposing modifications in a context C by removing some parts to create a subsequence $C' \subseteq C$ modifies the probability distribution of the answer. We can infer the loss or gain of information by analyzing the relationships between the distributions $P_M(\cdot | Q, C)$, $P_M(\cdot | Q, C')$, $P_M(\cdot | Q, C \setminus C')$, and $P_M(\cdot | Q, \emptyset)$. The first and fourth are constant with respect to the choice of C' . Together with the second and third, they produce the necessary and sufficient conditions to produce A when they are equal.

When selecting the candidate subsequence $C' \subseteq C$, we can work at different levels of granularity. In our case, the finest granularity would correspond to the sentence-level, where a sentence is defined by NLTK (Bird et al., 2009). But we can also consider coarser levels, where sentences are grouped by paragraphs, pages, etc. We propose a hierarchical approach that will start from a coarse subdivision of the context, and then, after selecting the best C' by removing complete groups of sentences, it will iteratively consider finer and finer granularity levels.

In this paper, we rank the chunks with the metric of Equation (7) and then choose to keep the t_k best ones to create C' .

1.3 Related Work

In the QA setting, multiple fundamental questions arise: Does the answer come from internal or external knowledge? Is it reliable? Can it be explained or traced back? When we limit ourselves to the context, the main interest in explainability is corroborative (elements supporting an answer) and contributive (elements responsible for the answer) attributions (Worledge et al., 2024a).

Corroborative Attribution Corroborative attributions can exist without a reference model because they correspond to the rationale for the construction of the answer. Although interesting, they fundamentally achieve a different goal from ours. While we want to discover what the LLM is using as support, they want to support the ideal answer. Both attributions can be found, and their overlap is an indicator of the correctness of the model. We provide an intuition for this in Appendix D.

Contributive Attribution Contributive attributions can be given by the model itself. For example, GopherCite (Menick et al., 2022) uses guided generation to provide exact quotes. SelfCite (Chuang et al., 2025) refines this through Reinforcement Learning using probability drops as reward signals. Both methods rely heavily on the model being trusted, which is what we want to avoid. We therefore focus on properties that are impartial, such as probability distributions or attention weights, without dependence on the statistical generation.

Reducing the context or prompt to eliminate redundant information can be done by paraphrasing, filtering (Hou et al., 2024) or encoding it (Cheng et al., 2024) so that its size is reduced but the content is largely retained in meaning (Li et al., 2025). (Feng et al., 2018) uses an iterative filtering in which they keep removing the word that has the least impact. In contrast to us, they operate at a fixed grain and do not incorporate a tree approach.

The search for contributive attributions through the alteration of the context can be

207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254

found in the literature as a perturbation-based method for feature attribution (Zhao et al., 2024). Perturbations can be random, as in LIME or generated by models such as variational autoencoders (Alvarez-Melis and Jaakkola, 2017).

One recent work leveraging perturbations to find attributions is ContextCite (Cohen-Wang et al., 2025). To avoid computing the probability of a particular answer after the removal of an atom of text, ContextCite generates random ablation vectors and extrapolates the scores of each sentence with a linear model. In contrast, we avoid using such a model by finding important chunks of successively smaller size directly from the answer probabilities. ContextCite does not compute probabilities directly; instead, they compute the logit function $\log \frac{p}{1-p}$ of the difference of the chosen logits with a smooth approximation of the maximum of the rest of the logits (*LogSumExp*). This slightly modifies the objective.

TracLLM (Wang et al., 2025b) uses a tree-based method. They create chunks and remove the ones that have the lowest metric permanently from the context before continuing deeper. To compute their metric, they sample and aggregate a portion of all possible marginal probability drops. Similarly to ContextCite, they use linear models to compute a score for each chunk at every depth.

As emphasized in (Su, 2025), these approaches leveraging probability drop, working as a consequence of the answer being modeled probabilistically, are fundamentally robust to architectural modifications that are currently the focus of research for stronger LLMs (Assran et al., 2025; Wang et al., 2025a).

2 TreeFinder

The complete algorithm is described in Algorithm 1. Here is a short summary of how it works. It first divides the context into chunks. Afterwards, it computes both scores for each individual chunk in the running as R , and filters them to keep only the most promising ones. It then repeats with smaller chunks until they all reach the size of one sentence. This process is summarized in Figure 1.

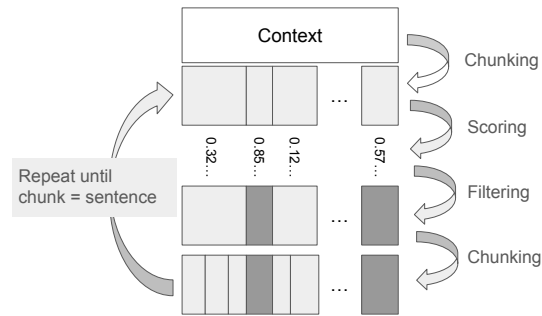


Figure 1: TreeFinder overview. The three repeating steps are represented in light gray. Dark-gray boxes represents the filtered out sentences that will not be subdivided or rescored in further repetitions.

We have defined two intermediary functions in Algorithms 2 and 3: Score and Chunkify. The first computes our metric $a(Q, C, A, r)$ for each given chunk r . The second simply cuts the context into approximately k chunks, returning a list of chunks and a mapping from the indices of sentences in the global context to the indices of the chunks that contain them. Since it is a straightforward implementation, it has been moved to the Appendix. A Top-k filter is present to ensure a logarithmic time complexity with the size of the context $|C|$. The filter can be modified to be more restrictive for faster code execution.

3 Experiments

In this section, we first choose three datasets to use in our experiments. We then describe the LLM models and parameters used for the experiments. Next, we compare our algorithm with ContextCite and TracLLM on the chosen datasets and explore the influence of the necessity weight parameter α of Algorithm 1.

3.1 Datasets

We work with the following three datasets. The first is HotpotQA (Yang et al., 2018), a widely used dataset for corroborative attribution in QA. Its purpose is to benchmark LLM for their ability to make multi-hop reasoning. Its “distractor” mode includes various paragraphs of related articles, which makes it quite challenging for LLMs.

