

CAVE: Controllable Authorship Verification Explanations

Sahana Ramnath[♥] Kartik Pandey[♥] Elizabeth Boschee[♣] Xiang Ren[♥]

[♥]Department of Computer Science, University of Southern California

[♣]Information Sciences Institute, University of Southern California

sramnath@usc.edu

Abstract

Authorship Verification (AV) (do two documents have the same author?) is essential in many real-life applications. AV is often used in privacy-sensitive domains that require an offline proprietary model that is deployed on premises, making publicly served online models (APIs) a suboptimal choice. Current offline AV models however have lower downstream utility due to limited accuracy (eg: traditional stylometry AV systems) and lack of accessible post-hoc explanations. In this work, we address the above challenges by developing a *trained, offline* model CAVE (Controllable Authorship Verification Explanations): CAVE generates free-text AV explanations that are controlled to be (1) *accessible* (uniform structure that can be decomposed into sub-explanations grounded to relevant linguistic features), and (2) easily verified for explanation-label *consistency*. We generate silver-standard training data grounded to the desirable linguistic features by a prompt-based method **PROMPT-CAVE**. We then filter the data based on rationale-label consistency using a novel metric **CONS-R-L**. Finally, we fine-tune a small, offline model (LLAMA-3-8B) with this data to create our model CAVE. Results on three difficult AV datasets show that CAVE generates high quality explanations (as measured by automatic and human evaluation) as well as competitive task accuracy¹.

1 Introduction

Authorship Verification (AV) (Koppel and Winter, 2014) is the NLP task of determining if two input documents were written by the same author. AV is used for tasks such as plagiarism detection, forensic analysis (often in support of law enforcement), analysis of the spread of misinformation. Given the sensitive nature of this task, it is imperative to develop AV methods that are explainable, accessible and secure. Explainability ensures that the system’s

¹github.com/INK-USC/Controllable-AV-Explanations

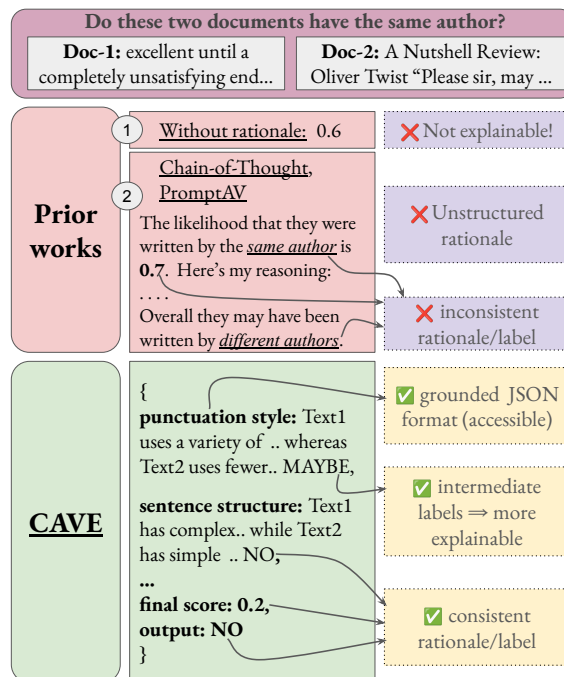


Figure 1: CAVE generates uniformly structured free-text explanations grounded in relevant linguistic features, that can be automatically verified for consistency.

decisions can be understood and trusted by the user. Accessibility ensures easier downstream usability; for example, a well-structured free-text explanation with clearly defined topic points is more easily understood than unstructured explanations such as long lists of low-level features. Finally, security is necessary to safely use the AV system for proprietary data; a locally-hosted model on the user’s server is more secure than an API call to an online model, which could potentially store the user’s sensitive data. AV is also an extremely difficult task for humans, in contrast to other NLP tasks such as machine translation; this further necessitates the need for automatic, high-quality AV systems.

Early works performed AV by comparing hand-crafted features such as n-grams (van Dam, 2013), POS tags (Moreau and Vogel, 2013), LIWC fea-

```

{
  punctuation style: Both texts use a variety of punctuation, including commas, periods, and quotation marks, but Text 1 uses
  more diverse punctuation such as parentheses and hyphens. MAYBE
  special characters style, capitalization style: Text 1 uses continuous capitalization for emphasis (e.g., 'WONDERFUL',
  'THRILLED'), which is not observed in Text 2. NO ,
  acronyms and abbreviations: Neither text makes significant use of acronyms or unusual abbreviations. YES ,
  writing style: Text 1 has a more personal, reflective style, sharing personal opinions and feelings about the movie. Text 2
  provides a more detached, narrative-style review without personal input. NO ,
  expressions and idioms: Both texts avoid colloquial expressions and idioms, opting for a more formal format. YES ,
  tone and mood: Text 1 has a more varied tone, from enthusiasm to disappointment, while Text 2 maintains a consistent,
  somewhat formal and analytical tone. NO ,
  sentence structure: Text 1 features a mix of short and long sentences with more complex structures, while Text 2 tends to use
  more uniformly structured, intermediate-length sentences. MAYBE ,
  any other relevant aspect: The approach to movie critique is different; Text 1 is more about the impact on the viewer, while
  Text 2 focuses on plot summary and cinematic elements. NO ,
  final score: 0.375 ,
  output: NO
}

```

Table 1: Output structure of CAVE, PROMPT-CAVE: The text (JSON keys) highlighted in blue depict the linguistic features used in the analysis, as well as the keys for final confidence score and predicted task output. The text underlined and highlighted in yellow depict the intermediate labels with respect to each linguistic feature, the confidence score value, and the predicted task label. We use the predicted label to measure task accuracy and all of them to measure rationale-label consistency.

tures (Uchendu et al., 2020); while these are explainable, they are hard to scale. Subsequent works employed neural architectures to improve scalability and accuracy, such as the usage of Siamese networks (Araujo-Pino et al., 2020; Najafi and Tavvan, 2022; Boenninghoff et al., 2019) to assign authorship via distance between document embeddings. These methods, however, have little or no post-hoc explainability. This poses a huge problem: one cannot, for instance, make legal assertions about the authorship of a set of questionable documents solely on the basis of “a system said so” - decision makers (judges, university officials, intelligence analysts, etc.) need to know *why*. Recent years have seen the surge of large language models and their self-rationalization capabilities: Hung et al. (2023); Huang et al. (2024b) generate free-text rationales² for AV from models such as GPT-4 (OpenAI, 2024b) in a zero/few-shot manner. These methods however face issues such as unreliability of online models (varying results/details across runs, deprecation of APIs), inconsistency of rationales with predicted label (C), and high expenses.

In this work, we propose our model CAVE (Controllable Authorship Verification Explanations), an *in-house* model that generates

authorship predictions along with high-quality, *free-text rationales* that explain the former post-hoc. The rationales are structured as a set of sub-explanations (refer Table 1) grounded to AV-relevant linguistic features (Hung et al., 2023; Boenninghoff et al., 2019); these sub-explanations have corresponding (intermediate) labels that provide further structure and also serve as a means to verify overall consistency with CAVE’s authorship prediction.

Since there is no human-written AV explanation data that can be used to train our in-house model, we adopt the model distillation strategy: (West et al., 2022; Li et al., 2023a, 2024) we engineer **PROMPT-CAVE** to generate silver-standard train data in the format described above by prompting a large, oracle model (GPT-4-TURBO OpenAI (2024a)), and use that to supervised-finetune a smaller language model (LLAMA-3-8B Meta-AI (2024)). We also introduce our automatic metric **CONS-R-L** that measures the consistency between the rationale and the model’s final label; CONS-R-L is used to filter the silver train data before finetuning (Ramnath et al., 2023; Li et al., 2023a), and also assess CAVE during inference. Our experiments with three difficult AV datasets show that CAVE leads to competitive task performances as

²We use terms explanation and rationale interchangeably.

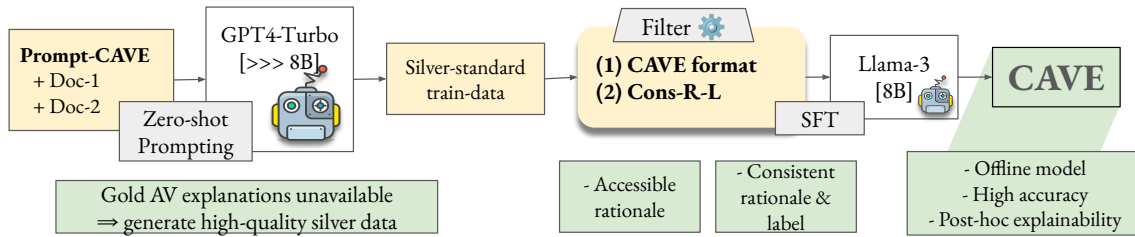


Figure 2: Pipeline to train CAVE: We obtain silver train data from GPT-4-TURBO using PROMPT-CAVE, filter it according to CONS-R-L and our output format. We then supervised-finetune a LLAMA-3-8B with the filtered data.

well as high quality of rationales.

2 Our Method

To enable and improve downstream utility of authorship verification systems, we present **CAVE**, our trained, *offline* model that generates accessible free-text rationales (well-structured, and uniform structure across all datapoints). We further propose a novel metric **CONS-R-L** that leverages the structure of CAVE’s rationales to measure consistency between a rationale and its task (AV) label. Figure 2 demonstrated the visual pipeline to train CAVE.

Silver-standard data. There exists no gold-standard, human-written datasets for AV *explanations*; hence, the first step to develop an AV model that can generate explanations is to create *silver-standard* data for it. We propose our zero-shot prompt, **PROMPT-CAVE**, that can be used with LLMs such as GPT-4-TURBO to generate silver training data. We adapt and extend PROMPTAV (Hung et al., 2023) to engineer PROMPT-CAVE. We then use this silver data to train an in-house language model CAVE to generate AV explanations.

PROMPT-CAVE and PROMPTAV. Prior work PROMPTAV proposes a zero/few-shot prompt to generate free-text rationales that contain linguistic features relevant for AV³. To our best knowledge, PROMPTAV was the first work in contemporary literature to generate natural language AV explanations (the second and latest work is LIP (Huang et al., 2024b), which uses a different compilation of linguistic features)⁴. While PROMPTAV leads to high AV accuracies, its rationales are (1) inaccessible towards automatic evaluation⁵, and (2) potentially *inconsistent* with their final label (detailed discussion and example in Appendix C),

³adapted and compiled from Boenninghoff et al. (2019)

⁴We discuss a contemporaneous work (Hu et al.) in Appendix B.1

⁵this is reinforced by their reporting only task accuracy and no rationale metrics

both of which makes the rationales unusable. In this work, we propose PROMPT-CAVE, an extension of PROMPTAV that addresses these issues: (1) PROMPT-CAVE generates rationales that are uniformly structured across datapoints; it generates a sub-rationale and label for each linguistic feature of analysis, under a fixed JSON format. This ensures ease of understanding, and enables any automatic quality measurement. (2) The aforementioned intermediate labels, along with the final score and the final label (also prompted for by PROMPT-CAVE) allow for an *automatic consistency evaluation* (described below), which leads to higher transparency and awareness. Table 10 shows PROMPT-CAVE.

