

TRACE: TRansformer-based Attribution using Contrastive Embeddings in LLMs

Anonymous ACL submission

Abstract

The rapid evolution of *large language models* (LLMs) represents a substantial leap forward in natural language understanding and generation. However, alongside these advancements come significant challenges related to the accountability and transparency of LLM outputs. Reliable source attribution is essential to adhering to stringent legal and regulatory standards, including those set forth by the General Data Protection Regulation. Despite the well-established methods in source attribution within the computer vision domain, the application of robust attribution frameworks to natural language processing remains underexplored. To bridge this gap, we propose a novel and versatile TRansformer-based Attribution framework using Contrastive Embeddings called TRACE that, in particular, exploits contrastive learning for source attribution. We perform an extensive empirical evaluation to demonstrate the performance and efficiency of TRACE in various settings and show that TRACE significantly improves the ability to attribute sources accurately, making it a valuable tool for enhancing the reliability and trustworthiness of LLMs.

1 Introduction

The recent era has seen a significant rise in the prevalence of *large language models* (LLMs) (Ouyang et al., 2022; Touvron et al., 2023) which have demonstrated an array of remarkable capabilities. However, studies (Huang et al., 2023; Liu et al., 2024; Wang et al., 2023a) have highlighted a critical concern on the accountability of LLMs. Considering the widespread usage and such a concern, it has brought to the forefront a critical need for source attribution that involves identifying the specific training data that contributes to generating part or all of an LLM’s output, which is crucial for legal and regulatory compliance and enhances the reliability of LLMs. Various regulations mandate transparency and accountability in data

usage, especially regarding intellectual property and privacy. For instance, the General Data Protection Regulation (GDPR) in the European Union requires that individuals have the right to be informed when their personal data is used. Proper source attribution ensures compliance with such legal frameworks, mitigating the risk of legal disputes and penalties.

A related topic would be that of *membership inference* (MI) (Miresghallah et al., 2022) whose task is to determine whether a given piece of data was used during the training of a machine learning model. While MI and source attribution share some similarities, they differ significantly in their granularity: MI typically only involves a binary classification task and does not require identifying a specific data provider. In contrast, source attribution requires to identify one or more data providers.

Though there are some studies on source attribution (Marra et al., 2018; Yu et al., 2022), a majority of them are situated within the computer vision domain. Techniques developed for computer vision tasks cannot be directly applied to LLMs due to the fundamental differences in the data and model architectures. To the best of our knowledge, effective source attribution for LLMs still remains an open and underexplored problem.

While numerous properties are important to a source attribution framework, we identify **accuracy**, **scalability**, and **interpretability** as the most crucial components. These three attributes are fundamental to ensuring the effectiveness and applicability of the framework across various contexts. Accuracy is essential to guaranteeing that the framework consistently produces reliable results. Scalability ensures that the framework can handle increasing volumes of data and complexity without a significant performance degradation, making it suitable for large-scale applications. Interpretability is equally critical as it enables stakeholders to un-

083 derstand and trust the attribution outcomes, hence
084 fostering transparency and facilitating informed de-
085 cision making.

086 This paper presents a novel TRansformer-
087 based Atribution framework using Contrastive
088 Embeddings (TRACE) to achieve source attribution
089 while satisfying the above three important proper-
090 ties. By detailing our methodology and presenting
091 empirical results, we seek to demonstrate the accu-
092 racy, scalability, and interpretability of TRACE.

093 Our contributions can be summarized as follows:

- 094 • We propose the novel TRACE framework based
095 on contrastive learning, which is designed to
096 achieve effective source attribution. TRACE
097 differs from traditional contrastive learning by
098 using source information as the label. Fig. 1
099 illustrates the TRACE framework.
- 100 • We have performed an extensive empirical
101 evaluation of TRACE to demonstrate its accu-
102 racy, scalability, and interpretability.

103 2 Preliminaries

104 **Contrastive Learning and NT-Xent Loss.** Con-
105 trastive learning is a conventional technique com-
106 monly used in representation learning (Arora et al.,
107 2019; Hadsell et al., 2006). Its underlying idea is
108 that similar objects should exhibit a closer distance
109 in the embedding space while dissimilar objects
110 should repel each other. This technique has been
111 widely employed in computer vision tasks due to its
112 convenient implementation to augment image input
113 to form a self-supervised problem. Models using
114 contrastive learning have achieved state-of-the-art
115 performances (Cui et al., 2021; Tian et al., 2020).
116 Apart from the attention it receives in computer
117 vision, new approaches using contrastive learning
118 in natural language processing (Meng et al., 2021;
119 Wu et al., 2020) have also started gaining attention
120 and showcasing great capabilities.

121 Our TRACE framework assigns the same label to
122 all the data from the same source, hence naturally
123 forming a *supervised contrastive learning* problem.
124 In particular, TRACE utilizes NT-Xent Loss (Sohn,
125 2016) for supervised contrastive learning:

$$126 \mathcal{L} = \sum_{i \in I} \frac{-1}{|P_i|} \sum_{p \in P_i} \log \left(\frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_p / \tau)}{\sum_{a \in A_i} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)} \right)$$

127 where the set I ($P_i \subset I$) contains indices of the sen-
128 tences in the given batch (sharing the same label as

129 sentence i , but does not include i), $A_i = I \setminus \{i\}$, \mathbf{z}_i
130 denotes the embedding of sentence i , and $\tau \in \mathbb{R}^+$
131 is a temperature parameter. Minimizing \mathcal{L} would
132 maximize the similarity between embeddings (of
133 sentences) from the same source while minimizing
134 the similarity between embeddings from different
135 sources.

