# Threads of Subtlety: Detecting Machine-Generated Texts Through Discourse Motifs

Zae Myung Kim<sup>1</sup> and Kwang Hee Lee<sup>2</sup> and Preston Zhu<sup>1</sup> and Vipul Raheja<sup>3</sup> and Dongyeop Kang<sup>1</sup>
University of Minnesota Twin Cities<sup>1</sup>, Kumoh National Institute of Technology<sup>2</sup>, Grammarly<sup>3</sup>
{kim01756, zhu00604, dongyeop}@umn.edu, kwanghee@kumoh.ac.kr, vipul.raheja@grammarly.com

#### **Abstract**

With the advent of large language models (LLM), the line between human-crafted and machine-generated texts has become increasingly blurred. This paper delves into the inquiry of identifying discernible and unique linguistic properties in texts that were written by humans, particularly uncovering the underlying discourse structures of texts beyond their surface structures. Introducing a novel methodology, we leverage hierarchical parse trees and recursive hypergraphs to unveil distinctive discourse patterns in texts produced by both LLMs and humans. Empirical findings demonstrate that, although both LLMs and humans generate distinct discourse patterns influenced by specific domains, human-written texts exhibit more structural variability, reflecting the nuanced nature of human writing in different domains. Notably, incorporating hierarchical discourse features enhances binary classifiers' overall performance in distinguishing between human-written and machine-generated texts, even on out-of-distribution and paraphrased samples. This underscores the significance of incorporating hierarchical discourse features in the analysis of text patterns. The code and dataset are available at https://github.com/ minnesotanlp/threads-of-subtlety.

#### 1 Introduction

The emergence of powerful instruction-tuned large language models (LLMs) (Ouyang et al., 2022; Muennighoff et al., 2023; Kopf et al., 2023) has led to an explosion of machine-generated texts in both offline and online domains. Consequently, discerning the authorship of texts has become a significant challenge, spanning from educational settings to the landscape of online advertising (Extance, 2023; Dalalah and Dalalah, 2023; Gołąb-Andrzejak, 2023). Indeed, many efforts have been made to tackle this issue by constructing corpora of machine-generated and human-authored texts (Dou

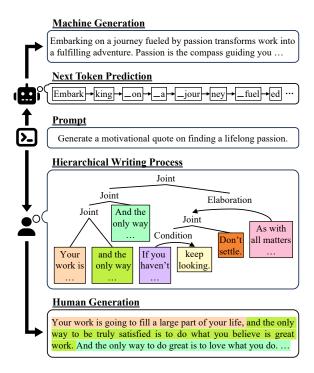


Figure 1: Human writers often employ hierarchical linguistic structures in writing whereas LLMs primarily operate by the sequential next token prediction task.

et al., 2022; Guo et al., 2023; Li et al., 2023) and developing models and benchmarks to tell them apart (Wu et al., 2023; Verma et al., 2023; Su et al., 2023; Chakraborty et al., 2023). The consensus seems to be that while classifiers that make use of the presence of LLM-specific signatures in the generated texts perform relatively well on in-domain texts, their accuracy drops significantly with out-of-domain samples. Furthermore, these detectors can be fooled easily with "paraphrasing attacks" even with in-domain samples (Sadasivan et al., 2023; Krishna et al., 2023).

This raises interesting questions on the underlying nature of human-written texts: "Are there any discernible, unique properties within texts crafted by humans?" and if so, "Might these distinctive signatures manifest at levels beyond surface struc-