Output structure (Table 1). We define the *correct* and *complete* output structure of CAVE/PROMPT-CAVE to be a JSON structure where the linguistic features are keys; the model has to discuss the authorship of the texts grounded to each feature independently, via corresponding explanations and intermediate labels (YES/NO/MAYBE). The model also has to provide a score $\in [0, 1]$ that indicates its overall confidence in the shared authorship of the documents, and an overall predicted task label (YES/NO).

CONS-R-L, automatic consistency metric. Rationales that are inconsistent with the model’s predicted label are risky and unusable in downstream applications (Lyu et al., 2023). Prior works in self-rationalization check for consistency either by human analysis (Chen et al., 2023b) or by training models (Hase et al., 2020; Chen et al., 2023a; Wiegrefe et al., 2021) for the same. In our output format, the presence of intermediate labels allows the automatic verification of rationale-label consistency. We define a consistency metric **CONS-R-L** using two binary variables CS-1 and CS-2: $\text{CONS-R-L} \iff (\text{CS-1} \wedge \text{CS-2})$.

CS-1: Is the final score (which indicates the model’s confidence in shared authorship) consistent with the final label? If the final label is ‘YES’, is

the final score ≥ 0.5 ? Similarly, if the final label is ‘NO’, is the final score ≤ 0.5 ?

CS-2: Are the intermediate labels faithful (overall) to the final label? Since the model output is structured as a JSON, we can automatically extract the intermediate labels for each feature and count them. If the final label is ‘YES’ and $\text{count}(\text{YES}) + \text{count}(\text{MAYBE}) > \text{count}(\text{NO})$, then the labels are consistent. Similarly, if the final label is ‘NO’ and $\text{count}(\text{NO}) + \text{count}(\text{MAYBE}) > \text{count}(\text{YES})$ then the labels are consistent⁶.

We define that is, a given output is consistent (CONS-R-L = 1) only if the final score, final label and the aggregate of the intermediate labels all align with each other.

Training CAVE. We generate silver training data (D_{train}) from GPT-4-TURBO (T) using PROMPT-CAVE (P).

$$D_{train} = \{(x_i, y_i^t)\} \sim T(x_i|P) \quad (1)$$

where x_i represents the input document pairs and y_i^t represents the outputs (rationale + label) created from T . We then filter D_{train} according to rationale-label consistency (CONS-R-L) and conformity to our desired structure, to get our final train set F_{train} . Finally, we supervised fine-tune a smaller, offline language model (S) with F_{train} ; for a decoder-only model like LLAMA-3-8B, we train the model only on the output tokens.

$$E_{(x,y^t) \sim F_{train}} S(y^t|x) \quad (2)$$

3 Experiments and Results

3.1 Datasets, Training, Inference

We conduct experiments on three popular AV datasets: IMDB62 (Seroussi et al., 2014), BLOG-AUTH (Schler et al., 2006) and FANFICTION (via the PAN-2020 Authorship Verification). Our datasets are diverse in domain (IMDB62: movie reviews, BLOG-AUTH: website blogs, FANFICTION: movie/book fanfiction); we carefully select these datasets to test the generalizability of our method. For IMDB62 we use the test set of 1k text pairs from Hung et al. (2023). For BLOG-AUTH, we create our own test set by sampling 1k document pairs from the original blog authorship corpus. For FANFICTION, we sample 1k document pairs from the PAN-2020 test set (which had 14k pairs). Our test sets have equal split of YES/NO labels.

⁶at this point, we assume all linguistic features have equal weightage towards shared authorship.

Creating the train set. For IMDB62 and BLOG-AUTH, the original datasets consist of documents and their anonymous author IDs; we sample document pairs with equal amounts of same/different authors to create trainsets. For FANFICTION, we take the ‘small’ train set provided by PAN-2020, and sample data from the same. We start off with a train set of $\sim 1.8k$ documents pairs for each dataset⁷ and combine them with PROMPT-CAVE and get silver rationale data from GPT-4-TURBO (temperature 0.0 and 2 responses per sample). Finally, we filter them according to the criteria in §2 and use them for training. Appendix D details the data statistics after each stage of filtering.

Training and Inference. We supervised, fine-tune a LLAMA-3-8B with LoRA (Hu et al., 2021). We report all hyperparameters in Appendix E. We use greedy decoding during inference.

3.2 Baselines

We present and compare with multiple baselines:

Traditional stylometry. We present a weak baseline (Ikae, 2021) from traditional stylometry. We represent input document pairs as the concatenation of their TF-IDF or word-count⁸ features, and train a Support Vector Machine (Cortes and Vapnik, 1995) to classify their shared authorship status (hyperparameters in Appendix E).

Zero-shot, prior work. We present three zero-shot baselines with GPT-4-TURBO (temperature 0.0) and LLAMA-3-8B-INSTRUCT (greedy decoding). The first baseline is chain-of-thought COT (Wei et al., 2022), a standard self-rationalization method that does not include any AV specific instructions. The second baseline is PROMPTAV (Hung et al., 2023). We use the prompts for COT and PROMPTAV reported in Hung et al. (2023) (Appendix C). Both these baselines report confidence scores (0-1, higher score means more confidence in shared authorship) in lieu of labels. We extract this confidence score by string matching for decimal numbers in the output; we measure accuracy by thresholding these scores at 0.5. *Consistency cannot be calculated for these two baselines.* The third zero-shot baseline is PROMPT-CAVE itself.

Fine-tuned. Similar to how we train CAVE, we generate silver COT and PROMPTAV data from GPT-4-TURBO, filter them according to task cor-

⁷we ensure that there is no leakage between train/test datasets.

⁸sklearn.feature_extraction.text TfidfVectorizer, CountVectorizer

Model	Method	Test →	IMDB62		BLOG-AUTH		FANFICTION	
		Train ↓	ACC.	CONS-R-L	ACC.	CONS-R-L	ACC.	CONS-R-L
SVM	-	Fine-tune	59.7	-	56.2	-	55.7	-
GPT-4-TURBO [>>> 8B]	CoT	Zero-shot	73.6	-	62.0	-	54.2	-
	PROMPTAV	Zero-shot	73.0	-	62.2	-	57.0	-
	PROMPT-CAVE	Zero-shot	67.6	0.86	57.7	0.73	52.2	0.67
	PROMPT-CAVE	Few-shot	71.3	0.94	64.0	0.97	61.7	0.86
LLAMA-3-8B -INSTRUCT	CoT	Zero-shot	62.4	-	60.7	-	57.5	-
	PROMPTAV	Zero-shot	62.8	-	57.9	-	53.8	-
	PROMPT-CAVE	Zero-shot	36.7	0.16	38.8	0.25	45.5	0.17
	PROMPT-CAVE	Few-shot	53.8	0.98	62.3	0.99	48.1	0.93
LLAMA-3-8B	CoT	Fine-tune	65.7	-	64.6*	-	60.5	-
	PROMPTAV	Fine-tune	76.2*	-	61.3	-	58.7	-
	CAVE	Fine-tune	74.1	0.99*	60.6	1.00*	62.6*	0.99*

Table 2: Empirical results comparing baselines and CAVE. We note that fine-tuned CAVE is competitive with strong baselines such as zero-shot and few-shot GPT-4-TURBO, while having higher rationale quality (via consistency).

rectness and train a LLAMA-3-8B with them. We report final train set sizes in Appendix E.

Few-shot. Finally, we present a few-shot baseline with PROMPT-CAVE on GPT-4-TURBO and LLAMA-3-8B-INSTRUCT. This baseline is strong, but undesirable owing to the high expenses of performing few-shot AV⁹. We have submitted the few shot prompts in supplementary material.

3.3 Empirical Results

Table 2 shows the empirical results. We report accuracy (ACC.) for all models, and consistency (CONS-R-L) for CAVE and PROMPT-CAVE. For all 3 datasets, CAVE beats all relevant baselines in terms of rationale-label consistency. For FANFICTION, CAVE beats all baselines (including few-shot GPT-4-TURBO by 0.9%) in accuracy; for IMDB62 and BLOG-AUTH, CAVE obtains competitive accuracies (2% and 4% less than the highest accuracy, which were PROMPTAV-SFT and CoT-SFT respectively).

Why show two zero/few-shot baselines? In addition to the strong (online) baseline of GPT-4-TURBO, we also show zero/few-shot results using a less-strong local model, LLAMA-3-8B-INSTRUCT. The significant performance degradation there (for all cases except few-shot BLOG-AUTH) shows that it is insufficient merely to take existing approaches and run them using a local instruction-tuned model; our proposed distillation strategies are essential to get a high performance

⁹Each document in our datasets has *on-average* 1000 tokens; performing even a 4-shot baseline is highly expensive.

with local, offline models that face the disadvantage of having a *much lower* parameter size.

PROMPT-CAVE versus PROMPTAV, CoT.

We investigate the drop in accuracy as we go from CoT/PROMPTAV to PROMPT-CAVE in the zero-shot baseline GPT-4-TURBO. We observe that a majority of this drop is caused by instruction following errors, such as not generating in a JSON format, or generating a ‘MAYBE’ as the final label instead of ‘YES’ or ‘NO’. GPT-4-TURBO might still does produce the right *score* in the prediction (which is the indicator for the predicted label used by PROMPTAV and CoT), but since our evaluation requires the JSON format and a clear YES/NO label at the end, GPT-4-TURBO’s performance decreases with PROMPT-CAVE. For IMDB62, BLOG-AUTH and FANFICTION respectively, the number of test datapoints (out of 1000) where GPT-4-TURBO doesn’t follow the output format with zero-shot PROMPT-CAVE are 47, 153 and 148 respectively; the number of datapoints where it produces a label of ‘MAYBE’ are 59, 45 and 103 respectively. However, as we move to *few-shot* PROMPT-CAVE or *fine-tuned* CAVE, we find that instruction-following errors are almost non-existent, and there is a huge improvement in accuracy as well as consistency! Appendix F further discusses practical advantages of CAVE over GPT-4-TURBO, as well as experiments with more local models.