Sentence Encoder. Similar to the concepts
136 of Word2Vec (Mikolov et al., 2013) and
137 GloVe (Brochier et al., 2019) which produce
138 meaningful vector representations of words, such
139 techniques can be applied to larger text units
140 such as sentences. A straightforward way is to
141 take the average of word embeddings within a
142 sentence, but this often results in embeddings
143 that lack semantic depth. Several models have
144 been developed to address this issue, including
145 InferSent (Conneau et al., 2018), Universal
146 Sentence Encoder (Cer et al., 2018), and
147 Sentence-BERT (SBERT) (Reimers and Gurevych,
148 2019). Given its superior performance and
149 efficiency, SBERT is chosen to generate sentence
150 embeddings in TRACE. SBERT leverages a pre-
151 trained BERT network and utilizes Siamese and
152 triplet network structures to produce semantically
153 meaningful sentence embeddings.
154

155 3 TRACE Framework

156 3.1 Source-Specific Semantic Distillation

157 Projecting every piece of data from each provider
158 into the embedding space is desirable but would
159 incur considerable computational costs. Moreover,
160 it is prudent to recognize that not all information
161 carries equal importance: For example, sentences
162 that occur less frequently typically tend to be more
163 representative of the document. So, we propose
164 to extract principal sentences from each source
165 by leveraging the *Term Frequency-Inverse Doc-*
166 *ument Frequency* (TF-IDF) which is effective for
167 identifying significant sentences within documents.
168 It is generally recommended to select 10-20% of
169 the sentences, thereby striking a balance between
170 complexity and performance; these sentences are
171 subsequently defined as *principal sentences*. The
172 length of these sentences is specified by a param-
173 eter called WINDOW_SIZE. Section 4.7 presents an
174 ablation study examining the effect of different
175 WINDOW_SIZEs on accuracy.

176 SBERT (Reimers and Gurevych, 2019) has
177 proven effective in deriving high-quality sentence

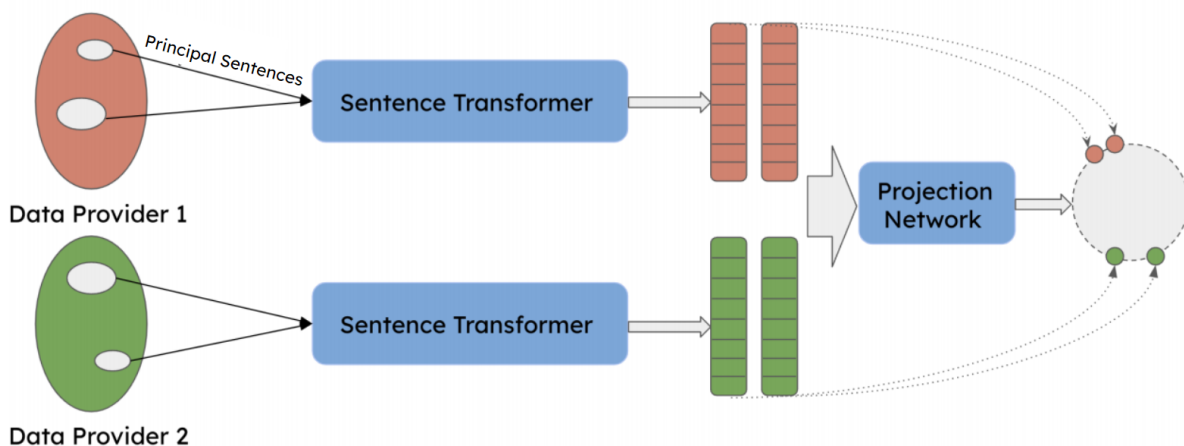


Figure 1: Illustration of TRACE framework.

178 representations. However, to enhance its suitability for TRACE, we propose several modifications
 179 inspired by the work of SimCLR (Chen et al., 2020).
 180 A key finding from SimCLR is that adding a non-
 181 linear projection head significantly improves the
 182 representation quality. Following this insight, we
 183 incorporate a projection network at the end of the
 184 traditional SBERT architecture. This projection net-
 185 work is trained together with the base SBERT model,
 186 thus encouraging the learned representations to be
 187 more discriminative in the embedding space.
 188

189 3.2 Supervised Contrastive Embedding 190 Training for Source-Coherent Clustering

191 Unlike the other contrastive learning frameworks
 192 in computer vision whose tasks are typically de-
 193 fined to be auto-regressive due to the availability of
 194 various data augmentation techniques to generate
 195 positive samples, TRACE aims to achieve source-
 196 coherent clustering. In our case, we already pos-
 197 sess the label of each sentence indicating its source.
 198 So, we can frame our task as a *supervised con-*
 199 *trastive learning* problem. The supervision is de-
 200 rived from the label information which corresponds
 201 to the source. Contrastive learning aligns with our
 202 objective to form clusters based on these various
 203 sources.

204 SimCLR has demonstrated that NT-Xent Loss
 205 outperforms other contrastive loss functions such
 206 as logistic loss (Mikolov et al., 2013) and margin
 207 loss (Schroff et al., 2015). So, we employ NT-Xent
 208 Loss as the loss function for TRACE.