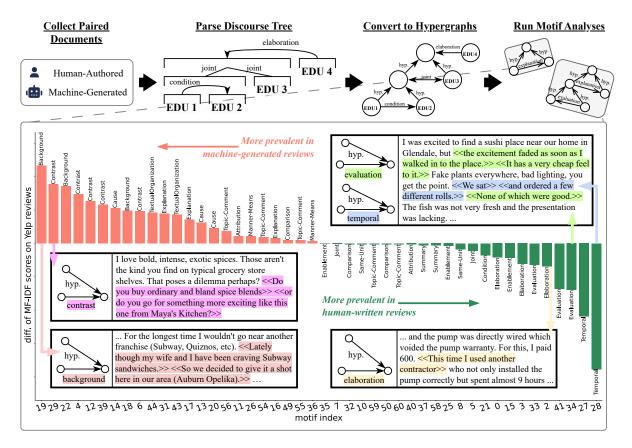


Figure 2: Difference in motif distribution of machine-generated and human-written texts for Yelp domain. Below are the top discourse motifs for each authorship and their corresponding discourse relations with examples. Text enclosed in angle brackets denotes the EDUs involved in the relation.

ture?" Undeniably, how we write varies greatly depending on the domain and intent addressed by the texts. Consider, for instance, the distinction between academic writing and a casual response to a Reddit post. The former pays much more attention to the logical progression of arguments and overall structural coherence, while the latter is characterized by a more spontaneous and less structured thought process.

On the contrary, LLMs, in consideration, are autoregressive models that generate the next token based on the previous sequence, i.e.  $P(x_t|x_1...x_{t-1})$ , without explicitly modeling the hierarchical structures in the process. Figure 1 illustrates this structural difference in the writing processes of LLMs and humans. We note that LLMs may also internally capture these hierarchical structures to some degree as a by-product, especially in their self-attention matrices (Xiao et al., 2021; Huber and Carenini, 2022). However, our interest lies in exploring whether there exist distinctive hierarchical structures that can aid in distinguishing

their authorship.<sup>1</sup>

To answer these inquiries systematically, we draw inspiration from discourse analysis in linguistics (Mann and Thompson, 1987), which uncovers structures found within a sentence, between sentences, and among the paragraphs of a document. Specifically, Figure 2 highlights our main approach: Given texts written by humans and LLMs, we construct a hierarchical parse tree for each document. Subsequently, we transform these trees into recursive hypergraphs, allowing us to perform network motif analysis of discourse relationships. The essence of our analysis lies in computing the difference in the distribution of these "discourse motifs" between texts generated by machines and those crafted by humans, identifying discernible discourse patterns within each authorship.

We summarize our contributions as follows:

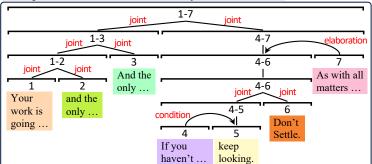
 We demonstrate that hierarchical discourse structures, as defined by RST, can be effectively modeled using recursive hypergraphs.

<sup>&</sup>lt;sup>1</sup>We provide a more comprehensive comparison with prior studies in Appendix A.

#### **Quote Example**

Your work is going to fill a large part of your life, and the only way to be truly satisfied is to do what you believe is great work. And the only way to do great work is to love what you do. If you haven't found it yet, keep looking. Don't settle. As with all matters of the heart, you'll know when you find it. — Steve Jobs

#### **Example of Discourse Tree Defined by RST Framework**



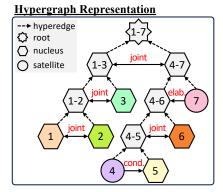


Figure 3: A quote from Steve Jobs and its RST tree converted into a hypergraph form. A hexagonal node represents the "nucleus" node, while a circular one denotes the "satellite" node. Each node is labeled with a span of EDU indices that it covers. The star-shaped node is the root node of the graph, encompassing all subgraphs and EDUs.

- To the best of our knowledge, we are the first to demonstrate that integrating hierarchical discourse-level features into the authorship detection task results in enhanced and more robust performance.
- Our empirical findings suggest that the robustness against paraphrasing attacks stems from the preservation of higher-level discourse structures, despite significant variations at the sentence level.
- Addressing a critical gap in existing research, where machine-generated texts typically fall within the range of 200 to 500 tokens, we construct a dataset specific to the creative writing domain. This dataset features 20 stories with lengths of up to 8K tokens, providing a more comprehensive perspective on longerform machine-generated content.