3.4 In-house human study:

AV + explanations is a highly complex task; annotators require a considerable amount of training and practice to be skillful in their evaluation. Hence, for our evaluation, we opted to conduct a *small but high-quality in-house human study*, instead of a large-scale study with platforms such as Amazon MTurk.

We pick 50 document pairs per dataset (total of 150 pairs) and evaluate the rationales generated by their corresponding CAVE. We analyze the strengths & weaknesses of CAVE’s rationales, which help to understand the utility of these rationales to downstream applications (Joshi et al., 2023) and to learn about potential areas for improvement (van der Lee et al., 2019). We analyze the generated rationales (for each linguistic feature as in Table 1) via three properties as described below; we employ 3 distinct annotators per datapoint. We also require the annotators to comment on why a rationale is unsatisfactory, if they find it to be so with respect to any property; this helps to perform a more detailed analysis at the instance level.

1. Detail-Consistency: Are the details mentioned in the rationale consistent with the input documents, or are they hallucinated? (eg:- if the rationale mentions that both documents use parantheses, but the documents don’t, then the details are *hallucinated*)

2. Factual-Correctness: Are the rationales factually correct? (eg:- if the rationale says that the text was informal in tone, while it was actually formal, or if it misinterprets an acronym to be author-specific slang, it is incorrect)

3. Label-Consistency CAVE’s output format has a distinct rationale and intermediate label (YES/NO/MAYBE) for each linguistic feature. As defined in §2, we use these intermediate labels to check the overall consistency of the rationale with the final label. In this human evaluation however, we check whether each individual rationale is consistent with its corresponding label.

In Table 3, we show number of examples (out of 150) for which the three annotators agree that the rationale aligns with the property being analyzed. We also show aggregate statistics for each property across all linguistic features, i.e., $\frac{sum(P1)}{150*8}$ where 8 is the number of linguistic features being analyzed.

Aggregate Analysis. Annotators find that CAVE’s rationales align with properties P1, P2, P3 **77.8%**, **84.6%** and **90.4%** of the cases respectively. That is, CAVE hallucinates, presents incorrect details

Linguistic Features	P1	P2	P3
punctuation style	63	103	125
special characters & capitalization style	87	116	136
acronyms and abbreviations	110	124	135
writing style	136	137	145
expressions and idioms	130	137	138
tone and mood	140	138	145
sentence structure	126	133	141
any other relevant tabs	142	127	120
aggregate	934	1015	1085
aggregate as %	77.8%	84.6%	90.4%

Table 3: **Human Pilot:** This table shows the number of examples out of 150 where all 3 annotators agree that the rationale conforms to the property being analyzed. P1 - ‘Detail-Consistency’, P2 - ‘Factual-Correctness’, P3 - ‘Label-Consistency’. We report these numbers for all linguistic features used by CAVE.

and produces inconsistent intermediate labels in 22.2%, 15.4% and 9.6% of the cases respectively.

Spurious patterns. Annotators noticed that documents in the FANFICTION dataset all used double instead of single quotes for apostrophes (eg:- it’s instead of it’s), leading to CAVE reasoning that these double quotes were indicative of shared authorship. Similarly, some documents in the BLOG-AUTH dataset used the term ‘urlLink’ to denote the presence of hyperlinks (perhaps a formatting pattern used by the blog platform) - CAVE however, took the presence of this unique term as an indicator of shared authorship. Since annotators view many datapoints in their analysis, they are able to identify these spurious patterns; models such as CAVE and GPT-4-TURBO view these datapoints in isolation and hence do not dismiss these patterns as indicators of shared authorship.

Qualitative examples & Annotator comments. We report commonly observed rationale errors in Table 4. Hallucination of details was a common error, especially for punctuation style and special characters / capitalization; for example, CAVE mentioning the presence of punctuation marks such as ellipses, parantheses, etc., when those did not exist in the document. One hypothesis is that LLMs are bad at these linguistic features since they deal with individual characters and not broad pieces of text like the other features. Another common error was the CAVE’s usage of topic of texts as a reason to indicate or refute shared authorship of two documents; topic cannot be an absolute factor

Rationale / Label	Annotator’s Comment
punctuation style: Both texts use a variety of punctuation marks, including commas, parentheses, and quotation marks, but Text2 uses more complex structures such as dashes and colons.	Hallucination of punctuation marks when they don’t exist in the original documents.
any other relevant aspect: The subject matter of the two texts is vastly different, with Text1 focusing on technology and personal preference, while Text2 focuses on societal issues and personal identity. NO	Similarity/difference in topics or shared platforms does not definitively imply same/different authorship - the label should have been MAYBE for both these cases.
special characters style, capitalization style: Text1 shows a casual approach to capitalization, with instances of lower-case usage in contexts that typically require capitalization (e.g., ‘relly’ instead of ‘really’). Text2 adheres to standard capitalization rules.	Confusion between capitalization and misspelling/typo.
acronyms and abbreviations: Neither text makes significant use of acronyms or unusual abbreviations, maintaining a formal tone in dialogue without informal abbreviations.	Lack of acronyms does not indicate a formal tone - this rationale is even more incorrect when the original text is informal in tone.
acronyms and abbreviations: Neither text makes significant use of acronyms or abbreviations, which suggests a similarity in this aspect. However, the absence of such elements makes it difficult to draw a definitive conclusion. YES	The label should have been MAYBE owing to lack of information with this feature.
special characters style, capitalization style: Both texts use colloquial expressions (‘buttload of dishes’, ‘nagged the darn thing out of Shawn’). These expressions contribute to the informal tone of the texts.	Hallucination and factually incorrect: Both expressions were in Text-1, none in Text-2.

Table 4: Sample rationales from various examples, to show observed error types by CAVE

in deciding authorship (Wegmann et al., 2022a), and hence such labels must be ‘MAYBE’. We also discuss random errors owing to computational limitations (such as limited context window in LMs) in Appendix H.

We report the human study numbers for individual datasets in Table 19. We also release the [complete set of test-set rationales generated by CAVE](#).

3.5 Future steps

The first step towards handling these issues is to make downstream users aware about them, so that they can make an informed final decision with CAVE’s rationales; the independence between linguistic features make it easy for users to process them separately and use them as needed.

Ablation study with linguistic features. A straightforward solution is to simply identify the linguistic feature that is the most uncertain/problematic in the training data, and remove it from the datapoints: we perform this experiment as an ablation study. We identified from analyzing the silver data that ‘punctuation style’ has an intermediate label of ‘MAYBE’ in roughly 60-70% of the data for each dataset. We hypothesized that this feature was generally indecisive and unimportant; we removed it from the train data and fine-tuned a LLAMA-3-8B on the same - however, we found that the performance *reduced* upon removal of this

feature for all 3 datasets. This suggests that though this feature is generally indecisive, it does offer much needed information or context that the model depends upon; simply removing the feature will not improve the performance, and hence, we need more intricate methods such as weighing of linguistic features depending on the datapoint at hand.

Targeted training. As future steps, we propose to perform targeted training that directly addresses the problems of (1) *hallucinated details*, and (2) *trivial reasonings from unimportant features*. For the hallucination issue, given that we have the input documents, we can potentially have a reward model that checks the entailment between the sub-explanations and the documents and verifies whether a detail is real or a hallucination; we can further perform reward-based learning with algorithms such as PPO (Schulman et al., 2017), Quark (Lu et al., 2022). For the latter issue, we can have *dynamic* weighting of linguistic features for each example, which highlight only the features that contribute meaningful similarities or differences between the documents.

4 Discussion

4.1 Robustness of CS-1 and CS-2

CONS-R-L, the automatic metric used to measure rationale-label consistency, is defined based on two binary variables (§2): CS-1 which measures con-

sistency between the final score and the final label, and Cs-2 which does the same between the final label and the aggregated intermediate labels. These variables were designed with common sense logic (such as, a final score ≥ 0.5 must indicate a final label of YES). But how robust are they statistically? We analyze the silver train data from GPT-4-TURBO with relevant questions below. We use the term *accepted* for datapoints that satisfied the consistency condition being analyzed (CS-1 or CS-2), and *rejected* for those that were filtered out. We split the analysis into 5 categories: accepted datapoints with a final label of YES/NO, and rejected datapoints with a final label of YES/NO/MAYBE (when the predicted label is ‘MAYBE’, it is always filtered out).

Cs-1: Is 0.5 the best threshold? In Table 5 we present the average score obtained by accepted/rejected datapoints in the five categories. We note that for accepted data for all three datasets, GPT-4-TURBO assigns a final score of 0.86 or higher when predicting the final label as ‘YES’, and assigns a final score of 0.33 or lower when predicting the final label as ‘NO’ - these indicate that accepted datapoints are well away from the threshold boundary. For rejected data, GPT-4-TURBO on average assigns a score surrounding 0.6 for predicted labels of ‘NO’ and ‘MAYBE’ (there are no cases where the rejected data has a label of ‘YES’, i.e., where the label was ‘YES’ and the score < 0.5) - this shows that rejected datapoints were *not* rejected just because of a harsh thresholding¹⁰.

Dataset	Accepted		Rejected		
	YES	NO	YES	NO	MAYBE
IMDB62	0.89	0.33	-	0.63	0.65
BLOG-AUTH	0.86	0.18	-	0.62	0.58
FANFICTION	0.87	0.30	-	0.63	0.62

Table 5: Statistical analysis of Cs-1: Average ‘final score’ generated by GPT-4-TURBO + PROMPT-CAVE in the silver train data.

Cs-2: How much do the intermediate MAYBE’s contribute? We require that the intermediate labels as an aggregate (with ‘MAYBE’ as a buffer for both ‘YES’ and ‘NO’) side with the final label. But statistically, how much of this support comes from the MAYBE’s, and how much comes from YES and NO? In Table 6 we present

¹⁰harsh thresholding is when the final score is marginally crossing the boundary to be accepted, for example, a score of 0.51 for a label of ‘NO’ or 0.49 for a label of ‘YES’

Predicted Label →	Accepted		Rejected		
	YES	NO	YES	NO	MAYBE
Dataset ↓	% MAY. in MAY.+YES	% MAY. in MAY.+NO	%NO	%YES	%MAY.
IMDB62	15.7	38.8	-	51.2	45.7
BLOG-AUTH	22.1	24.6	-	50.0	54.3
FANFICTION	17.0	32.1	-	51.3	43.1

Table 6: Statistical analysis of Cs-2: % of intermediate labels

relevant proportions of intermediate labels for the five categories. For accepted data, we calculate the % of MAYBE in support of YES and NO respectively, that is, % of MAYBE in MAYBE + YES or MAYBE + NO. For rejected data, we calculate the % of NO’s in the intermediate labels for a predicted label of ‘YES’ (that is, how many NO’s were there in the intermediate labels to make them inconsistent with the final label?), % YES’s for a predicted label of ‘NO’, and % of MAYBE’s for a predicted label of ‘MAYBE’. We note that for accepted data, the support given by ‘MAYBE’ is on average $\sim 15\text{-}22\%$ for a final label of ‘YES’, and $\sim 24\text{-}39\%$ for a final label of ‘NO’ - this indicates that the main support for the final labels of YES/NO come respectively from intermediate labels of YES/NO. For rejected data where the predicted label was NO, the aggregate intermediate labels came roughly to $\sim 50\%$ YES, indicating that the inconsistency was not due to indecision of the model with MAYBE’s. For rejected data where the predicted label was MAYBE, the aggregate intermediate labels are $\sim 43\text{-}54\%$ MAYBE, indicating a higher level of uncertainty in the model.