209 3.3 Proximity-based Inference

210 Once the training phase is completed, we transition
 211 to the inference stage where each data source is rep-
 212 resented by its own set of contrastive embeddings.
 213 At this stage, when a language model generates a re-
 214 sponse, we employ the *k-Nearest Neighbor (kNN)*
 215 algorithm to assign the response to the closest data
 216 source in the embedding space, as demonstrated
 217 in Fig. 2. This ensures accurate source attribution
 218 by matching the generated response with its most
 219 similar source representation.

220 However, responses generated by language mod-
 221 els may not always be exclusively influenced by a
 222 single data source: there could be instances where
 223 information from multiple sources contributes to
 224 the generated text. To consider this possibility,
 225 we introduce the concept of *multi-source attribu-*
 226 *tion*. Multi-source attribution acknowledges and
 227 accounts for the potential influence of multiple data
 228 sources on the generated response.

229 We have developed three different implementa-
 230 tions for single-source attribution and multi-source
 231 attribution, which allow users to select the most
 232 appropriate inference method based on time con-
 233 straints and the number of sources. Section 4 pro-
 234 vides a comparison of these methods.

235 **Hard k NN (Single-Source Attribution).** *Hard*
 236 *kNN* follows the traditional k NN algorithm closely.
 237 Here, the attribution is determined by considering
 238 the k embeddings that are closest in distance to
 239 the query. The source that appears most frequently
 240 among these k neighbors is assigned as the source
 241 of the query.

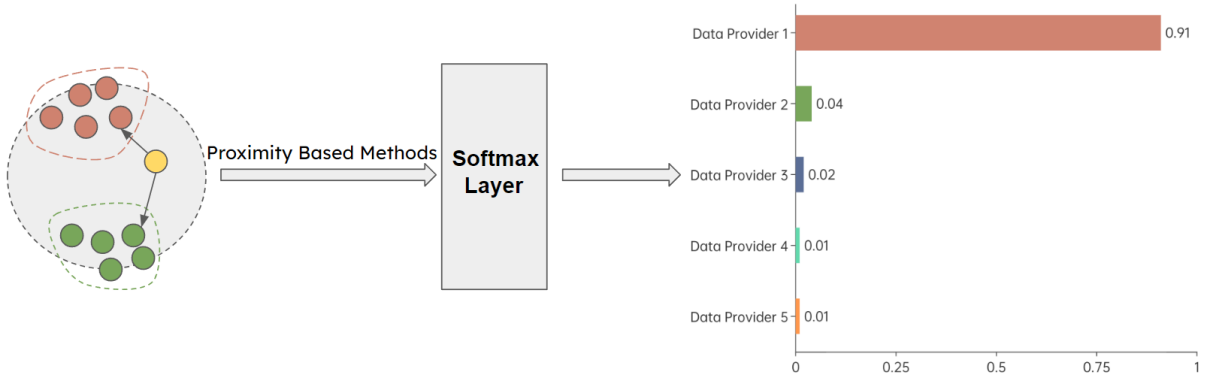


Figure 2: Illustration of the attribution step in TRACE framework.

Soft k NN (Multi-Source Attribution). To differentiate from traditional k NN where each query is assigned to a single source, we introduce *soft k NN*. Here, k represents the number of data sources rather than the number of closest neighbors. We rank the distances from the query to all other embeddings and select them in ascending order of distance until k distinct sources are covered.

Nearest Centroid (Single-Source Attribution). To reduce inference time, we employ the nearest centroid method. Here, the centroid of each cluster is calculated by adding the normalized embeddings within that cluster (i.e., corresponding to the sentences with the same label/source), as shown below. We then apply k NN using these centroids. This method significantly reduces inference time as it scales with the number of data providers rather than the volume of data from each source. We will demonstrate in the next section that this method maintains an impressively high accuracy.

Given a cluster of embeddings z_1, z_2, \dots, z_k with the same label/source, a good representative of the cluster would be the centroid \bar{z} that maximizes the sum of its cosine similarity with every normalized embedding z_i for $i = 1, \dots, k$. Equivalently, \bar{z} minimizes the sum of its standard cosine distance with every normalized embedding:

$$\begin{aligned} \sum_{i=1}^k \left(1 - \frac{z_i \cdot \bar{z}}{\|z_i\| \|\bar{z}\|} \right) &= k - \sum_{i=1}^k \frac{z_i \cdot \bar{z}}{\|z_i\| \|\bar{z}\|} \\ &= k - \left(\sum_{i=1}^k \frac{z_i}{\|z_i\|} \right) \cdot \frac{\bar{z}}{\|\bar{z}\|} \geq k - \left\| \sum_{i=1}^k \frac{z_i}{\|z_i\|} \right\| \cdot \frac{\bar{z}}{\|\bar{z}\|} \\ &\geq k - \left\| \sum_{i=1}^k \frac{z_i}{\|z_i\|} \right\| \end{aligned}$$

by Cauchy-Schwarz inequality. The equality holds

when there exists some $\lambda \in \mathbb{R}$ such that

$$\sum_{i=1}^k \frac{z_i}{\|z_i\|} = \lambda \frac{\bar{z}}{\|\bar{z}\|}.$$

In other words, \bar{z} can be obtained by adding all normalized embeddings and setting $\lambda = 1$:

$$\bar{z} = \sum_{i=1}^k \frac{z_i}{\|z_i\|}.$$

4 Experiments

4.1 Experimental Setup

Data. We perform an extensive empirical evaluation of TRACE using three datasets: booksum (Kryściński et al., 2022), dbpedia_14 (Zhang et al., 2015), and cc_news (Hamborg et al., 2017); a summary of these datasets can be found in Table 5 in the appendix. In the booksum dataset, we treat different books as distinct data providers and vary the number of data providers from 10, 25, 50, to 100 to demonstrate TRACE’s scalability to a large number of data providers. Similarly, each class in dbpedia_14 or each domain in cc_news is considered a separate data provider. In this section, we primarily present the experimental results on the booksum dataset with 25 data providers. Section 4.6 provides additional results.