## **2** Modeling Hierarchical Structures

To explore unique patterns that may extend beyond surface-level structure, it is imperative to employ an expressive framework capable of representing hierarchical structures within texts. Human writers strategically utilize textual structures to systematically convey meaning, thereby augmenting the clarity, coherence, and persuasiveness of their written compositions. In the field of linguistics, the study of document structure is within the domain of discourse frameworks (Mann and Thompson, 1987; Miltsakaki et al., 2004; Lascarides and Asher, 2007). This paper adapts Rhetorical Structure Theory (RST) (Mann and Thompson, 1987) due to

its widespread acceptance and the availability of openly accessible pre-trained parsers. In this paper, we make use of the DMRST parser released by Liu et al. (2021).

#### 2.1 Rhetorical Structure Theory

RST assumes that any well-written document can be parsed into a (recursive) discourse tree. In this tree, each leaf node, called an Elementary Discourse Unit (EDU), corresponds to a phrase within a sentence, while higher-level nodes integrate these lower-level representations into more complex structures, such as phrases to sentences to paragraphs. It does so by assigning a discourse label (e.g., elaboration, contrast, cause, effect, etc.) to a relation (i.e., edge) between nodes. The linked nodes can be either "nuclei" (main ideas) or "satellites" (supporting details) depending on their relative importance in the relation. Figure 3 illustrates how a quote from Steve Jobs is segmented into seven EDUs, each marked in a distinct color, and parsed to form a discourse tree structure.

## 2.2 Conversion to Hypergraphs

Within the RST framework, discourse relations emerge not only through the linking of EDUs but also by encapsulating other relations. Specifically, RST relations are applied to a text "recursively" until every unit within the text becomes a constituent of an RST relation. This recursive application implies that a single EDU can form an RST relation with another individual EDU, and also become part of a relation involving a group of EDUs. For example, in Figure 3, EDU 4 and EDU 5 are linked

by a "condition" relation, while EDU 6 forms a "joint" relation with the collective unit of EDU 4 and EDU 5. Additionally, since the "joint" relation encompasses the "condition" relation, RST trees inherently demonstrate recursive encapsulations.

This recursive encapsulation characteristic of RST trees naturally corresponds with recursive hypergraphs, also referred to as "Ubergraphs" in mathematics (Joslyn and Nowak, 2017), where a discourse relation aligns with a *recursive hyperedge*. In short, hypergraphs generalize traditional graphs by allowing edges (or "hyperedges") to connect more than two nodes. The recursive hypergraphs further generalize the standard hypergraphs by allowing hyperedges to contain not only nodes but also other hyperedges. This leads to a hierarchical structure, where hyperedges at different levels of recursion represent relationships at varying levels of abstraction.

To facilitate ease of analysis, the hypergraph can be transformed into an isomorphic Levi graph (i.e., a standard traditional graph form) by introducing a dummy edge (depicted as a dotted line in Figure 3) that signifies the expanded hyperedge relation. By expanding the nested recursions, we are now able to run standard network analysis algorithms as well as train Graph Attention Networks (GATs) (Veličković et al., 2018), leveraging the inherent discourse structure and semantics of documents.

We note that although Yadati (2020) briefly explored the transformation of text into recursive hypergraphs, their methodology defined hyperedge relations as the relationships between entities (such as persons, places, etc.) occurring within the sentence-level organization of the text.

## 3 Proposed Method: Revealing Key Patterns via *Discourse Motifs*

This paper focuses on finding distinctive structural patterns, namely, hierarchical discourse relations based on the RST framework (§2). Our main investigative approach involves analyzing the distribution of discourse motifs, which are *recurring and statistically significant subgraph patterns* within a larger network. These motifs are often considered as the building blocks of larger networks and can provide insights into the organization and dynamics of the system (Milo et al., 2002; Takes et al., 2018). As illustrated in Figure 2, our proposed approach involves gathering documents written by humans and LLMs, parsing their RST trees, transforming

them into hypergraphs, and running motif analyses to understand how individuals (both humans and machines) craft texts. The following subsections illustrate how the discourse motifs are systematically formed (§3.1) and scored using our proposed metrics (§3.2).