Measuring uncertainty via score and labels: As per our logic, a final score close to 0.5, and/or a high proportion of ‘MAYBE’ intermediate labels indicate a higher uncertainty of the model in its final prediction. To validate this hypothesis, we calculate the correlation¹¹ between the fraction of MAYBE’s in the intermediate labels, and 1 minus the absolute distance of the final score from 0.5¹². Formally,

$$v_1 = \frac{\text{count}(\text{MAYBE})}{\text{num_intermediate_labels}}$$

$$v_2 = 1 - |\text{final score} - 0.5|$$

$$\text{correlation} = \text{PearsonCorr}(v_1, v_2)$$

We calculate this correlation for the GPT-4-

¹¹from `scipy.stats import pearsonr`

¹²a lesser distance to 0.5 implies high uncertainty, so it is subtracted from 1 to make the final value higher if there is higher uncertainty

TURBO train data; we obtain IMDB62: 0.66, BLOG-AUTH: 0.8, FANFICTION: 0.6, all with $p \ll 0.01$. The high correlation scores validate our hypothesis.

4.2 Analyzing signal strengths across linguistic features

In Table 7, we present a detailed analysis of the intermediate label (MAYBE) signal strengths across all linguistic features, as measured from the training data of each dataset. Punctuation style sees a high proportion of MAYBE’s (60-70%) for all datasets, whereas features such as writing style, tone & mood etc., have a very low proportion. On average across all linguistic features, the proportion of ‘MAYBE’ is 20-25%, reinforcing the conclusion from §4.1 that the majority support for the predicted label comes from decisive intermediate labels of ‘YES’ and ‘NO’, and not from the uncertain ‘MAYBE’.

Ling. Feature	IMDB62	BLOG-AUTH	FANFICTION
punctuation style	68.4 %	60.2 %	59.4 %
special char./ capitalization	21.2 %	23.9 %	20.7 %
acronyms & abbr.	26.6 %	44.5 %	29.1 %
writing style	6.5 %	5.9 %	7.1 %
expressions/idioms	24.6 %	14.4 %	23.2 %
tone & mood	5.1 %	7.0 %	7.2 %
sentence structure	28.0 %	9.8 %	18.8 %
any other	17.5 %	8.6 %	8.8 %
avg. across features	24.74 %	21.79 %	21.79 %

Table 7: Proportion of datapoints in the train set for each linguistic feature where the intermediate label is ‘MAYBE’. Punctuation style sees a high proportion of ‘MAYBE’ labels (60-70% across the datasets), whereas other features see much lower proportions.

4.3 CAVE versus GPT-4-TURBO

In Appendix D, we present analyses of the silver training data generated by GPT-4-TURBO as well as human evaluation results of the same. In addition, here we compare rationales generated by CAVE and GPT-4-TURBO on the test set in the following three aspects:

Does CAVE yield less MAYBE’s than GPT-4-TURBO? In terms of the final output label, since CAVE is fine-tuned on data that has only YES/NO labels as the final output, it always produces only YES/NO and never a MAYBE; this is not the case with GPT-4-TURBO since we zero/few-shot prompt it. In terms of the intermediate labels, we report the % of MAYBE’s averaged across linguistic

features for zero-shot GPT-4-TURBO and CAVE’s predictions on the test set in Table 8. The proportions are found to be similar, with zero-shot GPT-4-TURBO having slightly more MAYBE’s as its intermediate labels than CAVE.

Dataset	GPT-4-TURBO	CAVE
IMDB62	29.22 %	23.71 %
BLOG-AUTH	25.23 %	23.46 %
FANFICTION	26.64 %	20.34 %

Table 8: Proportion of intermediate labels MAYBE’s averaged across linguistic features for zero-shot GPT-4-TURBO and CAVE’s predictions on the test set.

Does CAVE yield less rejected samples with respect to CS-1/CS-2 than GPT-4-TURBO?

From Table 2, we see that CAVE has a higher CONS-R-L, which means that it yields less rejected samples than GPT-4-TURBO; in our individual CS-1 and CS-2 analysis, we find that CAVE has higher values for both. We report values below for CAVE and zero-shot GPT-4-TURBO on the test set below in Table 9.

Dataset	Model	CONS-R-L	CS-1	CS-2
IMDB62	GPT-4-T	0.86	0.86	0.86
	CAVE	0.99	0.99	1.0
BLOG-AUTH	GPT-4-T	0.73	0.73	0.73
	CAVE	1.0	1.0	1.0
FANFIC-TION	GPT-4-T	0.67	0.67	0.67
	CAVE	0.99	0.98	0.99

Table 9: CS-1 and CS-2 values for CAVE and zero-shot GPT-4-TURBO on the test sets.

5 Conclusion

Existing methods for authorship verification face a lack of *accessible explainability*. In this work, we train an in-house language model CAVE to generate structured and consistent explanations for AV. We present strong experimental results on three difficult AV datasets; we also present a human pilot study on the quality of generated rationales. Finally, we provide detailed discussions on the robustness of our consistency metric, and explore errors made by CAVE due to factors such as hallucination, dataset biases, etc. to guide research in this area. As future work, we would like to explore reward metrics and reward-based learning for AV explanations.

Limitations

Bias propagation from GPT-4-TURBO to CAVE, and mitigation: Since there is no human-written explanation data for AV, we used silver data generated by GPT-4-TURBO; however, this can lead to *biases*¹³ present in GPT-4-TURBO to unconsciously propagate to CAVE. The datasets we use in this paper are not highly bias-prone; however, researchers/users who use our framework on their dataset might deal with biases that stem from their dataset (such as gender bias, racial bias, etc., for example: ‘Document-1’s writing style is depictive of demographic-X whereas Document-2’s writing style is depictive of demographic-Y’). For users dealing with sensitive data for AV, we recommend they use bias mitigation strategies from relevant literature such as Counterfactual Role Reversal (Gupta et al., 2022) for gender bias, reducing bias in the upstream model before fine-tuning (UBM (Jin et al., 2021), strategic pretraining (Feng et al., 2023)), ensembling models with varying/opposing biases (partisan ensemble (Feng et al., 2023)), etc.

Hallucination. As we see in our human study, CAVE’s rationales sometimes have *hallucinated* details, especially for certain linguistic features. Hallucination is a common problem faced when using LMs (Ji et al., 2023; Huang et al., 2023); prior works (Zhang et al., 2023; Shi et al., 2024) have proposed task-specific strategies to alleviate the same. In our work, our proposed output format for CAVE allows for downstream users to analyze each linguistic feature independently, making it easier to *identify* hallucinations. We caution users to be aware of this issue, and to always manually read and verify the rationales before using them.

Completeness of rationale. We would like to emphasize here that our models have not been explicitly trained for *completeness*; that is, apart from any wrong information the generated explanation may have, it could also have missed some crucial similarity/dissimilarity between the documents. We discuss rationale properties such as completeness in Appendix K.

Trust. Lipton (2018) discusses that the definition of *trust* in a model is subjective: some users only trust well-understood models; other users might be inclined to trust well-performing mod-

els regardless of how (not) explainable they are. A recent work Sun et al. (2024) defines various *dimensions* for trust for language models, such as truthfulness, safety, robustness, etc. Our improvements in terms of accessibility of rationales and automatic quality evaluation contribute to improved trust in certain dimensions; as future work, we would like to evaluate and improve on other dimensions of trust, including higher accuracy, more datasets, more OOD evaluation, etc.

Practical limitations. Lastly our method has practical limitations faced by any NLP task; our models are restricted by the size of the context window allowed by the model (both input and output text) (also keeping in mind restrictions due to computational costs, financial costs and GPU size available). This can lead to data missed by the model in its reasoning.

Data. All the datasets that we use in our work are released publicly for usage and have been duly attributed to their original authors. Note that while we perform authorship verification, all the datasets we use have been *anonymized* by their respective creators. We do not have access to or use the names of individual people.

Reproducibility

For all our experiments, we report (1) the complete hyperparameter setting and any bounds explored (Appendix E) as well as the sizes and versions/pretrained-model links of all models used, (2) the time taken per experiment, and infrastructure used, (3) descriptions of datasets, and demonstrations used to sample rationales from GPT-4-TURBO. Further, we have released all our codes and our training/test data at github.com/INK-USC/Controllable-AV-Explanations.

LLAMA-3 Usage and License. Note that we were granted access to download and use LLAMA-3-8B from <https://huggingface.co/meta-llama/Meta-Llama-3-8B>. While we train a local LLAMA-3-8B, we do not release the trained model; we submit only the script and data used to train the model. Also note that we adhere to the intended use case for LLAMA-3-8B, i.e., “Llama 3 is intended for commercial and research use in English. Instruction tuned models are intended for assistant-like chat, whereas pretrained models can be adapted for a variety of natural language generation tasks.”

¹³we adopt the definition of bias from Jin et al. (2021); Blodgett et al. (2020); disparate model performance on different subsets of data which are associated with different demographic groups

Acknowledgements

This research is supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via the HIATUS Program contract #2022-22072200006. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein. We thank the human annotators for their help with the pilot study. Finally, we thank all our collaborators at USC INK Research Lab and ISI for their constructive feedback on this work.