The other two, LooGLE (Li et al., 2024) and LongBench-v2 (Bai et al., 2025), are made up of long contexts. We use two subdivi-

Algorithm 1 TreeFinder

```
1: Input:  $M$  (Model),  $C$  (Context),  $Q$  (Question),  $A$  (Answer),  $k$  (Initial Chunking Factor),  $t_k$  (Top-k),  $\alpha$  (Necessity weight)
2: Output:  $X$  (Per sentence score)
3:  $V \leftarrow \mathbf{1}^{|C|}$  ▷ Ablation
4:  $X \leftarrow \mathbf{0}^{|C|}$  ▷ Score
5:  $V_c \leftarrow \mathbf{1}^{|k|}$  ▷ Chunk ablation
6:  $X_c \leftarrow \mathbf{0}^{|k|}$  ▷ Chunk score
7: while  $k < t_k|C|$  and  $V_c \neq \mathbf{0}$  do
8:    $S, map \leftarrow \text{Chunkify}(C, k)$ 
9:   for  $0 \leq i < |S|$  do
10:     for  $j \in map[i]$  do
11:        $V[j] \leftarrow V[j] \wedge V_c[i]$ 
12:        $X[i] \leftarrow X_c[i]$  if  $V_c[i]$ 
13:     end for
14:   end for
15:    $X_c \leftarrow \text{Score}(C, Q, A, M, V, S, map, \alpha)$ 
16:    $V_c \leftarrow \text{Top-k}(X_c, map, t_k)$ 
17:    $k \leftarrow k * t_k$ 
18: end while
19: return  $X$ 
```

Algorithm 2 Score

```
1: Input:  $C$  (Context),  $Q$  (Question),  $A$  (Answer),  $M$  (Model),  $V$  (Ablation vector),  $S$  (Chunks),  $map$  (Local to global indices),  $\alpha$  (Necessity weight)
2: Output:  $scores$  (Per chunk score)
3:  $P_{total} \leftarrow \log P_M(A|Q, C)$ 
4:  $P_{\emptyset} \leftarrow \log P_M(A|Q, \emptyset)$ 
5: for  $0 \leq i < |S|$  do
6:   if  $\forall j \in map[i], V[j] = 1$  then
7:      $P_{rem}[i] \leftarrow \log P_M(A|Q, C \setminus S_i)$ 
8:      $P_{unique}[i] \leftarrow \log P_M(A|Q, S_i)$ 
9:   end if
10: end for
     $scores = \alpha(P_{rem} - P_{\emptyset}) + (1 - \alpha)(P_{total} - P_{unique})$ 
11: return  $scores$ 
```

sions (named 2a and 2b hereafter) of this LooGLE: short dependencies and long dependencies. These measure the distance between the different elements needed to answer in the context.

All datasets have been truncated to allow the algorithms to run on a maximum of two A100 40GB (one for HotpotQA). LooGLE and Longbench-v2 have had their maximum context length set to 20000 tokens, with a maximum of a thousand sample per for time and budget constraints.

3.2 Model and Parameters

The algorithm has 3 hyperparameters: the chunking factor k , the Top-k value $t_k = 3$ and the necessity weight α . We chose $k = 6$ with $t_k = 3$ and $\alpha = 0.25$ after having tried $k = 2$ and $k = 3$ which gave poor results.

The model M chosen for all experiments is Qwen-2.5-7B-Instruct-1M (Yang et al., 2025) for its stated long context understanding. Comparisons are carried out using the Transformers Python library (Wolf et al., 2020).

Corroborative attribution tests are run using the vLLM engine (Kwon et al., 2023) for fast inference in long contexts.

3.3 Results

Evaluation criteria We compute the necessity and sufficiency of the first sentences taken together in a given group of size $w \in [1, \dots, 5]$ with $\alpha = 0.5$ and report the average and standard error in the Figure 2 for HotpotQA. Since this is a greedy approach, “Ground Truth” is an upper bound.

Lastly, in Figures 3 and 4, we search for the first five “Ground Truth” sentences and report the median positions each method has assigned to them. Unfortunately, since TracLLM does not provide a score for all sentences, we cannot include them.

It might seem surprising that the evaluation does not use human annotations, but this is intended since the goal is to identify the sentences M relies on.

Comparison We now compare the rankings given by our algorithm with those provided by ContextCite and TracLLM. We can see on Fig-

337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383

384	ure 2, that for the first sentence, TracLLM and	Loogle (long). However, contributive attribu-	434
385	ContextCite reach lower average scores than	tions capture more than the support, and thus a	435
386	TreeFinder. However, as more sentences are in-	cluded, TreeFinder improves the average score,	436
387	cluded, TreeFinder improves the average score,		
388	showing the improved quality of the broader	Influence of the scores necessity weight	437
389	contributive attribution.	parameter In this section, we modify the	438
390	With the selected hyperparameters, we	weights of the necessity and sufficiency scores	439
391	achieve a similar time of 5.72s as ContextCite	α from 0 to 1. Choosing 0 or 1 nearly halves	440
392	(5.84s) despite an increase in the number of for-	the number of forward calls to the LLM since	441
393	ward calls (46.9 against 32). The reason is that	we do not compute one of the two metrics, any	442
394	$P(A Q, r)$ becomes cheaper in terms of tokens	other repartition does not have this benefit.	443
395	used compared to $P(A Q, C \setminus r)$ as r shortens.	From Figure 4, we see that the best value for	444
396	This leads to about half of the calls that account	α is between 0 and 1 since our original value of	445
397	for the entire compute time. Since TracLLM	$\alpha = 0.25$ achieves a lower sentence ranking in	446
398	discards chunks from the context and never uses	all cases. However, in HotpotQA and Loogle	447
399	them again, they also have a disproportionate	(long), $\alpha = 0$ is much faster than $\alpha = 0.25$	448
400	amount of calls (293.3) to the compute time	(2x and 8x, respectively) and barely modifies	449
401	(25.6s) compared to ContextCite. Their default	the ranking, so we could sacrifice a bit of accu-	450
402	parameters have three score estimations that	racy for speed in certain contexts.	451
403	compute many more marginal probabilities for		
404	the same initial chunking factor (their $K = 6$).	3.4 Discussions	452
405	For the score-based criterion, no method	When using probability drop to find contribu-	453
406	shows a clear advantage across all datasets (oth-	tive attribution, we must consider an inherent is-	454
407	ers in Appendix). In Appendix B, we provide	sue coming from the method itself. That is, the	455
408	results for Loogle (short) and Longbench where	alteration of the context can modify its meaning.	456
409	TracLLM is better, as well as Loogle (long)	Indeed, many sentences refer to previous ones,	457
410	where the result depends on α . In general,	such that, if a sentence modifying the subject is	458
411	while TracLLM is slower, its attributions are	removed, then the other sentences referring to	459
412	closer in value to the ground truth.	it might imply wrong facts.	460
413	Since the metric values are close, we provide	Additionally, the LLM can miss important	461
414	a last angle for comparison in Figure 3, where	facts in the initial long context that can be used	462
415	we take the positions at which the first k ground-	when it becomes truncated. According to our	463
416	truth sentences can be found in each method	definition, these should not be included in the	464
417	and report the median over the dataset. We find	attribution, but our relaxation can include them.	465
418	that our method misses the ranking of the most	If r contains them and the answer is correct,	466
419	important sentence, but the following are con-	then the left-hand-side of the sufficiency con-	467
420	sistently classified better than ContextCite. In-	dition in Equation (4) can become negative.	468
421	creasing the number of ablations (calls) for Con-	The Top-k filter t_k plays a strategic role in	469
422	textCite from 32 to match our methods' num-	the rapid termination of the algorithm. The	470
423	ber of calls, we find no significant improvement.	value and softmax thresholds can sometimes	471
424	This, we believe, highlights the need to incor-	prune the tree but cannot always do so at ev-	472
425	porate both the necessity and sufficiency scores	ery step, and thus do not bind the time complex-	473
426	in attribution finding algorithms to properly en-	ity of TreeFinder. Every node has a maximum	474
427	compass all the desired properties of an attribu-	of t_k children, so the maximum is reached at	475
428	tion.	$t_k * k$. The total number of nodes explored is	476
429	We note that our slice of HotpotQA has a me-	this result multiplied by the depth of the tree,	477
430	median of 2 sentences as support, while Loogle	which is the logarithm of the number of sen-	478
431	(long) has 3 and Loogle (short) has 1 (and Long-	tences, and the two calls per iteration with the	479
432	bench is not labeled). The ranking difference,	two initial calls for the empty and complete con-	480
433	therefore, can be argued to be only relevant on	texts. This leads to the number of forward calls	481