Model. Focusing primarily on the booksum dataset, we evaluate the performance of TRACE using three different LLMs of varying sizes: t5-small-booksum (Raffel et al., 2020), GPT-2 (Radford et al., 2019), and Llama-2 (Touvron et al., 2023). The t5-small-booksum model is readily available on Hugging Face,¹ while GPT-2

¹<https://huggingface.co/cnlicu/t5-small-booksum>.

and Llama-2 have been fine-tuned on a subset of the booksum dataset. This setup allows us to assess the performance of TRACE across LLMs of different scales. In App. A, we provide more details about our experiments.

4.2 Visualization of TRACE’s Embedding Space

After training for 150 epochs on booksum, a visualization tool such as UMAP (McInnes et al., 2020) can be used to view the distribution of principal sentences. Fig. 3 shows that after the contrastive learning step, the desired outcome has been achieved, i.e., data coming from the same source form clear and distinct clusters. This validates that our contrastive learning successfully groups different data providers. Supposing the responses from an LLM are projected into the embedding space without incorporating the contrastive learning step, the resulting neighborhood exhibits chaos and it is challenging to derive robust information. This further demonstrates the importance of the contrastive learning step.

4.3 Accuracy

Evaluating the accuracy of source attribution is particularly challenging due to the inherent difficulty in obtaining ground-truth test datasets. Even with a dataset, a language model, and specific inputs, pinpointing the exact parts of the training data that influence a particular output remains complex. Here are the key reasons:

- Lack of Explicit Traceability.** Language models like LLMs generate outputs based on patterns learned from vast amounts of data. However, these models do not provide explicit traceability back to the specific training data. This means we cannot directly observe which parts of the training data contribute to a given output.
- Intermixed Training Data.** The training data for LLMs is often a massive, intertwined collection of texts from various sources. Disentangling these sources to identify the precise contribution of each segment to the final output is nearly impossible due to the sheer volume and complexity.
- Influence of Pre-training Data.** It is also likely that the model generates outputs based on data encountered during the pre-training

stage, which comprises a vast and diverse corpus. This pre-training data is often not fully documented or accessible, making it difficult to determine its influence on specific outputs during fine-tuning or evaluation.

Due to these challenges, obtaining ground-truth test datasets that accurately reflect the contribution of specific training data to the outputs of LLMs is exceedingly difficult. To address this issue, our approach involves using training data where the source is known. We then use this known source as the ground-truth label and evaluate whether TRACE can correctly determine the source. This allows us to approximate the evaluation of source attribution by leveraging the known origins of the specific training data.

Single-Source Attribution Accuracy. In this case, accuracy is simply defined as the number of correct source attributions divided by the total number of attributions evaluated, the latter of which is 250 in our experimental setup.

Multi-Source Attribution Accuracy. In certain settings, providing multiple sources and allowing the user to determine the justification of the attribution is acceptable. For a successful *soft kNN* attribution in such cases, the ground-truth source must appear among the *top-k* sources returned by TRACE. Using the same setup as that of single-source attribution, we have evaluated TRACE on 250 instances. Table 1 below shows the results:

Model	acc.	Soft <i>kNN</i>		Hard <i>kNN</i>		Nearest Centroid
		top-3 acc.	top-5 acc.	<i>k</i> = 10	<i>k</i> = 20	
t5	84.4%	95.3%	97.3%	84.4%	84.4%	84.4%
GPT-2	81.3%	92.3%	94.0%	81.3%	81.3%	81.3%
Llama-2	86.2%	96.1%	97.2%	86.2%	86.2%	86.2%

Table 1: Source attribution accuracy for 25 data providers on booksum dataset using TRACE.

It can be observed that the accuracy for models of different sizes remains consistently high and significantly surpasses the random guess’ accuracy of 4%. Another notable observation from the results is that varying the values of *k* in the hard *kNN* approach has minimal impact on accuracy and yields results identical to that of the nearest centroid method, which we attribute to the highly compact nature of the embeddings learned under the TRACE framework. When a query is projected into the embedding space, it becomes closely associated with its nearest neighbors regardless of