References

- Emir Araujo-Pino, Helena Gómez-Adorno, and Gibran Fuentes Pineda. 2020. Siamese network applied to authorship verification. In *CLEF (Working Notes)*.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. [Language \(technology\) is power: A critical survey of “bias” in NLP](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. Association for Computational Linguistics.
- Benedikt Boenninghoff, Steffen Hessler, Dorothea Kolossa, and Robert M Nickel. 2019. Explainable authorship verification in social media via attention-based similarity learning. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 36–45. IEEE.
- John F. Burrows. 2002. [‘delta’: a measure of stylistic difference and a guide to likely authorship](#). *Lit. Linguistic Comput.*, 17:267–287.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. [e-snli: Natural language inference with natural language explanations](#). In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Aaron Chan, Maziar Sanjabi, Lambert Mathias, Liang Tan, Shaoliang Nie, Xiaochang Peng, Xiang Ren, and Hamed Firooz. 2022. Unirex: A unified learning framework for language model rationale extraction. In *International Conference on Machine Learning*, pages 2867–2889. PMLR.
- Aaron Chan, Zhiyuan Zeng, Wyatt Lake, Brihi Joshi, Hanjie Chen, and Xiang Ren. 2023. [Knife: Distilling reasoning knowledge from free-text rationales](#). *Preprint*, arXiv:2212.09721.
- Hanjie Chen, Faeze Brahman, Xiang Ren, Yangfeng Ji, Yejin Choi, and Swabha Swayamdipta. 2023a. [REV: Information-theoretic evaluation of free-text rationales](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2007–2030, Toronto, Canada. Association for Computational Linguistics.
- Wei-Lin Chen, An-Zi Yen, Hen-Hsen Huang, Cheng-Kuang Wu, and Hsin-Hsi Chen. 2023b. [Zara: Improving few-shot self-rationalization for small language models](#). *Preprint*, arXiv:2305.07355.
- Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20:273–297.
- Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023. [From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair NLP models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11737–11762, Toronto, Canada. Association for Computational Linguistics.
- John G. D. Grieve. 2007. [Quantitative authorship attribution: An evaluation of techniques](#). *Lit. Linguistic Comput.*, 22:251–270.
- Umang Gupta, Jwala Dhamala, Varun Kumar, Apurv Verma, Yada Pruksachatkun, Satyapriya Krishna, Rahul Gupta, Kai-Wei Chang, Greg Ver Steeg, and Aram Galstyan. 2022. [Mitigating gender bias in distilled language models via counterfactual role reversal](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 658–678, Dublin, Ireland. Association for Computational Linguistics.
- Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal. 2020. [Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language?](#) *arXiv preprint arXiv:2010.04119*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for nlp](#). *ArXiv*, abs/1902.00751.
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. [Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8003–8017, Toronto, Canada. Association for Computational Linguistics.
- J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *ArXiv*, abs/2106.09685.
- Yujia Hu, Zhiqiang Hu, Chun Wei Seah, and Roy Ka-Wei Lee. Instructav: Instruction fine-tuning large

- language models for authorship verification. In *First Conference on Language Modeling*.
- Zhengmian Hu, Tong Zheng, and Heng Huang. 2024. A bayesian approach to harnessing the power of llms in authorship attribution. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13216–13227.
- Baixiang Huang, Canyu Chen, and Kai Shu. 2024a. Authorship attribution in the era of llms: Problems, methodologies, and challenges. *arXiv preprint arXiv:2408.08946*.
- Baixiang Huang, Canyu Chen, and Kai Shu. 2024b. Can large language models identify authorship? *arXiv preprint arXiv:2403.08213*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*.
- Chia-Yu Hung, Zhiqiang Hu, Yujia Hu, and Roy Lee. 2023. Who wrote it and why? prompting large-language models for authorship verification. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14078–14084.
- Catherine Ikae. 2021. Unine at pan-clef 2021: Authorship verification. In *CLEF (Working Notes)*, pages 1995–2003.
- Sarthak Jain, Sarah Wiegrefe, Yuval Pinter, and Byron C Wallace. 2020. Learning to faithfully rationalize by construction. *arXiv preprint arXiv:2005.00115*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.
- Xisen Jin, Francesco Barbieri, Brendan Kennedy, Aida Mostafazadeh Davani, Leonardo Neves, and Xiang Ren. 2021. On transferability of bias mitigation effects in language model fine-tuning. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3770–3783, Online. Association for Computational Linguistics.
- Xisen Jin, Zhongyu Wei, Junyi Du, Xiangyang Xue, and Xiang Ren. 2019. Towards hierarchical importance attribution: Explaining compositional semantics for neural sequence models. *arXiv preprint arXiv:1911.06194*.
- Brihi Joshi, Ziyi Liu, Sahana Ramnath, Aaron Chan, Zhewei Tong, Shaoliang Nie, Qifan Wang, Yejin Choi, and Xiang Ren. 2023. Are machine rationales (not) useful to humans? measuring and improving human utility of free-text rationales. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7103–7128, Toronto, Canada. Association for Computational Linguistics.
- Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brahman, Chandra Bhagavatula, Ronan Le Bras, and Yejin Choi. 2022. Maieutic prompting: Logically consistent reasoning with recursive explanations. *arXiv preprint arXiv:2205.11822*.
- Akos Kádár, Grzegorz Chrupała, and Afra Alishahi. 2017. Representation of linguistic form and function in recurrent neural networks. *Computational Linguistics*, 43(4):761–780.
- Gregory R. Koch. 2015. Siamese neural networks for one-shot image recognition.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners. *Preprint*, arXiv:2205.11916.
- Moshe Koppel and Yaron Winter. 2014. Determining if two documents are written by the same author. *Journal of the Association for Information Science and Technology*, 65(1):178–187.
- Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2016. Rationalizing neural predictions. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 107–117, Austin, Texas. Association for Computational Linguistics.
- Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. Understanding neural networks through representation erasure. *arXiv preprint arXiv:1612.08220*.
- Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023a. Symbolic chain-of-thought distillation: Small models can also “think” step-by-step. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2665–2679, Toronto, Canada. Association for Computational Linguistics.
- Shiyang Li, Jianshu Chen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, Wenhui Chen, et al. 2024. Explanations from large language models make small reasoners better. In *2nd Workshop on Sustainable AI*.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023b. Making language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5315–5333.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. *Preprint*, arXiv:2305.20050.

- Zachary C Lipton. 2018. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57.
- Wei Liu, Jun Wang, Haozhao Wang, Ruixuan Li, Yang Qiu, YuanKai Zhang, Jie Han, and Yixiong Zou. 2023. Decoupled rationalization with asymmetric learning rates: A flexible lipschitz restraint. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23*, page 1535–1547, New York, NY, USA. Association for Computing Machinery.
- Ximing Lu, Sean Welleck, Jack Hessel, Liwei Jiang, Lianhui Qin, Peter West, Prithviraj Ammanabrolu, and Yejin Choi. 2022. Quark: Controllable text generation with reinforced unlearning. *Advances in neural information processing systems*, 35:27591–27609.
- Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In *Proceedings of the 31st international conference on neural information processing systems*, pages 4768–4777.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 305–329.
- Ana Marasovic, Iz Beltagy, Doug Downey, and Matthew Peters. 2022. Few-shot self-rationalization with natural language prompts. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 410–424, Seattle, United States. Association for Computational Linguistics.
- Meta-AI. 2024. Introducing meta llama 3: The most capable openly available llm to date. <https://ai.meta.com/blog/meta-llama-3/>. Accessed: (June 1, 2024).
- Erwan Moreau and Carl Vogel. 2013. Style-based distance features for author verification-notebook for pan at clef 2013. In *CLEF 2013 Evaluation Labs and Workshop-Working Notes Papers*, pages Online–proceedings.
- Frederick Mosteller and David L. Wallace. 1984. *Applied bayesian and classical inference : the case of the federalist papers*.
- Maryam Najafi and Ehsan Tavan. 2022. Text-to-text transformer in authorship verification via stylistic and semantical analysis. In *CLEF (Working Notes)*, pages 2607–2616.
- Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan. 2020. Wt5?! training text-to-text models to explain their predictions. *arXiv preprint arXiv:2004.14546*.
- OpenAI. 2024a. Gpt-4-turbo. <https://help.openai.com/en/articles/8555510-gpt-4-turbo-in-the-openai-api>. Accessed: (May 18, 2024).
- OpenAI. 2024b. Information about gpt-xx models. <https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo>. Accessed: (May 18, 2024).
- Ajay Patel, Nicholas Andrews, and Chris Callison-Burch. 2022. Low-resource authorship style transfer: Can non-famous authors be imitated?
- Ajay Patel, Delip Rao, Ansh Kothary, Kathleen McKeown, and Chris Callison-Burch. 2023. Learning interpretable style embeddings via prompting llms. *Preprint*, arXiv:2305.12696.
- Nina Poerner, Benjamin Roth, and Hinrich Schütze. 2018. Evaluating neural network explanation methods using hybrid documents and morphological agreement. *arXiv preprint arXiv:1801.06422*.
- Danish Pruthi, Bhuwan Dhingra, Livio Baldini Soares, Michael Collins, Zachary C Lipton, Graham Neubig, and William W Cohen. 2020. Evaluating explanations: How much do explanations from the teacher aid students? *arXiv preprint arXiv:2012.00893*.
- Sahana Ramnath, Brihi Joshi, Skyler Hallinan, Ximing Lu, Liunian Harold Li, Aaron Chan, Jack Hessel, Yejin Choi, and Xiang Ren. 2023. Tailoring self-rationalizers with multi-reward distillation. In *The Twelfth International Conference on Learning Representations*.
- Marco Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “why should I trust you?”: Explaining the predictions of any classifier. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*, pages 97–101, San Diego, California. Association for Computational Linguistics.
- Rafael A. Rivera-Soto, Olivia Elizabeth Miano, Juanita Ordonez, Barry Y. Chen, Aleem Khan, Marcus Bishop, and Nicholas Andrews. 2021. Learning universal authorship representations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 913–919, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Swarnadeep Saha, Peter Hase, and Mohit Bansal. 2023. Can language models teach weaker agents? teacher explanations improve students via theory of mind. *Preprint*, arXiv:2306.09299.
- Yunita Sari, Mark Stevenson, and Andreas Vlachos. 2018. Topic or style? exploring the most useful features for authorship attribution. In *International Conference on Computational Linguistics*.