#### 3.1 Discourse Motifs as Unions of Triads

One of the standard ways to construct motifs is to generate all possible non-isomorphic subgraphs of a fixed number of vertices. However, as the number of vertices increases, the number of possible patterns increases super-exponentially, rendering their counting extremely inefficient. Complicating matters further, our RST graphs are *directed* with their edges *labeled* with discourse relations, which further increases the complexity.<sup>2</sup>

With this in mind, we start by generating all possible motifs with 3 vertices, i.e., "triads." This is a minimal form of motifs other than the single edge case, "dyad." We select the triads to be the basis motifs as we find that when converting a recursive RST tree into a standard graph, it always forms *the union of subgraphs with 3 vertices*. We highlight the proof in Appendix B.

This process results in 69 non-isomorphic motifs. To create larger motifs, we pick two triads from all possible pairs, join them at all possible nodes, and only keep the non-isomorphic ones, resulting in unions of two triads. We term these configurations as "double-triads." We conduct this process one more round with pairs of a (single-)triad and a double-triad to produce a "triple-triad." We note that depending on how these multiples of triads are joined together, we can form discourse motifs with varying numbers of vertices such as 4, 5, and 7. Table 1 shows the number of identified motifs for each type. Graphical examples of motifs can be found in Appendix K.

Single-Triads	Double-Triads	Triple-Triads	
69	592	2,394	

Table 1: Number of non-isomorphic discourse motifs.

<sup>&</sup>lt;sup>2</sup>It is worth noting that some discourse relations can be bidirectional as well.

<sup>&</sup>lt;sup>3</sup>As we are joining motifs at their nodes, the total number of nodes for double-triads and triple-triads will be smaller than 6 and 9, respectively.

#### 3.2 Metrics for Analyses

Once we gather potentially useful discourse motifs, we simply count the number of isomorphic subgraphs for each motif across all documents in the datasets. However, as the graph isomorphism test is (possibly) NP-intermediate, the process requires heavy engineering feats to maximize computational efficiency, such as hashing the subgraphs (Shervashidze et al., 2011) and multiprocessing the datasets.<sup>4</sup> Afterward, we can identify distinctive motifs by examining their difference distributions based on their normalized counts (i.e., frequencies) (§3.2.1) as well as calculating our proposed metric termed "motif frequency-inverse document frequency (MF-IDF)" measure (§3.2.2). The rationale behind devising the MF-IDF scoring lies in pinpointing and prioritizing significant motifs, as counting the entire motif distribution is inefficient. An analogy can be drawn to crafting a localized tokenization scheme for LLM pre-training, wherein tokens represent motifs, and texts serve as hypergraphs in our context.

#### 3.2.1 Motif Difference Distributions

As a simple measure, we consider the difference distribution between features of motifs present in graphs. Specifically, we utilize their *motif frequencies*,  $\mathcal{F}(\cdot)$  and *weighted average depths*,  $\mathcal{D}(\cdot)$  for (document) graphs generated by machines  $(\mathbb{D}_{\text{machine}})$  or authored by humans  $(\mathbb{D}_{\text{human}})$ :

$$\mathcal{F}(\mathbb{D}_{\text{diff}}) = \mathcal{F}(\mathbb{D}_{\text{machine}}) - \mathcal{F}(\mathbb{D}_{\text{human}})$$
 (1)

$$\mathcal{D}(\mathbb{D}_{\text{diff}}) = \mathcal{D}(\mathbb{D}_{\text{machine}}) - \mathcal{D}(\mathbb{D}_{\text{human}})$$
 (2)

These distributions represent the discrepancies of motifs between the two classes of authorship and can be computed across different domains of texts as well.