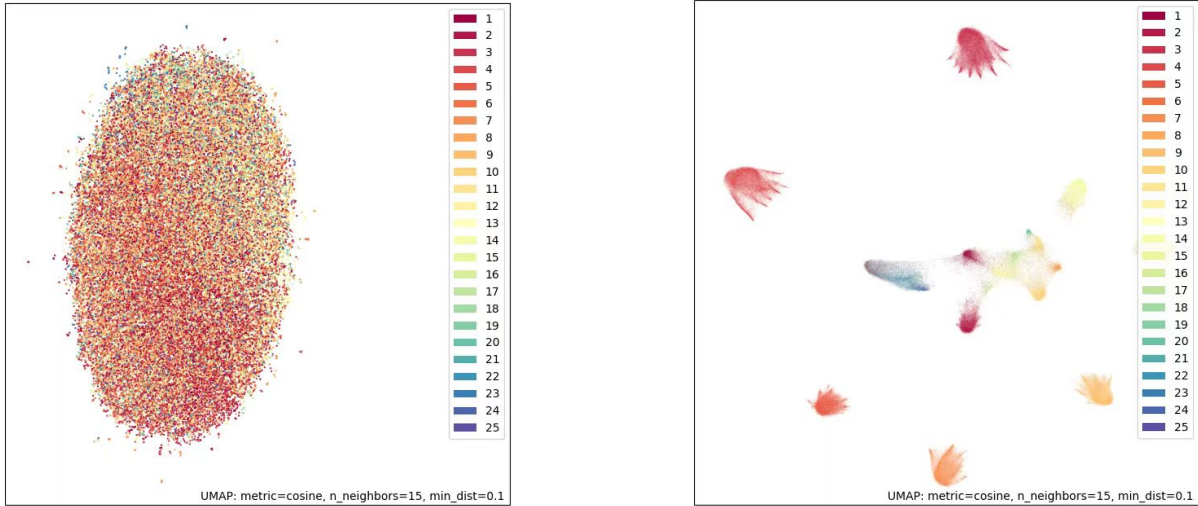


Figure 3: Visualization (using UMAP) of the embedding space before (left) and after (right) contrastive learning.

the specific value of k . This compactness suggests that the centroid of each cluster serves as an excellent representative of the entire cluster. Consequently, relying solely on these centroids can significantly reduce inference time. Even with 100 data providers as demonstrated in next subsection, the inference process remains almost instantaneous.

4.4 Scalability

Contemporary LLMs often necessitate substantial quantities of training data and the capability to manage a multitude of data providers. Hence, it is imperative to demonstrate the scalability of the TRACE framework under such settings. We assess the scalability of TRACE by selecting 10, 25, 50, and 100 distinct books from the booksum dataset, while maintaining a consistent experimental configuration. The results in Table 2 indicate a diminishing trend in accuracy with an increasing number of data providers, which is expected as the task complexity grows. However, despite this challenge, TRACE exhibits a relatively high level of accuracy across all settings, thus affirming its scalability.

4.5 Interpretability

The TRACE framework not only delivers accurate source attribution but also provides interpretability by offering additional insights into the attribution process. This interpretability is crucial for understanding the reasoning behind the model’s decisions and gaining confidence in its outputs. We illustrate the interpretability of TRACE using responses from the `t5-small-booksum` model as a demonstration.

Table 3 shows a summary of correctly

attributed single-source responses from the `t5-small-booksum` model. Each response is paired with the nearest principal sentence from the identified source. This pairing allows users to understand the specific evidence or context from the source text that influences the model’s attribution decision.

Moreover, TRACE offers interpretability through the inclusion of different similarity scores. These scores provide insights into the model’s confidence levels regarding the attribution outcomes. By examining the similarity scores, users can gauge the strength of the connection between the response and the identified source.

Overall, TRACE enhances interpretability by not only delivering the final attribution outcomes but also by providing supporting evidence from the source text and indicating the model’s confidence levels through similarity scores. This transparency and insight into the attribution process empower users to trust and understand the model’s outputs, which makes TRACE a valuable tool for source attribution tasks.

4.6 Additional Experimental Results

We conduct additional experiments to assess the performance of TRACE on alternative datasets, thereby evaluating its versatility. Table 4 summarizes the results. For a consistent comparison, we employ the same LLM across these datasets.

Our additional experiments affirm the adaptability of the TRACE framework across various datasets, thereby validating its applicability across various knowledge domains and settings.

n_books	t5			GPT2			Llama-2		
	acc.	top-3 acc.	top-5 acc.	acc.	top-3 acc.	top-5 acc.	acc.	top-3 acc.	top-5 acc.
10	87.5%	98.3%	99.4%	85.3%	96.8%	98.7%	88.2%	99.2%	99.5%
25	84.4%	95.3%	97.3%	81.3%	92.3%	94.0%	86.2%	96.1%	97.2%
50	73.1%	82.0%	84.0%	72.9%	82.9%	84.1%	70.3%	79.8%	82.2%
100	45.4%	74.8%	78.8%	49.0%	73.2%	77.7%	46.7%	76.8%	80.2%

Table 2: Source attribution accuracy for different no. of data providers on booksum dataset using TRACE.

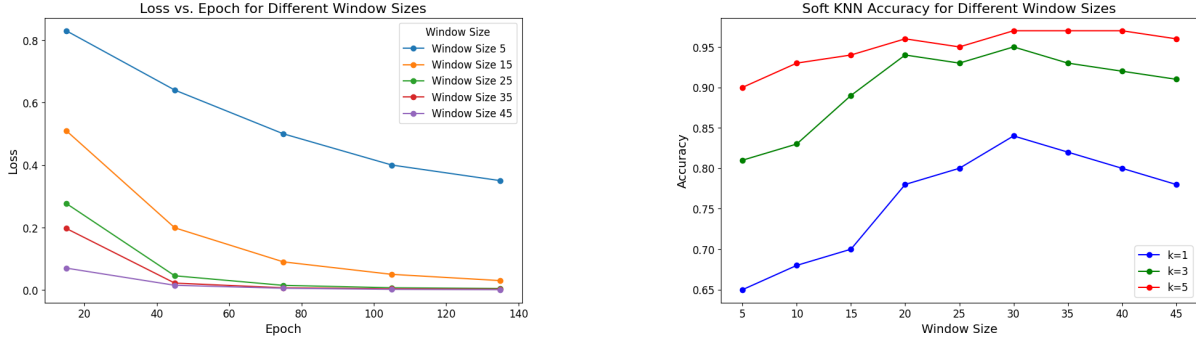


Figure 4: Contrastive loss (left) and soft k NN accuracy (right) with different WINDOW_SIZES. Note that the results for hard k NN (regardless of the value of k) are identical to that of soft k NN when $k = 1$.