- Jonathan Schler, Moshe Koppel, Shlomo Argamon, and James W Pennebaker. 2006. Effects of age and gender on blogging. In *AAAI spring symposium: Computational approaches to analyzing weblogs*, volume 6, pages 199–205.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. [Proximal policy optimization algorithms](#). *ArXiv*, abs/1707.06347.
- Yanir Seroussi, Ingrid Zukerman, and Fabian Bohnert. 2014. Authorship attribution with topic models. *Computational Linguistics*, 40(2):269–310.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih. 2024. Trusting your evidence: Hallucinate less with context-aware decoding. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 783–791.
- Xuelin Situ, Ingrid Zukerman, Cecile Paris, Sameen Maruf, and Gholamreza Haffari. 2021. Learning to explain: Generating stable explanations fast. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5340–5355.
- Joe Stacey, Yonatan Belinkov, and Marek Rei. 2022. Supervising model attention with human explanations for robust natural language inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11349–11357.
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pages 3319–3328. PMLR.
- Adaku Uchendu, Thai Le, Kai Shu, and Dongwon Lee. 2020. Authorship attribution for neural text generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8384–8395.
- Michiel van Dam. 2013. A basic character n-gram approach to authorship verification. In *Proceedings of CLEF*.
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, Sander Wubben, and Emiel Kraahmer. 2019. [Best practices for the human evaluation of automatically generated text](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 355–368, Tokyo, Japan. Association for Computational Linguistics.
- Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023. [SCOTT: Self-consistent chain-of-thought distillation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5546–5558, Toronto, Canada. Association for Computational Linguistics.
- Anna Wegmann, Marijn Schraagen, and Dong Nguyen. 2022a. Same author or just same topic? towards content-independent style representations. *ACL 2022*, page 249.
- Anna Wegmann, Marijn Schraagen, and Dong Nguyen. 2022b. [Same author or just same topic? towards content-independent style representations](#). In *Proceedings of the 7th Workshop on Representation Learning for NLP*, pages 249–268, Dublin, Ireland. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.
- Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022. [Symbolic knowledge distillation: from general language models to commonsense models](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4602–4625, Seattle, United States. Association for Computational Linguistics.
- Sarah Wiegrefe, Ana Marasović, and Noah A. Smith. 2021. [Measuring association between labels and free-text rationales](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10266–10284, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. *arXiv preprint arXiv:1908.04626*.
- Shuo Zhang, Liangming Pan, Junzhou Zhao, and William Yang Wang. 2023. Mitigating language model hallucination with interactive question-knowledge alignment. *arXiv preprint arXiv:2305.13669*, 3.

A PROMPT-CAVE

Table 10 shows the full PROMPT-CAVE.

B Related Work

B.1 Contemporaneous work

INSTRUCTAV (Hu et al.) is a contemporaneous work that also extends PROMPTAV to generate silver training data to fine-tune a smaller language model. INSTRUCTAV uses a *gold-label-aware*

PROMPT-CAVE:

Task: On a scale of 0 to 1, with 0 indicating low confidence and 1 indicating high confidence, please provide a general assessment of the likelihood that Text 1 and Text 2 were written by the same author. Your answer should reflect a moderate level of strictness in scoring. Here are some relevant variables to this problem.

1. punctuation style (e.g. hyphen, brackets, colon, comma, parenthesis, quotation mark)
2. special characters style, capitalization style (e.g. Continuous capitalization, capitalizing certain words)
3. acronyms and abbreviations(e.g. Usage of acronyms such as OMG, Abbreviations without punctuation marks such as Mr Rochester vs. Mr. Rochester, Unusual abbreviations such as def vs. definitely)
4. writing style
5. expressions and Idioms
6. tone and mood
7. sentence structure
8. any other relevant aspect

Provide the answer in a Python JSON format. Use the following keys for your dictionary: 'punctuation style', 'special characters style, capitalization style', 'acronyms and abbreviations', 'writing style', 'expressions and Idioms', 'tone and mood', 'sentence structure', 'any other relevant aspect', 'final score'. Apart from the 'final score', everything else must have a text value; also, the text should include a concluding YES/NO/MAYBE about whether the two texts are similar or not with respect to the key at hand. Finally, provide an 'output' key in your dictionary, which says YES if the two texts are written by the same author, and NO otherwise.

Text1: excellent until a completely unsatisfying ending This movie really got me hooked. The plot about...

Text2: A Nutshell Review: Oliver Twist "Please sir, may I have some more?" Ah, the immortal line from...

Table 10: PROMPT-CAVE: We extend [Hung et al. \(2023\)](#) to improve controllability in terms of the structure/parsing of the output and automatic verification of the consistency between the rationale and the label. The highlighted part is our extension to PROMPTAV.

zero-shot prompt to get silver data from ChatGPT, which are then filtered for quality (they perform consistency verifications based on identifying specific phrases in the generated rationale). In contrast, we prompt GPT-4-TURBO (without providing the gold label, in order to obtain unbiased rationales from GPT-4-TURBO) to generate explanations that are structured with respect to the various linguistic features, as well as intermediate labels (which allow for easy verification of consistency); we filter the silver data by checking of consistency between the different steps of the rationale and the final label. Our method CAVE differs from INSTRUCTAV in the specific format of the output used, as well as the metrics we subsequently use to verify/ensure rationale quality - **in addition to creating accurate and high-quality models, we also have the goal of making the process of authorship verification more transparent and accessible for end-users.**

B.2 Authorship Analysis

Classical AV systems use inter-textual distance to quantify the likelihood for two documents to share the same author ([Mosteller and Wallace, 1984](#); [Burrrows, 2002](#); [Grieve, 2007](#)). These verification techniques rely on expert-derived hand-crafted features to define the distance metric. Since such stylometric methods are grounded in hand-written rules, they are inherently explainable. However, these

methods are known to not scale well, especially in the scenario when for a given query document the correct target author needs to be mined from a large collection of candidate documents¹⁴.

Different from classical stylometric techniques, current AV systems typically employ neural architectures which perform better at scale ([Rivera-Soto et al., 2021](#); [Wegmann et al., 2022b](#)). Although significant amount of work has been devoted to advancing neural AV architectures, most of the current techniques lack interpretability. For example, a typical neural AV architecture comprises of a Siamese network ([Koch, 2015](#)), where each document is separately encoded into a vector and the model is trained to ensure that the document vectors are close if they are written by the same author and far otherwise ([Rivera-Soto et al., 2021](#); [Wegmann et al., 2022b](#)). In such a network, what each dimension of the document vector *means* is completely opaque.

To bridge this dichotomy between usability and interpretability between classical stylometric techniques and neural AV systems, ([Boenninghoff et al., 2019](#)) proposed to use attention heat-maps. This technique enabled them to determine which tokens in the input text are attended to (and not attended

¹⁴See Table 1 of ([Rivera-Soto et al., 2021](#)) for a comparison between scores from neural transformer/convolutional models and classical *tf-idf* technique.

to) by the neural model when performing AV. Different from attention-heat maps, model-agnostic methods, such as LIME (Ribeiro et al., 2016), have also been used to probe feature importance (Sari et al., 2018). Such explanations are helpful, however, they are inherently local¹⁵ and thus cannot capture higher level abstractions which can be informative of authorship, for example *fluency*, *tone and mood*, etc.

Motivated by the recent advancements made possible using LLMs, (Hung et al., 2023) and (Huang et al., 2024b) investigated whether LLMs can perform the task of authorship analysis. They found that, when prompted correctly, the LLMs can reliably determine whether or not two documents are written by the same author. LLMs have also been used to construct interpretable authorship style vectors (Patel et al., 2023) and subsequently used to rewrite text from a source author in the writing style of a target author (Patel et al., 2022). Similarly, LLMs have also been used by works (Hu et al., 2024) to perform authorship *attribution*, where, given a text, the task is to identify *who* wrote the text (i.e., which human or LLM from a predetermined set of authors (Huang et al., 2024a)).

Our work is the first attempt to distill free-text explanations for authorship verification from LLMs like GPT-4-TURBO into smaller LMs; we provide explicit focus to the joint explainability and security of AV systems via our model CAVE.

B.3 Self-rationalization

Explainability for neural network based models come in two formats: extractive rationales (includes pre-LLM era) and free-text rationales. Extractive rationales (Li et al., 2016; Sundararajan et al., 2017; Lundberg and Lee, 2017; Jin et al., 2019) involve analyzing the influence of input tokens on the predicted output via various methods such as gradient-based analysis of input tokens (Sundararajan et al., 2017; Lundberg and Lee, 2017), input perturbation (Poerner et al., 2018; Kádár et al., 2017), attention heatmap analysis (Pruthi et al., 2020; Stacey et al., 2022; Wiegrefe and Pinter, 2019), and trained models for this purpose (Lei et al., 2016; Chan et al., 2022; Jain et al., 2020; Situ et al., 2021; Liu et al., 2023). However, extractive rationales have limited applicability as discussed previously; hence we focus on free-text rationales.

¹⁵For example, heat-maps can only capture features at the token/span level.

The initial set of approaches for free-text rationales used gold standard human-written rationales to train rationale generation models (Camburu et al., 2018; Narang et al., 2020; Wiegrefe et al., 2021). Following the advent of emergent self-rationalization in LLMs (Wei et al., 2022), the research area moved to prompt large LMs with the help of curated templates with or without rationale demonstrations (Wei et al., 2022; Kojima et al., 2023; Li et al., 2023b; Jung et al., 2022; Lightman et al., 2023). Other approaches include leverage few-shot training approaches with a handful of gold rationales (Marasovic et al., 2022; Chen et al., 2023b), or generating silver standard training data from large LMs to distill small LMs to be better at the task or better rationalizers. (Ramnath et al., 2023; Li et al., 2023a; Chan et al., 2023; Wang et al., 2023; Saha et al., 2023; Hsieh et al., 2023).

C PROMPTAV and COT

Tables 11 and 12 show the prompts used to apply COT and PROMPTAV respectively. The prompts are taken from Hung et al. (2023); further note that the prompt is uniform across datasets.

PROMPTAV (and by extension, COT) which serve as the inspiration for CAVE produce outputs which cannot be automatically checked for consistency; these prompts produces a free-text explanation followed by a final score which can be thresholded to produce a YES/NO AV label. As observed in the experiments of (Hung et al., 2023), models tend to produce low confidence scores with PROMPTAV, which means that for better accuracy, the optimal threshold in their experiments with GPT-3.5-TURBO was found to be 0.2-0.3, which will most probably not be consistent with the free text rationale (that is, if the rationale says that it thinks it is not confident that the authors are the same, and the confidence is 0.25, the automatic label prediction will still assume that the predicted label is YES). Further, with just the score as the label, it is still possible that we misconstrue the label because the confidence score and the model’s predicted label in the text are different, as shown in the examples in Tables 13 and 14.

D GPT-4-TURBO silver data analysis

In Table 15 we show train dataset sizes at every stage of data generation and preprocessing.