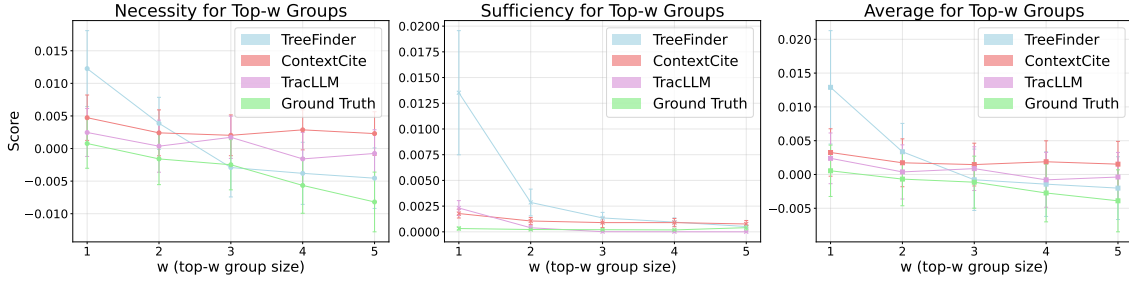


Figure 2: Sufficiency and necessity for the top sentences jointly in HotpotQA. Lower is better. TreeFinder lags behind for the first three sentences but quickly achieves a better overall attribution quality when more sentences are taken into account.

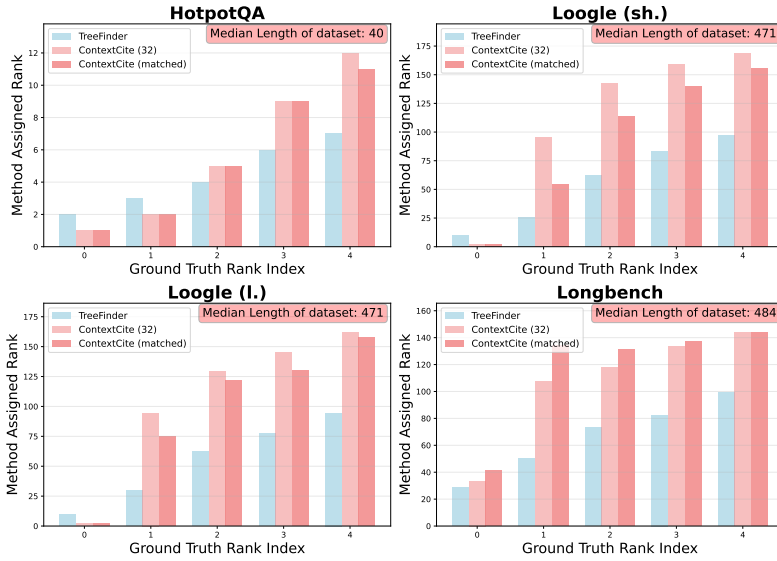


Figure 3: Median positions the k^{th} ground truth sentences ($\alpha=0.25$). Median context length in sentences shown per dataset. Lower median positions are better. The sentences beyond the first are almost always found earlier in the ranking of TreeFinder. Increasing the number of ablations in ContextCite from 32 to match TreeFinder (50 for HotpoQA, 100 for the others) does not close the gap.

482 $\#M_{forward} = \mathcal{O}(4 + 2(k * t_k) * \log_{t_k}(|C|)).$

483 There can be an order of magnitude between
 484 the values of necessity and sufficiency scores.
 485 Therefore, we believe that balancing α such
 486 that both scores contribute meaningfully to the
 487 choice of attribution is difficult. A possible remedy
 488 to this issue could be the normalization of
 489 the scores at each depth in the tree.

490 **4 Conclusion**

491 This work presented a metric to measure the
 492 quality of contributive attributions with prob-
 493 ability drops. Using this metric, we have
 494 developed a method, namely TreeFinder, to
 495 iteratively extract contributive attribution in
 496 question-answer environments using out-of-
 497 the-box LLMs without requiring task-specific

training or surrogate models. TreeFinder em-
 498 ploys a tree-based context pruning strategy, pro-
 499 gressively refining the contributive attribution
 500 from coarse-grained chunks to individual sen-
 501 tences. All code is available at: [https://anony-
 502 mous.4open.science/r/TreeFinder-F3FD](https://anonymous.4open.science/r/TreeFinder-F3FD)
 503

Experimental results on HotpotQA demon-
 504 strate that TreeFinder ranks attributions, as
 505 commonly defined in previous works, better
 506 than ContextCite when they are composed of
 507 multiple sentences. This highlights the benefits
 508 of enforcing the necessity and sufficiency crite-
 509 ria in the search process.
 510

Future work can focus on improving the
 511 chunking to group complementary sentences or
 512 incorporate meaningful structures of the con-
 513 text. Indeed, syntactically meaningful tree rep-
 514