4.7 Ablation Study

The most important factor in TRACE is the user-defined WINDOW_SIZE. If the WINDOW_SIZE is too small, the principal sentences cannot capture sufficient contextual information, hence deteriorating the performance. However, an exceedingly large WINDOW_SIZE will not only require more computational resources and time to train but also the meaning will be diluted by other redundant information. This presents a natural trade-off between source attribution performance and computational efficiency. Therefore, in this subsection, we will analyze this trade-off and present the results in Fig. 4.

It can be observed that a larger WINDOW_SIZE facilitates faster model convergence. However, model loss alone is not a comprehensive indicator of the clustering quality. So, we evaluate the source attribution accuracy on the test dataset. When the WINDOW_SIZE is set to 30, our TRACE framework achieves its highest accuracy. We hypothesize that this is primarily because the WINDOW_SIZE of 30 is sufficient to capture essential contextual information without excessively diluting it.

5 Related Work

Source Attribution. Though source attribution remains relatively underexplored in the domain of natural language processing, WASA (Wang et al., 2023b) stands out as a notable framework.² Oper-

²Note that neither the source code nor comprehensive details of the experimental setup have been provided in (Wang

ating on the principle of watermarking, WASA embeds distinct source identifiers within the training data to ensure that responses convey pertinent data provider information. However, WASA necessitates extensive manipulation of training data and training the entire LLM from scratch, which is a time-consuming process given their sizes. In contrast, TRACE distinguishes itself by being model-agnostic, i.e., requiring no knowledge about the model. This characteristic enhances efficiency and adaptability.

In the context of identifying information sources for quotes, Quobert (Vaucher et al., 2021) is a minimally supervised framework designed for extracting and attributing quotations from extensive news corpora. Additionally, Spangher et al. (2023) have developed robust models for identifying and attributing information in news articles. However, these approaches are primarily focused on specific domains such as news. In contrast, TRACE is designed to handle knowledge across a wide range of domains and hence provides a more generalized and versatile solution for source attribution tasks.

Information Retrieval. A related topic to our work here is information retrieval. Traditional retrieval techniques like BM25 (Robertson et al., 1994) hinge heavily on frequency-based rules which prove to be inadequate when dealing with responses that share semantic similarities without significant lexical overlap. More contemporary

et al., 2023b), making a fair comparison with WASA infeasible.

Response	Nearest Principal Sentence
Morel is in Sheffield, and he feels guilty towards Dawes, who is suffering and despairing, too. And besides, they had met in Nottingham in a way that is more or less responsible.	on his knees, feeling so awkward in presence of big trouble. Mrs. Morel did not change much. She stayed in Sheffield
But Emma thought at least it would turn out so. Mrs. Elton was first seen at church: but although devotion might be interrupted, curiosity could not be satisfied by a bride in. Pew, and it must be left for the visits in form which were then paid, to settle whether she was very pretty indeed, or only rather pretty at all.	or any thing just to keep my boot on.” Mr. Elton looked all happiness at this proposition; and nothing could exceed
to marry Lord Warburton. Isabel enquired. “Your uncle’s not an English nobleman,” said Mrs. Touchett in her smallest, sparest voice. The girl asked if the correspondent of the Interviewer was to take the party to London under Ralph’s escort. It was just the sort of plan, she said, that Miss Stackpole would be sure to suggest, and Isabel said that she did right to refuse him then.	he told Ralph he’s engaged to be married.” “Ah, to be married!” Isabel mildly exclaimed. “Unless he breaks it off. He seemed

Table 3: Sample responses with correct single-source attribution from t5-small-booksum model.

Dataset	Data Providers	Soft k NN			Hard k NN		Nearest Centroid
		$k = 1$	$k = 3$	$k = 5$	$k = 10$	$k = 20$	
booksum	10	85.3%	96.8%	98.7%	85.3%	85.3%	
dbpedia_14	10	88.2%	94.1%	97.2%	88.2%	88.2%	
booksum	25	81.3%	92.3%	94.0%	81.3%	81.3%	
cc_news	25	83.1%	90.8%	92.1%	83.1%	83.1%	

Table 4: Source attribution accuracy on dbpedia_14 and cc_news datasets using TRACE.

512 methods, such as ANCE (Xiong et al., 2020) and
513 Contriever (Izacard et al., 2022), opt for generat-
514 ing compact, dense representations of documents
515 rather than long, sparse ones. Thus, they tend to
516 achieve better results.

517 While information retrieval and TRACE both use
518 dense representations to measure text similarity,
519 they differ in objectives and applications. Informa-
520 tion retrieval aims to rank relevant documents for a
521 user’s query. In contrast, TRACE focuses on identi-
522 fying and attributing the original source of specific
523 information, hence ensuring accurate credit and
524 authenticity.