Further, we also perform manual analysis of

CoT:

Task: On a scale of 0 to 1, with 0 indicating low confidence and 1 indicating high confidence, please provide a general assessment of the likelihood that Text 1 and Text 2 were written by the same author. Your answer should reflect a moderate level of strictness in scoring. **Let's think step by step.**

Text1: [T1]

Text2: [T2]

Table 11: CoT prompt for the task of AV, taken from [Hung et al. \(2023\)](#)

PROMPTAV:

Task: On a scale of 0 to 1, with 0 indicating low confidence and 1 indicating high confidence, please provide a general assessment of the likelihood that Text 1 and Text 2 were written by the same author. Your answer should reflect a moderate level of strictness in scoring. **Here are some relevant variables to this problem.**

1. punctuation style (e.g. hyphen, brackets, colon, comma, parenthesis, quotation mark)
2. special characters style, capitalization style (e.g. Continuous capitalization, capitalizing certain words)
3. acronyms and abbreviations(e.g. Usage of acronyms such as OMG, Abbreviations without punctuation marks such as Mr Rochester vs. Mr. Rochester, Unusual abbreviations such as def vs. definitely)
4. writing style
5. expressions and Idioms
6. tone and mood
7. sentence structure
8. any other relevant aspect

First step: Understand the problem, extracting relevant variables and devise a plan to solve the problem. Then, carry out the plan and solve the problem step by step. Finally, show the confidence score

Text1: [T1]

Text2: [T2]

Table 12: PROMPTAV prompt for the task of AV, taken from [Hung et al. \(2023\)](#)

GPT-4-TURBO’s silver rationales. For the most part, GPT-4-TURBO generates high-quality rationales; the most frequent issue observed is *hallucination* of details, especially punctuation marks. Out of the 60 samples analyzed (20 from each dataset), 10 samples showed hallucination in their rationale for punctuation style, 16 samples showed factual incorrectness across varying linguistic features, and 6 samples showed inconsistency of rationale with intermediate label in one linguistic feature (varying across samples) each.

E Hyperparameters

We download the pre-trained LLAMA-3-8B from HuggingFace (<https://huggingface.co/meta-llama/Meta-Llama-3-8B>) and train it for all our experiments. We train it with Low-Rank Adaption (LoRA; [Hu et al., 2021](#)), a parameter-efficient training method which offers lightweight training no additional inference latency, unlike other existing parameter-efficient methods ([Houlsby et al., 2019](#)). We train LLAMA-3-8B for 10 epochs. We use no validation set, and instead use the last checkpoint obtained from training; since we have no gold-standard explanations,

we opt against using the GPT-4-TURBO generated data as the validation set. We experiment with different values of ‘r’ and ‘alpha’ for LoRA (refer to https://huggingface.co/docs/peft/main/en/conceptual_guides/lora for their definitions). We use r=128 and alpha=256. Note that all our training experiments for LLAMA-3-8B take ~12 hours, and all our test set inferences take ~7-8 hours on a single NVIDIA Quadro RTX 8000 GPU. We report single-run results for all our experiments.

Why not LLAMA-3-8B-INSTRUCT for fine-tuning? As explained in the documentation for LLAMA-3-8B-INSTRUCT in huggingface <https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>, “Instruction tuned models are intended for assistant-like chat, whereas pretrained models can be adapted for a variety of natural language generation tasks”. Hence, we provide zero-shot and few-shot baselines with LLAMA-3-8B-INSTRUCT whereas we provide fine-tuned results on LLAMA-3-8B.

Document truncation. When we truncate documents for GPU memory considerations in the training of CAVE, we always ensure that equal amounts of both documents are present in the input. Further,

Sample CoT explanation that is inconsistent with label:

Based on the provided texts, I assess the likelihood that Text 1 and Text 2 were written by the same author as 0.7 (moderate confidence).

Here's my reasoning:

1. Both texts exhibit a personal and introspective tone, which suggests that the author is comfortable sharing their thoughts and feelings.
2. The language and sentence structure in both texts are similar, with a mix of short and long sentences, and a conversational tone.
3. Both texts touch on themes of self-reflection, personal growth, and the struggle to overcome negative thoughts and emotions.
4. The author's writing style is informal and conversational, with a focus on expressing their thoughts and feelings rather than presenting a polished or formal argument.

However, there are some differences that suggest the texts may not have been written by the same author:

1. The tone of Text 1 is more lighthearted and humorous, while Text 2 is more serious and introspective.
2. The language and vocabulary used in Text 2 are more formal and philosophical, with words like 'suffering,' 'unrealistic standards,' and 'setbacks.'
3. The structure and organization of the texts are different, with Text 1 consisting of short, disconnected paragraphs, and Text 2 featuring longer, more cohesive paragraphs.

Overall, while there are some similarities between the two texts, the differences in tone, language, and structure suggest that they may have been written by different authors.

Table 13: Sample output by CoT where the predicted label and the predicted confidence score are inconsistent with each other. The model predicts a high confidence score of 0.7 for same authorship (by usual standards of 0.5 as the threshold), but it predicts that the final answer is different authorship.

we always use the first chunk of the document as input and not a random chunk.

For our baseline SVM, we use TF-IDF feature vectors for BLOG-AUTH (words) and FANFICTION (characters), and word-count features for IMDB62 (characters). All 3 datasets' models use the RBF kernel. IMDB62 and FANFICTION use n-grams of minimum size 1 and maximum size 2, whereas BLOG-AUTH uses only n-grams of size 2. IMDB62 uses a C-value of 5.0. We used sklearn version 1.2.2 for this baseline.

For our baseline of LLAMA-3-8B finetuned with CoT and PROMPTAV, our final training dataset sizes are: for IMDB62, BLOG-AUTH and FANFICTION, the CoT sizes - 2.7k, 2.9k, 2.4k and the PROMPTAV sizes - 2.6k, 3k, 2.6k.

F CAVE with different offline models, Practical advantages of CAVE

We report experiments with local models of varying sizes in Table 20. After LLAMA-3-8B, MISTRAL-7B is the next best performing model.

Practical Advantages: CAVE also demonstrates practical advantages over GPT-4-TURBO such as (1) Energy: LLAMA-3-8B is just 8B parameters, whereas GPT-4-TURBO is orders of magnitude bigger - CAVE hence consumes much less energy

per datapoint, and hence is more affordable and efficient, (2) Financial cost: CAVE is also more financially affordable since it can be hosted on a local server; GPT-4-TURBO or any API based language model would cost much more since it charges for every input/output token, (3) Speed: On average, CAVE takes 3 minutes per datapoint for inference, whereas GPT-4-TURBO takes 4.5 minutes.

F.1 Human Pilot

Prior works in this field largely focus on task accuracy, with much lesser focus on explainability (Patel et al., 2023) and by extension, human study of the explanations. In this work, we take the *first step towards human analysis of AV explanations* by providing a pilot study on CAVE's rationales.

G Human Evaluation

Table 16 shows the instructions given to the annotators regarding the properties to be analyzed and Table 17 shows a sample human evaluation. We perform an in-house human evaluation, and annotators used were proficient in English and well aware of the high complexity of the task.

Table 19 shows the human pilot study numbers for individual datasets. The aggregate of these numbers were used for the overall analysis in Appendix

Sample PROMPTAV explanation that is inconsistent with label:

To assess the likelihood that Text 1 and Text 2 were written by the same author, we will analyze the texts based on the provided variables:

1. Punctuation Style:

- Text 1 uses a variety of punctuation, including colons, commas, and quotation marks, with a distinctive style of spacing before and after punctuation marks (e.g., 'this : THIS'). It also uses apostrophes in plural forms, which is unconventional (e.g., 'one's').

- Text 2 uses standard punctuation without the unconventional spacing or misuse of apostrophes seen in Text 1.

2. Special Characters and Capitalization Style: - Text 1 features continuous capitalization for emphasis (e.g., 'THIS IS BETTER THAN THE ORIGINAL') and capitalizes certain words for emphasis.

- Text 2 does not exhibit continuous capitalization or the same pattern of capitalizing certain words for emphasis.

3. Acronyms and Abbreviations:

- Text 1 does not prominently feature acronyms or abbreviations.

- Text 2 also does not prominently feature acronyms or abbreviations, making this variable less useful for differentiation.

4. Writing Style:

- Text 1 has a detailed, narrative style, sharing personal experiences and opinions in depth.

- Text 2 has a more straightforward, opinion-based style without the narrative depth or personal anecdotes found in Text 1.

5. Expressions and Idioms:

- Text 1 uses expressions and a more vivid language (e.g., 'take 15 showers afterwards').

- Text 2 uses simpler language and fewer idiomatic expressions.

6. Tone and Mood:

- Text 1 has a passionate, intense tone, especially when discussing the movie's impact and the author's feelings.

- Text 2 has a lighter, more casual tone, even when expressing disappointment.

7. Sentence Structure:

- Text 1 features complex sentences with multiple clauses and detailed descriptions.

- Text 2 uses simpler, more straightforward sentence structures.

8. Other Aspects:

- The subject matter and interest areas differ significantly between the two texts, with Text 1 focusing on horror films and Text 2 on a Disney film. This could indicate different interests or purposes for writing, though it's not definitive evidence of different authorship.

Given these observations, the differences in punctuation style, capitalization, writing style, tone, and sentence structure suggest that Text 1 and Text 2 were likely written by different authors. However, without more direct evidence (e.g., consistent unique spelling errors or very distinctive stylistic quirks that appear in both texts), absolute certainty is not possible.

Confidence Score: 0.8

This score reflects a high level of confidence that the texts were written by different authors, based on the analysis of the provided variables, but it leaves room for the possibility of the same author adapting their style to different subjects or formats."

Table 14: Sample output by PROMPTAV where the predicted label in the text and the predicted confidence score are inconsistent with each other. The model predicts a high confidence score of 0.8 (by usual standards of 0.5 as the threshold), but it predicts that the final answer is different authorship. Unless a human (or a trained model) can manually analyze the text, it is impossible to determine that the predicted label via the confidence score and the actual predicted label are different.

F.1.

H Random Errors due to Computational Limitations

We acknowledge random errors (i.e., not systematic) that could have occurred in our results owing to limitations of language models and computational resources. Our method uses a language model to read two input documents and provide an output; however, in practice, language models have a fixed maximum length of input due to model limitations as well as GPU memory limitations. This means that if the two input documents together (along with the instructions) do not fit into the in-

put, it becomes necessary to truncate them, which leads to loss of information of their latter parts. This issue is exacerbated when we move towards using in-house models on user servers (which typically have a shorter context window than API-based models like GPT-4-TURBO). In our experiments with LLAMA-3-8B, we used just 2 GPUs (to have a realistic user environment), which means that while training we had to restrict the input to be of maximum length 1280¹⁶, which could have affected the training for the longer documents. Fur-

¹⁶ note: during inference, we did not limit the size of the input, but this could pose a problem if the documents became much longer

Dataset	Initial size	Filter for Structure	Filter for float score	Filter for CONS-R-L
IMDB62	3400	3396	3255	2354
BLOG-AUTH	3600	3600	3146	2724
FANFICTION	3600	3383	2887	2159

Table 15: Train dataset sizes at every stage of data generation and preprocessing. The ‘Initial size’ column shows the number of datapoints generated from GPT-4-TURBO, and each succeeding column shows the number of datapoints remaining after the different filtering methods (in order of application).