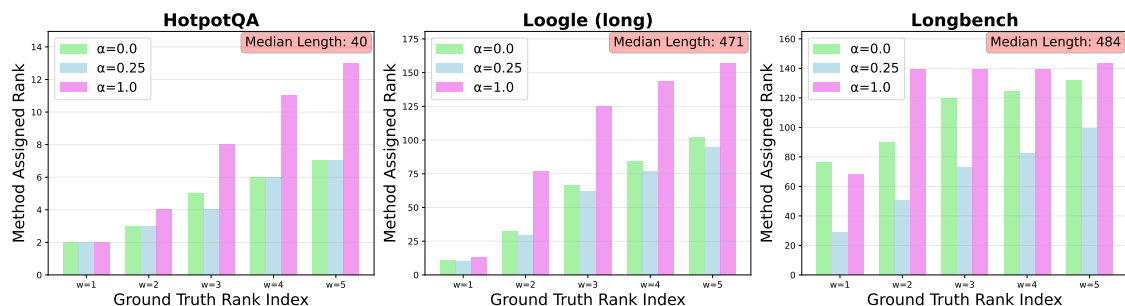


Figure 4: Ablation study of Tree Finder across $\alpha = 0.0, 0.25, 1.0$. Median context length in sentences shown per method. Lower median positions are better. $\alpha = 0$ can be enough to find attributions but Longbench necessitates both necessity and sufficiency.

515 representations have been shown to improve the
 516 results of textual classification tasks (Nallapati
 517 and Allan, 2002). We believe that more accu-
 518 racy and speed can be achieved by using textual
 519 segmentation techniques (Ghinassi et al., 2024).
 520 Alternative pruning decisions and node score
 521 estimations can potentially guide the process to
 522 be faster and more accurate.

523 Additionally, dropping some sentences X
 524 from the context for the remainder of the algo-
 525 rithm, as does TracLLM, can also improve
 526 termination speed. However, this changes
 527 the marginal probability from $P(A|Q, C)$ to
 528 $P(A|Q, C \setminus X)$ and could have unforeseen con-
 529 sequences on the attribution found.

530 As a final word, we believe that TreeFinder
 531 represents a step towards building more trust-
 532 worthy and explainable QA systems, enabling
 533 users to not only receive reliable answers but
 534 also understand why those answers are gener-
 535 ated, fostering confidence, and facilitating fact-
 536 checking.

537 Limitations

538 Due to the time and compute required to pro-
 539 duce certain graphs, our work only explores re-
 540 sults for a single model. This can impact the
 541 generalization of our findings.

542 Since TracLLM and ContextCite already
 543 compare themselves to traditional methods
 544 such as attention, shapley and others, we only
 545 compared ourselves to them to save time, com-
 546 pute and paper length. There may be cases
 547 where these traditional methods would have
 548 added insights to our analysis.

Acknowledgments

Omitted for submission.

References

- David Alvarez-Melis and Tommi S. Jaakkola. 2017. A causal framework for explaining the predictions of black-box sequence-to-sequence models. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 412–421. Association for Computational Linguistics.
- Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Mojtaba, Komeili, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, Sergio Arnaud, Abha Gejji, Ada Martin, Francois Robert Hogan, Daniel Dugas, Piotr Bojanowski, Vasil Khalidov, and 11 others. 2025. V-JEPA 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv: 2506.09985*.
- Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2025. LongBench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 3639–3664. Association for Computational Linguistics.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. O’Reilly Media, Inc.”.
- Marc Brinner and Sina Zarrieß. 2023. Model interpretability and rationale extraction by input

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586	mask optimization. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13722–13744, Toronto, Canada. Association for Computational Linguistics.	language models? <i>arXiv preprint arXiv:2404.06654</i> .	641
587			642
588			
589			
590	Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan Zhao. 2024. xrag: Extreme context compression for retrieval-augmented generation with one token . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Lei Huang, Xiaocheng Feng, Weitao Ma, Yuxuan Gu, Weihong Zhong, Xiachong Feng, Weijiang Yu, Weihua Peng, Duyu Tang, Dandan Tu, and Bing Qin. 2024. Learning fine-grained grounded citations for attributed large language models . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 14095–14113, Bangkok, Thailand. Association for Computational Linguistics.	643 644 645 646 647 648 649 650 651
591			
592			
593			
594			
595			
596			
597			
598			
599	Yung-Sung Chuang, Benjamin Cohen-Wang, Shannon Zejiang Shen, Zhaofeng Wu, Hu Xu, Xi Victoria Lin, James Glass, Shang-Wen Li, and Wentau Yih. 2025. SelfCite: Self-supervised alignment for context attribution in large language models . <i>arXiv preprint arXiv:2502.09604</i> .	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with PagedAttention . In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> .	652 653 654 655 656 657 658
600			
601			
602			
603			
604			
605	Benjamin Cohen-Wang, Harshay Shah, Kristian Georgiev, and Aleksander Madry. 2025. ContextCite: Attributing model generation to context . <i>Advances in Neural Information Processing Systems</i> , 37:95764–95807.	Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua, Devendra Singh Sachan, Michael Boratko, Yi Luan, Sébastien M. R. Arnold, Vincent Perot, Siddharth Dalmia, Hexiang Hu, Xudong Lin, Panupong Pasupat, Aida Amini, Jeremy R. Cole, Sebastian Riedel, Iftekhar Naim, Ming-Wei Chang, and Kelvin Guu. 2024. Can long-context language models subsume retrieval, RAG, SQL, and more? <i>arXiv preprint arXiv:2406.13121</i> .	659 660 661 662 663 664 665 666 667
606			
607			
608			
609			
610	Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C. Wallace. 2020. ERASER: A benchmark to evaluate rationalized NLP models . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4443–4458, Online. Association for Computational Linguistics.	Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2016. Rationalizing neural predictions . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 107–117, Austin, Texas. Association for Computational Linguistics.	668 669 670 671 672 673
611			
612			
613			
614			
615			
616			
617			
618	Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, and Jordan Boyd-Graber. 2018. Pathologies of neural models make interpretations difficult . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3719–3728, Brussels, Belgium. Association for Computational Linguistics.	Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2024. LooGLE: Can long-context language models understand long contexts? In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 16304–16333. Association for Computational Linguistics.	674 675 676 677 678 679 680 681
619			
620			
621			
622			
623			
624			
625			
626	Iacopo Ghinassi, Lin Wang, Chris Newell, and Matthew Purver. 2024. Recent trends in linear text segmentation: A survey . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 3084–3095, Miami, Florida, USA. Association for Computational Linguistics.	Zongqian Li, Yinhong Liu, Yixuan Su, and Nigel Collier. 2025. Prompt compression for large language models: A survey . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2025 - Volume 1: Long Papers, Albuquerque, New Mexico, USA, April 29 - May 4, 2025</i> , pages 7182–7195. Association for Computational Linguistics.	682 683 684 685 686 687 688 689 690 691
627			
628			
629			
630			
631			
632	Haowen Hou, Fei Ma, Binwen Bai, Xinxin Zhu, and Fei Yu. 2024. Enhancing and accelerating large language models via instruction-aware contextual compression . <i>arXiv preprint arXiv:2408.15491</i> .	Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	692 693 694 695
633			
634			
635			
636			
637			
638			
639			
640			