525 **Membership Inference Attack.** The concept of
526 membership inference attack was first introduced
527 by Shokri et al. (2017). The primary objective of
528 this attack is to ascertain whether a specific piece
529 of information was part of the training data for a
530 given machine learning model. Various assump-
531 tions about the available information lead to dif-
532 ferent attack models. For instance, some models
533 assume access to hard labels (Li and Zhang, 2021),
534 the model’s confidence scores (Watson et al., 2022;
535 Mattern et al., 2023), or the internal parameters
536 of the model (Leino and Fredrikson, 2020). Wei

537 et al. (2024) have achieved membership inference
538 by inserting watermarks into data. Despite the vari-
539 ations, these attacks fundamentally seek to answer
540 a binary question, i.e., whether the information was
541 included in the training dataset or not.

542 In contrast, source attribution entails mapping
543 the response to distinct and specific sources rather
544 than simply determining the presence or absence
545 of the data in the training set. Additionally, TRACE
546 adheres to a black-box setting: It does not require
547 access to internal information such as confidence
548 scores or model parameters. Instead, TRACE only
549 necessitates the output from a LLM.

550 6 Conclusion

551 This paper describes a novel TRACE framework
552 which effectively achieves source attribution. By
553 selecting principal sentences and projecting them
554 into the embedding space via source-coherent con-
555 trastive learning, TRACE enhances the interpretabil-
556 ity of responses generated by LLMs. This enhance-
557 ment also conforms to regulations that aim to pro-
558 tect the privacy of users. After evaluating TRACE
559 on various datasets, we have demonstrated the ac-
560 curacy and effectiveness of our framework.

Limitations. Our experiments are subject to some limitations that can be addressed in the future work to ensure a comprehensive interpretation of results. Firstly, the balanced distribution of data across different sources may impact the final inference of TRACE given its reliance on the k NN algorithm. This uniformity in data volume may not be representative of real-world settings, which potentially limits the generalizability of our findings. Secondly, information within each source is quite distinct with no overlapping data. Future works can verify the setting where data sources contain similar information. These limitations underscore the importance of future research in addressing such challenges to enhance the robustness of TRACE across varied data environments.

Ethical Considerations. Our TRACE framework introduces a method for achieving source attribution. Utilizing this framework, a malicious actor may potentially identify the sources of data providers and reveal sensitive information about them. Therefore, the application of TRACE within this context necessitates meticulous handling to mitigate privacy concerns.

References

- Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. 2019. A theoretical analysis of contrastive unsupervised representation learning. arXiv:1902.09229.
- Robin Brochier, Adrien Guille, and Julien Velcin. 2019. Global vectors for node representations. In *Proc. WWW*, pages 2587–2593.
- Daniel Cer, Yinfei Yang, Sheng yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Universal sentence encoder. arXiv:1803.11175.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. arXiv:2002.05709.
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. 2018. Supervised learning of universal sentence representations from natural language inference data. arXiv:1705.02364.
- Jiequan Cui, Zhisheng Zhong, Shu Liu, Bei Yu, and Jiaya Jia. 2021. Parametric contrastive learning. In *Proc. ICCV*, pages 715–724.
- Raia Hadsell, Sumit Chopra, and Yann Lecun. 2006. Dimensionality reduction by learning an invariant mapping. In *Proc. CVPR*, pages 1735–1742.
- Felix Hamborg, Norman Meuschke, Corinna Breiterger, and Bela Gipp. 2017. news-please: A generic news crawler and extractor. In *Proc. ISI*, pages 218–223.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. arXiv:2311.05232.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. arXiv:2112.09118.
- Wojciech Kryściński, Nazneen Rajani, Divyansh Agarwal, Caiming Xiong, and Dragomir Radev. 2022. BookSum: A collection of datasets for long-form narrative summarization. In *Proc. EMNLP Findings*, pages 6536–6558.
- Klas Leino and Matt Fredrikson. 2020. Stolen memories: Leveraging model memorization for calibrated white-box membership inference. In *Proc. SEC*, pages 1605–1622.
- Zheng Li and Yang Zhang. 2021. Membership leakage in label-only exposures. arXiv:2007.15528.