Human Annotation Instructions:		
Given two documents, the task is to say whether they are written by the same author (label: YES) or not (label: NO). The model produces a set of explanations (corresponding to different aspects like punctuation, writing style, etc.) for the same each with an intermediate YES/NO/MAYBE label, and finally a YES/NO output label for the task. The human evaluation looks at the three properties described below, and gives a score of 1, 0.5, 0 or -1 for each of the individual explanations. The "Sample human eval" section gives an example evaluation. If you give a 0.5 or a 0, please write a short explanation what was wrong with the rationale!		
Criteria	Definition:	Options
consistency with details	Is the detail consistent with the documents or hallucinated? (for example, explanation mentions parentheses when it doesn't exist in input document)	1 - all details are consistent, 0.5 - some details are consistent, some are hallucinated, 0 - all details are hallucinated, -1 - I don't know
factual correctness	Are the details factually correct? (for example, it mentions serious writing style when writing style is actually humorous)	1 - all details are factual, 0.5 - some details are factual, 0 - no details are factual, -1 - I don't know
consistency with predicted label	Is the statement faithful with the YES/NO/MAYBE at the end?	1 - yes it is faithful, 0.5 - some details are faithful, 0 - no it is not faithful, -1 - I don't know

Table 16: Instructions given to annotators for the human pilot study. Sample human evaluations are in Table 17.

ther, even for generating the training data from GPT-4-TURBO for FANFICTION, the excessive size of the documents made the GPT-4-TURBO generations economically infeasible; we cut our the input to to a maximum of 300 words per document, which could have created errors/incomplete rationales that eventually propagated to our downstream training of LLAMA-3-8B.

I Out-of-domain inference

Finally, we evaluate CAVE models on the other datasets’ test sets (refer Table 21); we notice robust out-of-domain performance especially by the IMDB62 and FANFICTION models. This is a significant result, as it shows the generalizability of this approach to unseen AV datasets: although we use fine-tuning to improve performance over a zero-shot approach, it is not fine-tuning specifically to the domain that is most critical but rather fine-tuning to the structured format of CAVE’s output.

J Post-hoc Interpretability

We finally clarify that CAVE was designed to serve *post-hoc rationalization*, and not mechanistic (i.e., low-level mathematical) understanding of the language model’s workings. We quote from Lipton (2018) to emphasize that CAVE’s improved post-hoc explainability leads to improved utility to (1)

downstream users, and (2) researchers, who can leverage CAVE’s accessible outputs for detailed analyses, to address errors and develop improved models: “*while post-hoc explanations often do not elucidate precisely how a model works, they may nonetheless confer useful information for practitioners and end users of machine learning*”.

K More rationale properties

In addition to consistency, there are several properties necessary for high quality rationales (Joshi et al., 2023; Ramnath et al., 2023). Some of these are easily verifiable by humans, even if there exist no prevalent automatic methods to measure them: for eg:- factual correctness, hallucination in the rationale, etc. that we measured in our human evaluation (§F.1). There are also properties that are hard to verify even by human evaluators: for eg:- *completeness* of the rationale, i.e., has the rationale missed any significant (dis)similarity between the documents?. In fact, completeness does not have a comprehensive definition even for simpler tasks like multi-choice QA (Ramnath et al., 2023). In the future, we would like to work on automatic metrics and model training for these properties specifically for AV.

Sample evaluation:

Text1: One last update. I've decided to wear my khaki pants all day. I get off work at 4:30. I see no need to rush to the 5:05 train. So I'll take my time and go on the 5:35. I just blowdried my hair and I think it's going to be okay. I need to eat something. I've eaten my two favorite foods, of late, for breakfast: pop tarts and pears. I think it's going to be rice and peas for lunch.

Text2: Hey guys! I'd like to tell you all that RICE is organising an event, 'Where is the love'. Go to urlLink my site for more info, or go to the urlLink RICE site .It's going to be great! there will be around 800 people there, and this event will be about how the world's gone wrong with terrorism, war... etc. Date: 10 July 2004 Place: Willoughby Civic Centre Time: 7pm Ask Ben, Adro or me for more info. We'd like All of you to come! urlLink

GOLD LABEL: NO

Rationale	P1	P2	P3
punctuation style: Both texts use a variety of punctuation marks such as periods, commas, and ellipses. However, Text 1 uses more complex structures like colons and parentheses, which are less prevalent in Text 2. MAYBE	0.5 (ellipses / parantheses hallucinated)	0.5 (colons in both)	1
special characters style, capitalization style: Text 1 shows a more natural capitalization style without unusual emphasis on continuous capitalization or certain words. Text 2 capitalizes 'RICE' and uses all caps for 'URL', indicating a different approach. NO	0.5 (url not capitalized in Text2)	0.5 (RICE is implied as an organization, it is not a reflection of the capitalization approach)	1
acronyms and abbreviations: Text 1 does not prominently feature acronyms or abbreviations, while Text 2 includes 'URL' as an abbreviation. This suggests a difference in the use of abbreviations. Conclusion: NO	1	1	1
writing style: Text 1 has a more personal, reflective style, discussing personal routines and thoughts. Text 2 adopts a more informative and promotional tone, inviting people to an event. The styles are quite distinct. NO	1	1	1
expressions and idioms: Neither text heavily relies on idiomatic expressions, but the context in which they write differs significantly, with Text 1 focusing on personal narrative and Text 2 on event promotion. NO	1	1	1
tone and mood: The tone of Text 1 is casual and somewhat contemplative, while Text 2 is more enthusiastic and outward-facing, aiming to promote an event. These differences suggest different authors. NO	1	1	1
sentence structure: Text 1 uses a mix of short and long sentences, with a more conversational flow. Text 2 also uses a mix but is more structured, likely due to its promotional nature. The sentence structures show some similarities but also differences in context and purpose. MAYBE	1	1	1
any other relevant aspect: The subject matter of the two texts is vastly different, with Text 1 focusing on personal habits and reflections, and Text 2 on organizing an event and expressing concerns about global issues. This difference in content and focus further supports the likelihood of different authors. NO	1	0.5 (difference in content doesn't mean different authors)	0 (should be MAYBE. Difference in content is not difference in authorship)

Table 17: Sample human evaluation shown to human annotators

Linguistic Features	IMDB62			BLOG-AUTH			FANFICTION		
	P1	P2	P3	P1	P2	P3	P1	P2	P3
punctuation style	6	13	14	12	14	19	12	15	16
special characters/ capitalization style	11	17	17	9	15	18	12	16	15
acronyms/abbr.	17	17	18	12	17	17	18	18	13
writing style	18	18	19	19	19	20	17	16	20
expressions/idioms	18	19	19	17	15	18	18	18	17
tone and mood	18	19	19	18	18	20	19	16	20
sentence structure	17	19	16	17	18	20	14	16	19
any other	19	19	18	19	18	14	19	18	20

Table 18: **Human Pilot for individual datasets:** This table shows the number of examples (out of 20) where all 3 annotators agree that the rationale (corresponding to the linguistic feature) conforms to the property being analyzed. P1 is ‘Detail-Consistency’, P2 is ‘Factual-Correctness’ and P3 is ‘Label-Consistency’.

Linguistic Features	IMDB62			BLOG-AUTH			FANFICTION		
	P1	P2	P3	P1	P2	P3	P1	P2	P3
punctuation style	15	31	43	25	38	45	23	34	37
special characters/ capitalization style	31	40	47	23	37	46	33	39	43
acronyms/abbr.	38	41	47	25	37	45	47	46	43
writing style	47	45	46	44	48	49	45	44	50
expressions/idioms	44	48	49	40	43	45	46	46	44
tone and mood	47	48	48	46	47	48	47	43	49
sentence structure	45	48	46	42	45	48	39	40	47
any other	48	43	38	46	38	37	48	46	45

Table 19: **NEW Human Pilot for individual datasets:** This table shows the number of examples (out of 20) where all 3 annotators agree that the rationale (corresponding to the linguistic feature) conforms to the property being analyzed. P1 is ‘Detail-Consistency’, P2 is ‘Factual-Correctness’ and P3 is ‘Label-Consistency’.

L Analyzing cases where PROMPTAV succeeds but CAVE fails

We consider the 88 samples of the BLOG-AUTH test set where PROMPTAV succeeds but CAVE fails. From qualitative analysis of a few samples, we find out the following (non-exhaustive) patterns:

- Cases where the gold label is ‘YES’, and both methods fail: CAVE wrongly predicts a label of ‘NO’ with a final score of < 0.5 . PROMPTAV gives the wrong reasoning, but concludes by predicting a confidence score > 0.5 that the authors are *different* which leads to a misinterpretation that PROMPTAV was right. We provide an example of this error by PROMPTAV in Appendix C.

- Differences in perceiving patterns such as idioms, acronyms: Both methods generate similar explanations for abstract features such as writing style, tone & mood. But when it comes to token-level features such as acronyms, idioms, etc., the methods sometimes predict different subsets, which leads to differing final predictions. This comes under the issue of *completeness* of rationale, which we discuss in the Limitations section and in Appendix K. In this particular situation, both models have provided different incomplete rationales, which eventually leads to different labels.

Model	Test →	IMDB62		BLOG-AUTH		FANFICTION	
	Train ↓	ACC.	CONS.	ACC.	CONS.	ACC.	CONS.
LLAMA-3-8B	Fine-tune	74.1	0.99*	60.6	1.00*	62.6*	0.99*
MISTRAL-7B	Fine-tune	73.7	0.99	59.2	0.99	57.7	0.98
Stable-LM-Zephyr-3B	Fine-tune	50.8	0.71	47.7	0.81	42.0	0.68

Table 20: Empirical Results comparing CAVE trained with local models of different sizes. MISTRAL-7B is the next best performing model after LLAMA-3-8B. STABLE-LM-ZEPHYR-3B however has a lower performance.

Test →	IMDB62		BLOG-AUTH		FANFICTION	
Train ↓	ACC.	CONS.	ACC.	CONS.	ACC.	CONS.
IMDB62	—	—	59.9	1.0	61.2	0.99
BLOG-AUTH	71.3	0.96	—	—	45.3	0.77
FANFICTION	68.8	1.0	59.3	1.0	—	—

Table 21: Out of domain performance by CAVE