696	Jacob Menick, Maja Trebacz, Vladimir Mikulik,	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	751
697	John Aslanides, Francis Song, Martin Chadwick,	Chaumond, Clement Delangue, Anthony Moi,	752
698	Mia Glaese, Susannah Young, Lucy Campbell-	Pierric Cistac, Tim Rault, Rémi Louf, Morgan	753
699	Gillingham, Geoffrey Irving, and Nat McAleese.	Funtowicz, Joe Davison, Sam Shleifer, Patrick	754
700	2022. Teaching language models to support an-	von Platen, Clara Ma, Yacine Jernite, Julien Plu,	755
701	swers with verified quotes. <i>arXiv preprint arXiv:</i>	Canwen Xu, Teven Le Scao, Sylvain Gugger, and	756
702	<i>2203.11147</i> .	3 others. 2020. Transformers: State-of-the-art	757
703	Ramesh Nallapati and James Allan. 2002. Captur-	natural language processing . In <i>Proceedings of</i>	758
704	ing term dependencies using a language model	<i>the 2020 Conference on Empirical Methods in</i>	759
705	based on sentence trees. In <i>Proceedings of the</i>	<i>Natural Language Processing: System Demon-</i>	760
706	<i>eleventh international conference on Information</i>	<i>strations</i> , pages 38–45, Online. Association for	761
707	<i>and knowledge management</i> , pages 383–390.	Computational Linguistics.	762
708	Lize Pirene, Samy Mokeddem, Damien Ernst, and	Theodora Worledge, Judy Hanwen Shen, Nicole	763
709	Gilles Louppe. 2025. Exploration of rationale-	Meister, Caleb Winston, and Carlos Guestrin.	764
710	extraction methods for closed-domain question	2024a. Unifying corroborative and contributive	765
711	answering with a new sentence-level rationale	attributions in large language models . In <i>2024</i>	766
712	dataset. In <i>Natural Language Processing and In-</i>	<i>IEEE Conference on Secure and Trustworthy Ma-</i>	767
713	<i>formation Systems</i> , pages 3–13, Cham. Springer	<i>chine Learning (SaTML)</i> , pages 665–683, Los	768
714	Nature Switzerland.	Alamitos, CA, USA. IEEE Computer Society.	769
715	Marco Tulio Ribeiro, Sameer Singh, and Carlos	Theodora Worledge, Judy Hanwen Shen, Nicole	770
716	Guestrin. 2016. "Why should I trust you?": Ex-	Meister, Caleb Winston, and Carlos Guestrin.	771
717	plaining the predictions of any classifier . In	2024b. Unifying corroborative and contributive	772
718	<i>Proceedings of the 22nd ACM SIGKDD Interna-</i>	attributions in large language models . In <i>2024</i>	773
719	<i>tional Conference on Knowledge Discovery and</i>	<i>IEEE Conference on Secure and Trustworthy Ma-</i>	774
720	<i>Data Mining</i> , KDD '16, page 1135–1144, New	<i>chine Learning (SaTML)</i> , pages 665–683.	775
721	York, NY, USA. Association for Computing Ma-	An Yang, Bowen Yu, Chengyuan Li, Dayiheng Liu,	776
722	chinery.	Fei Huang, Haoyan Huang, Jiandong Jiang, Jian-	777
723	Ramprasaath R. Selvaraju, Michael Cogswell,	hong Tu, Jianwei Zhang, Jingren Zhou, Junyang	778
724	Abhishek Das, Ramakrishna Vedantam, Devi	Lin, Kai Dang, Kexin Yang, Le Yu, Mei Li, Min-	779
725	Parikh, and Dhruv Batra. 2019. Grad-CAM:	min Sun, Qin Zhu, Rui Men, Tao He, and 9 oth-	780
726	Visual explanations from deep networks via	ers. 2025. Qwen2.5-1M technical report. <i>arXiv</i>	781
727	gradient-based localization . <i>International Jour-</i>	<i>preprint arXiv: 2501.15383</i> .	782
728	<i>nal of Computer Vision</i> , 128:336–359.	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-	783
729	Weijie Su. 2025. Do large language models (re-	gio, William Cohen, Ruslan Salakhutdinov, and	784
730	ally) need statistical foundations? <i>arXiv preprint</i>	Christopher D. Manning. 2018. HotpotQA: A	785
731	<i>arXiv: 2505.19145</i> .	dataset for diverse, explainable multi-hop ques-	786
732	Yi Su, Dian Yu, Linfeng Song, Juntao Li, Haitao Mi,	tion answering . In <i>Proceedings of the 2018 Con-</i>	787
733	Zhaopeng Tu, Min Zhang, and Dong Yu. 2025.	<i>ference on Empirical Methods in Natural Lan-</i>	788
734	Crossing the reward bridge: Expanding RL with	<i>guage Processing</i> , pages 2369–2380, Brussels,	789
735	verifiable rewards across diverse domains. <i>arXiv</i>	Belgium. Association for Computational Linguis-	790
736	<i>preprint arXiv: 2503.23829</i> .	tics.	791
737	Pablo Villalobos, Anson Ho, Jaime Sevilla, Tamay	Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang,	792
738	Besiroglu, Lennart Heim, and Marius Hobbahn.	Minbyul Jeong, and Jaewoo Kang. 2024. Com-	793
739	2022. Will we run out of data? Limits of LLM	pact: Compressing retrieved documents actively	794
740	scaling based on human-generated data. <i>arXiv</i>	for question answering . In <i>Proceedings of the</i>	795
741	<i>preprint arXiv: 2211.04325</i> .	<i>2024 Conference on Empirical Methods in Nat-</i>	796
742	Guan Wang, Jin Li, Yuhao Sun, Xing Chen,	<i>ural Language Processing, EMNLP 2024, Mi-</i>	797
743	Changling Liu, Yue Wu, Meng Lu, Sen Song,	<i>ami, FL, USA, November 12–16, 2024</i> , pages	798
744	and Yasin Abbasi Yadkori. 2025a. Hierarchi-	21424–21439. Association for Computational	799
745	cal reasoning model. <i>arXiv preprint arXiv:</i>	Linguistics.	800
746	<i>2506.21734</i> .	Mo Yu, Shiyu Chang, Yang Zhang, and Tommi	801
747	Yanting Wang, Wei Zou, Runpeng Geng, and	Jaakkola. 2019. Rethinking cooperative rational-	802
748	Jinyuan Jia. 2025b. TracLLM: A generic frame-	ization: Introspective extraction and complement	803
749	work for attributing long context LLMs. <i>arXiv</i>	control . In <i>Proceedings of the 2019 Conference</i>	804
750	<i>preprint arXiv: 2506.04202</i> .	<i>on Empirical Methods in Natural Language Pro-</i>	805
		<i>cessing and the 9th International Joint Confer-</i>	806
		<i>ence on Natural Language Processing (EMNLP-</i>	807