638	Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li.	Florian Schroff, Dmitry Kalenichenko, and James Philbin.	691
639	2024. Trustworthy LLMs: A survey and guideline for evaluating large language models' alignment.	2015. FaceNet: A unified embedding for face recognition and clustering. In <i>Proc. CVPR</i> ,	692
640	arXiv:2308.05374.	pages 815–823.	693
641			694
642			
643			
644	Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi.	Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov.	695
645	2018. Do GANs leave artificial fingerprints? arXiv:1812.11842.	2017. Membership inference attacks against machine learning models. In <i>Proc. IEEE S&P</i> ,	696
646		pages 3–18.	697
647			698
648	Justus Mattern, Fatemehsadat Miresghallah, Zhijing Jin, Bernhard Schölkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick.	Kihyuk Sohn.	699
649	2023. Membership inference attacks against language models via neighbourhood comparison. arXiv:2305.18462.	2016. Improved deep metric learning with multi-class N-pair loss objective. In <i>Proc. NIPS</i> .	700
650			
651		Alexander Spangher, Nanyun Peng, Emilio Ferrara, and Jonathan May.	701
652	Leland McInnes, John Healy, and James Melville.	2023. Identifying informational sources in news articles. In <i>Proc. EMNLP</i> ,	702
653	2020. UMAP: Uniform manifold approximation and projection for dimension reduction. arXiv:1802.03426.	pages 3626–3639.	703
654			704
655		Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola.	705
656	Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song.	2020. What makes for good views for contrastive learning? arXiv:2005.10243.	706
657	2021. COCO-LM: Correcting and contrasting text sequences for language model pretraining. arXiv:2102.08473.		707
658		Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample.	708
659		2023. LLaMA: Open and efficient foundation language models. arXiv:2302.13971.	709
660	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean.		710
661	2013. Efficient estimation of word representations in vector space. arXiv:1301.3781.		711
662			712
663	Fatemehsadat Miresghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri.	Timoté Vaucher, Andreas Spitz, Michele Catasta, and Robert West.	713
664	2022. Quantifying privacy risks of masked language models using membership inference attacks. arXiv:2203.03929.	2021. Quotebank: A corpus of quotations from a decade of news. In <i>Proc. WSDM</i> ,	714
665		pages 328–336.	715
666			716
667		Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, Sang T. Truong, Simran Arora, Mantas Mazeika, Dan Hendrycks, Zinan Lin, Yu Cheng, Sanmi Koyejo, Dawn Song, and Bo Li.	717
668	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe.	2023a. DecodingTrust: A comprehensive assessment of trustworthiness in GPT models. In <i>Proc. NeurIPS</i> .	718
669	2022. Training language models to follow instructions with human feedback. arXiv:2203.02155.		719
670		Jingtang Wang, Xinyang Lu, Zitong Zhao, Zhongxiang Dai, Chuan-Sheng Foo, See-Kiong Ng, and Bryan Kian Hsiang Low.	720
671	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.	2023b. WASA: Watermark-based source attribution for large language model-generated data. arXiv:2310.00646.	721
672	2019. Language models are unsupervised multitask learners . Technical report, OpenAI.		722
673		Lauren Watson, Chuan Guo, Graham Cormode, and Alex Sablayrolles.	723
674		2022. On the importance of difficulty calibration in membership inference attacks. arXiv:2111.08440.	724
675			725
676		Johnny Tian-Zheng Wei, Ryan Yixiang Wang, and Robin Jia.	726
677	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu.	2024. Proving membership in LLM pretraining data via data watermarks. arXiv:2402.10892.	727
678	2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>JMLR</i> , 21(1):5485–5551.		728
679		Zhuofeng Wu, Sinong Wang, Jiatao Gu, Madian Khabsa, Fei Sun, and Hao Ma.	729
680		2020. CLEAR: Contrastive learning for sentence representation. arXiv:2012.15466.	730
681	Nils Reimers and Iryna Gurevych.		731
682	2019. SentenceBERT: Sentence embeddings using Siamese BERT-networks. arXiv:1908.10084.		732
683			733
684			734
685			735
686			736
687			737
688	Stephen Robertson, Steve Walker, Susan Jones, Michelle Hancock-Beaulieu, and Mike Gatford.		738
689	1994. Okapi at TREC-3. In <i>Proc. TREC</i> , pages 109–126.		739
690			740

745 Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang,
746 Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold
747 Overwijk. 2020. Approximate nearest neighbor neg-
748 ative contrastive learning for dense text retrieval.
749 arXiv:2007.00808.

750 Ning Yu, Vladislav Skripniuk, Sahar Abdelnabi, and
751 Mario Fritz. 2022. Artificial fingerprinting for gener-
752 ative models: Rooting deepfake attribution in training
753 data. arXiv:2007.08457.

754 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.
755 Character-level convolutional networks for text clas-
756 sification. In *Proc. NeurIPS*.

A Experimental Setup 757

Data Preparation. From booksum, we have ran- 758
759 domly selected subsets of 10, 25, 50, and 100
760 books. From dbpedia_14, we chose 10 distinct
761 classes. Additionally, we have extracted text sam-
762 ples from 25 diverse domains within the cc_news
763 dataset. 764

Before proceeding with the analysis, we have 765
766 performed standard preprocessing steps which in-
767 clude converting all text to lowercase and removing
768 punctuation to ensure uniformity and cleanliness
769 in the data. 770

Model. For sentence embedding, we have 771
772 opted for SBERT (Reimers and Gurevych,
773 2019). Leveraging the pre-trained model
774 xlm-r-distilroberta-base-paraphrase-v1
775 that is readily accessible on Hugging Face, we
776 have fine-tuned it within our TRACE framework.
777 Moreover, we have augmented the model with
778 additional feed-forward layers which serve as
779 the projection network. The dimension for the
780 embeddings is set as 64. 781

Training Details. The hyperparameters utilized 782
783 in our experimental setup are configured as follows:
784 the learning rate is 1×10^{-5} , the batch size is 64,
785 the number of epochs is 150, and the temperature
786 in the NT-Xent Loss is 0.1. Notably, all train-
787 ing procedures are conducted on a single NVIDIA
788 L40 GPU, obviating model or data parallelism tech-
789 niques. The results were obtained by averaging the
790 outcomes of three executions, each with a different
791 random seed. 792

Statistic	booksum	dbpedia_14	cc_news
Number of Documents	405 (books)	560,000	149,954,415
Languages Covered	English	English	English
Domains	Books	Encyclopedic	News

Table 5: Statistics of booksum, dbpedia_14, and cc_news datasets.