808 *IJCNLP*), pages 4094–4103, Hong Kong, China.
809 Association for Computational Linguistics.

810 Wenbo Zhang, Tong Wu, Yunlong Wang, Yong
811 Cai, and Hengrui Cai. 2023. [Towards trustworthy explanation: On causal rationalization](#). In
812 *Proceedings of the 40th International Confer-*
813 *ence on Machine Learning*, volume 202 of *Pro-*
814 *ceedings of Machine Learning Research*, pages
815 41715–41736. PMLR.

817 Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao
818 Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang,
819 Dawei Yin, and Mengnan Du. 2024. Explainabil-
820 ity for large language models: A survey. *ACM*
821 *Transactions on Intelligent Systems and Technol-*
822 *ogy*, 15(2):1–38.

823 Andy Zou, Zifan Wang, Nicholas Carlini, Milad
824 Nasr, J. Zico Kolter, and Matt Fredrikson. 2023.
825 Universal and transferable adversarial attacks on
826 aligned language models. *arXiv preprint arXiv:*
827 *2307.15043*.

828 A Algorithms

829 [Algorithm 3](#) separates the context, seen as
830 a list of sentences, into a maximum of k
831 chunks, each of which contains a minimum of
832 $total_length/k$ characters except the last. It
833 also produces a mapping from the block indices
834 to the set of sentences indices they contain.

835 B Results on other datasets

836 In this section are provided the necessity
837 and sufficiency scores for the first sentences,
838 grouped and not grouped, of the remaining
839 datasets (Longbench, Loogle).

840 We compute our metric from [Equation \(7\)](#)
841 with $\alpha = 0.5$ for every sentence in a dataset and
842 get a “Ground Truth” ranking. These values are
843 also computed for the first ten sentences given
844 by each method. The box plots for these values
845 are given in [Figures 5, 6, 8 and 10](#) for datasets
846 1, 2a, 2b, and 3, respectively.

847 B.1 Hotpot-QA

848 When the sentences are scored individually for
849 HotpotQA in [Figure 5](#), we can see that the first
850 is ranked better by TracLLM, then the following
851 ones are ranked better by TreeFinder.

852 B.2 LongBench-v2

853 For the sufficiency in [Figure 7](#), we can observe
854 that ContextCite is leading for the first sentence,

Algorithm 3 Chunkify

```
1: Input: Context  $C = (s_1, \dots, s_n)$ , char-  
   character lengths  $\text{len}(s_i)$ , maximum number of  
   chunks  $k > 1$   
2: Output:  $(\mathcal{C}_k, \text{map})$  where  $\mathcal{C}_k =$   
    $[C^{(1)}, \dots, C^{(m)}]$  and  $\text{map} =$   
    $[I^{(1)}, \dots, I^{(m)}]$  with  $I^{(j)} = [i \mid s_i \in$   
    $C^{(j)}]$   
3:  $L \leftarrow \sum_{i=1}^n \text{len}(s_i)$   
4:  $B \leftarrow \lceil \frac{L}{k} \rceil$  ▷ Char. per chunk  
5:  $\mathcal{C}_k \leftarrow []$ ,  $\text{map} \leftarrow []$   
6:  $C' \leftarrow []$ ,  $I' \leftarrow []$ ,  $\ell \leftarrow 0$   
7: for  $i \leftarrow 1$  to  $n$  do  
8:   append  $s_i$  to  $C'$   
9:   append  $i$  to  $I'$   
10:   $\ell \leftarrow \ell + \text{len}(s_i)$   
11:  if  $\ell > B$  then  
12:    append  $C'$  to  $\mathcal{C}_k$   
13:    append  $I'$  to  $\text{map}$   
14:     $C' \leftarrow []$ ,  $I' \leftarrow []$ ,  $\ell \leftarrow 0$   
15:  end if  
16: end for  
17: if  $C'$  is not empty then  
18:   append  $C'$  to  $\mathcal{C}_k$   
19:   append  $I'$  to  $\text{map}$   
20: end if  
21: return  $(\mathcal{C}_k, \text{map})$ 
```

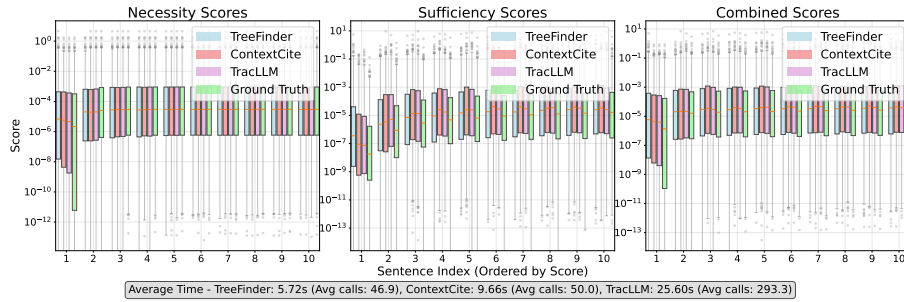


Figure 5: Sufficiency, necessity and average box plots for the top sentences individually in HotpotQA. Lower is better. The first sentence found with ContextCite and TracLLM is of high quality, while the following ones tend to bring little. TreeFinder misses the first but ranks the following ones better.

855 then TracLLM takes the lead. Afterwards, the
 856 curves interlace and the differences become too
 857 small to interpret meaningfully. In Figure 6,
 858 TracLLM has a pretty clear advantage in suffi-
 859 ciency.

860 B.3 Loogle Short

861 TracLLM consistently finds groups and indi-
 862 vidual sentences with lower necessity and suf-
 863 ficiency, as can be seen in Figures 8 and 9.
 864 When seen as individual sentences, the suffi-
 865 ciency scores remain far from the ground truth
 866 for all methods. Interestingly, when taken as a
 867 group, by the fourth sentence TracLLM reaches
 868 the ground-truth value.

869 We note that despite TreeFinder being worse
 870 in these graphs, it outperforms ContextCite in
 871 Figure 3.

872 B.4 Loogle Long

873 In Figure 10, TracLLM seems to perform bet-
 874 ter. In Figure 11, TreeFinder retains its slight
 875 advantage in necessity scores while TracLLM
 876 leads in sufficiency scores.

877 C Comparison with a specific example

878 In this section, we provide a sample from Hot-
 879 potQA where the answer provided by the model
 880 is correct along with the support sentences pro-
 881 vided in the dataset. We highlight the five best
 882 sentences as given by each method and match
 883 them with the support if possible in Tables 1
 884 to 3. The support sentences should not neces-
 885 sarily be matching the contributive attributions
 886 since they are corroborative attributions.

887 We observe that all methods find contributive

888 attributions that are aligned with the corrobora-
 889 tive attributions in these cases.

890 Question: All: “What Golden Globe nom-
 891 inated American actor from Juilliard School
 892 played a role in the 1986 romantic drama film,
 893 Children of a Lesser God, directed by Randa
 894 Haines?”

895 Support: “Children of a Lesser God is a
 896 1986 American romantic drama film directed
 897 by Randa Haines and written by Hesper An-
 898 derson and Mark Medoff.”, “An adaptation of
 899 Medoff’s Tony Award-winning stage play of
 900 the same name, the film stars Marlee Matlin (in
 901 an Oscar-winning performance) and William
 902 Hurt as employees at a school for the deaf: a
 903 deaf custodian and a hearing speech teacher,
 904 whose conflicting ideologies on speech and
 905 deafness create tension and discord in their de-
 906 veloping romantic relationship.”, “William Mc-
 907 Chord Hurt (born March 20, 1950) is an Ameri-
 908 can actor.”, “He received his acting training at
 909 the Juilliard School and began acting on stage in
 910 the 1970s.”, “Hurt made his film debut in 1980
 911 as a troubled scientist in Ken Russell’s science-
 912 fiction feature ”Altered States”, for which he re-
 913 ceived a Golden Globe nomination for New Star
 914 of the Year.”

915 Answers: TreeFinder and TracLLM: “The
 916 Golden Globe nominated American actor from
 917 the Juilliard School who played a role in the
 918 1986 romantic drama film *Children of a Lesser
 919 God*, directed by Randa Haines, is **William
 920 Hurt**. William Hurt received his acting train-
 921 ing at the Juilliard School and made his film de-
 922 but”. ContextCite: “The Golden Globe nomi-
 923 nated American actor from the Juilliard School
 924 who played a role in the 1986 romantic drama

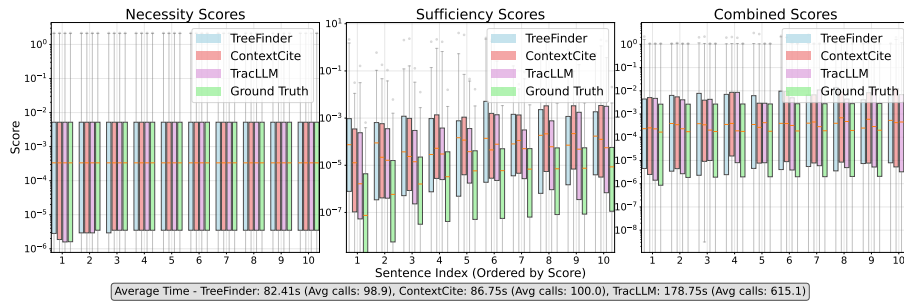


Figure 6: Sufficiency, necessity and Average box plots for the top sentences individually in LongBench-v2. Lower is better. The large error bars obscure any conclusion for the necessity. The sufficiency is better for TracLLM.

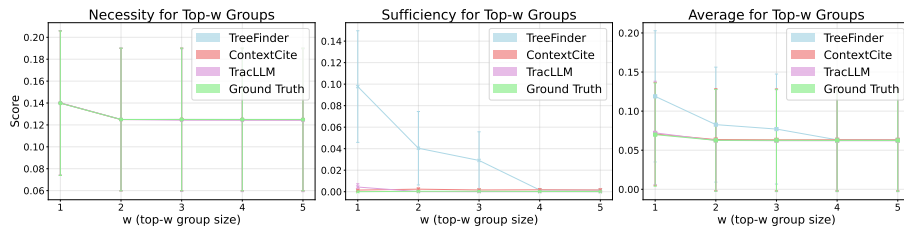


Figure 7: Sufficiency and necessity for the top sentences jointly in LongBench-v2. Lower is better. The large error bars obscure any conclusion for the necessity. The sufficiency of the rationale by TreeFinder is lacking a first but quickly converges to the same values as ContextCite and TracLLM.

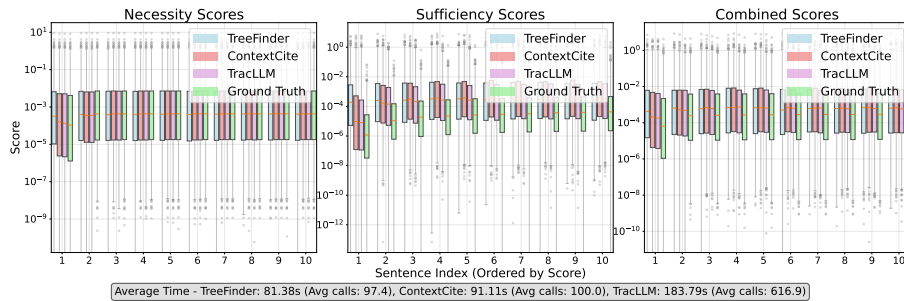


Figure 8: Sufficiency, necessity and average box plots for the top sentences individually in Loogle (short). Lower is better. TrackLLM consistently outperforms other methods.

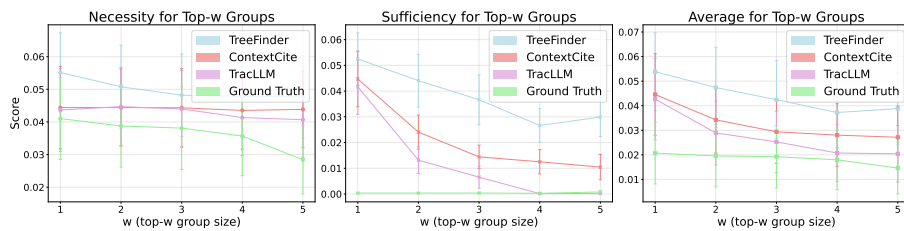


Figure 9: Sufficiency and necessity for the top sentences jointly in Loogle (short). Lower is better. TrackLLM consistently outperforms other methods.

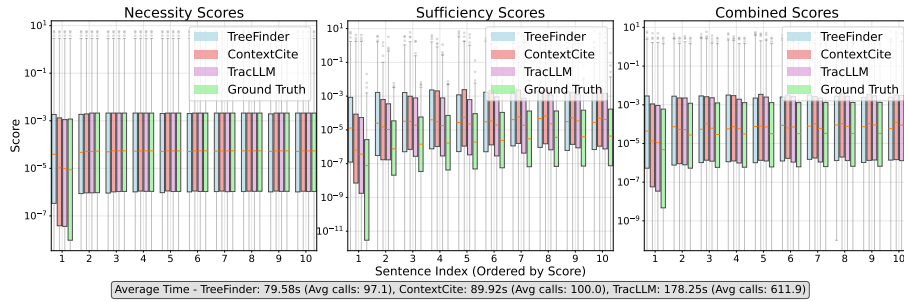


Figure 10: Sufficiency, necessity and average box plots for the top sentences individually in Loogle (long). Lower is better. TracLLM leads for the eight first sentences then is tied with TreeFinder. Both remain far from the ground truth.

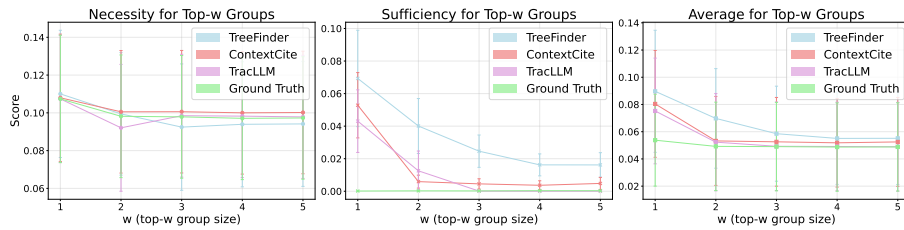


Figure 11: Sufficiency and necessity for the top sentences jointly in Loogle (long). Lower is better. There is a slight advantage in necessity for TreeFinder, but TracLLM has a clear lead in sufficiency.

925 film *Children of a Lesser God*, directed by
 926 Randa Haines, is **William Hurt**.*### Explanation:
 927 - **William Hurt** is mentioned as
 928 having received his acting training at the”

Rank	Sentence
1	#2
2	#3
3	#4
4	Two years later he made his breakthrough by starring in the romantic comedy ”Risky Business” (1983), which garnered Cruise his first nomination for the Golden Globe Award for Best Actor – Motion Picture Musical or Comedy.
5	In 1986, Cruise played a fighter pilot in the Tony Scott-directed action drama ”Top Gun” (the highest-grossing film that year), and also starred opposite Paul Newman in the Martin Scorsese-directed drama ”The Color of Money”.

Table 1: Top 5 sentences from TreeFinder.

Rank	Sentence
1	#3
2	#2
3	#1
4	The film was directed by Randa Haines, and was released directly on television.
5	#4

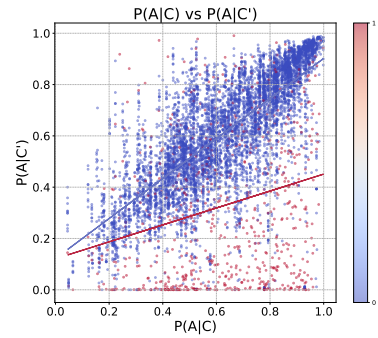
Table 2: Top 5 sentences from ContextCite.

Rank	Sentence
1	#3
2	#4
3	He subsequently played a leading role, as a lawyer who succumbs to the temptations of Kathleen Turner, in the neo-noir "Body Heat" (1981), and, as Arkady Renko, in Gorky Park (1983).
4	Tom Cruise is an American actor and producer who made his film debut with a minor role in the 1981 romantic drama "Endless Love".
5	She is perhaps most famous for directing the critically acclaimed feature film "Children of a Lesser God" (1986), which starred William Hurt and Marlee Matlin, for which Matlin won the 1987 Academy Award as Best Actress.

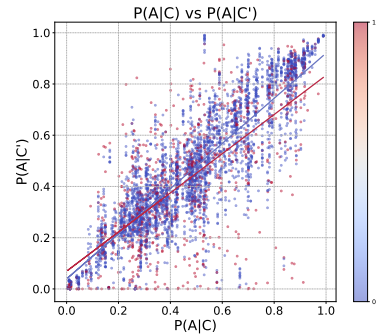
Table 3: Top 5 sentences from TracLLM.

D Proximity to corroborative attribution

We provide here an illustration of why the corroborative and contributive attributions are not the same, especially as the model has more difficulties responding correctly. To do this, we plot the probabilities of dropped chunks in Figure 12 against their original probabilities for each data point ($x = const$ in the figure) and check that we can separate the chunks that possess a corroborative attribution (red) from those that do



(a) Loogle (short)



(b) Loogle (long)

Figure 12: $P(A|C')$ against $P(A|C)$ in the case where $C' \subset C$ has a chunk removed. It is colored red if it contained part of the corroborative attribution and blue in the opposite. The chunk length is $|C'| = 8000$ characters. The difference of distribution shows the ability of the probability drop method to find corroborative attributions for a given model and dataset. Linear interpolations are drawn for reference.

929

930

931

932

933

934

935

936

937

938

939

940 not (blue).

941 In Loogle (short), the chunks that do not con-
942 tain a corroborative attribution have their proba-
943 bilities only slightly modified and remain close
944 to $x = y$, while the other are shifted toward
945 $y = 0$. This means that corroborative attribu-
946 tions are aligned with contributive attributions.

947 However, in Loogle (long), the two distribu-
948 tions are too similar, as indicated by both trend-
949 lines following $x = y$; removing a corrobo-
950 rative attribution no longer changes the prob-
951 ability of the answer. This can be the result
952 of the model missing these sentences in the
953 global context in the first place. Here, cor-
954 roborative and contributive attributions are no
955 longer aligned.

956 We highlight that discrepancies between the
957 two types of attributions are key to understand-
958 ing model uncertainty and should not be con-
959 fused for each other.