Motif Frequency (MF). For each motif m in a set of identified motifs M, we compute motif frequency by counting its occurrences in a graph g and normalize it by the total number of motifs in the graph. The frequency is then averaged over the corresponding dataset  $\mathbb{D}$ :

$$\mathcal{F}(m, \mathbb{D}) = \frac{1}{|\mathbb{D}|} \sum_{g \in \mathbb{D}} \frac{\operatorname{count}(m, g)}{\sum_{m' \in M} \operatorname{count}(m', g)} \quad (3)$$

Weighted Average Depth (WAD). To better capture the hierarchical nature of discourse graphs, we also calculate the average depths of each motif m appearing in a graph g. Let S(m,g) be a collection of subgraphs of a graph g that are isomorphic to a motif m. Then, the depth of a motif m in a graph g is measured by the *mean position* 

$$\overline{d}(m,g) = \sum_{g' \in S(m,g)} \sum_{v \in V_{g'}} \frac{d(v, v_{root})}{|S(m,g)||V_{g'}|}$$
(4)

where  $V_{g'}$  denotes the set of nodes of a subgraph g' and  $d(v, v_{root})$  denotes the distance from a node v to the root  $v_{root}$ , calculated in g. It is then weighted by the count of the corresponding motif and normalized by the total number of counts:

$$\mathcal{D}(m, \mathbb{D}) = \frac{1}{|\mathbb{D}|} \sum_{g \in \mathbb{D}} \frac{\operatorname{count}(m, g) \cdot \overline{d}(m, g)}{\sum_{m' \in M} \operatorname{count}(m', g)}$$
(5)

This effectively measures the average position of each discourse motif appearing in the graphs.

## **3.2.2** Motif Frequency-Inverse Document Frequency (MF-IDF)

Inspired by the term frequency-inverse document frequency (TF-IDF) measure from information retrieval domain (Sparck Jones, 1972), we propose "motif frequency-inverse document frequency" (MF-IDF), noting that a motif (m) can be considered as a vocabulary of documents that are graphs  $(q \in \mathbb{G})$  in our case:

$$MF-IDF(m, g, \mathbb{G}) = \mathcal{F}(m, \{g\}) \cdot IDF(g, \mathbb{G})$$

where:

$$\mathcal{F}(m, \{g\}) = \frac{\operatorname{count}(m, g)}{\sum_{m' \in M} \operatorname{count}(m', g)}$$
$$\operatorname{IDF}(g, \mathbb{G}) = \log \frac{1 + |\mathbb{G}|}{1 + |g \in \mathbb{G} : m \in g)|}$$

## 4 Experimental Setup

Following the generation of three sets of discourse motifs for single-, double-, and triple-triads (§3.1), we evaluate them using the MF-IDF metric and retain those with scores surpassing at least one standard deviation. This yields a total of 207 motifs, with a detailed breakdown presented in Table 2.

Section 4.1 outlines the utilization of our proposed features (§3) within both the baseline models and analyses; and Section 4.2 describes the datasets designated and proposed for experimentation.

<sup>&</sup>lt;sup>4</sup>The code and the processed datasets for the experiments can be found at https://github.com/minnesotanlp/threads-of-subtlety.

Single-Triads	Double-Triads	Triple-Triads	
31	96	80	

Table 2: Number of selected discourse motifs with MF-IDF scores exceeding the threshold of at least one standard deviation.

#### 4.1 Baseline Models

To evaluate the usefulness of discourse motifs, we integrate them into three baseline models for authorship detection (§5.1): Random Forest (RF) (Breiman, 2001), Graph Attention Network (GAT) (Veličković et al., 2018), and Longformer (LF) (Beltagy et al., 2020). Figure 4 depicts an overall process of various types of input features being handled by the different models with varying numbers of model parameters and granularity of inputs.

RF is chosen to evaluate the baseline performance solely based on the discourse motif features (defined in §3.2.1) without considering the semantic content of the texts.

GAT receives the EDU-level representation of the texts as initial node embeddings, along with the discourse edge labels as edge embeddings. The texts for both EDUs and discourse labels are embedded by using Sentence-BERT (Reimers and Gurevych, 2019) model, ALL-MINILM-L6-v2. The discourse motif features are then appended to the mean-pooled representation of the graph.

LF is initialized with pretrained weights<sup>5</sup> and fed with the entire sequence of texts, in addition to the discourse motif features concatenated to its [CLS] token. Specifically, we align our use of LF with the methodology described by Li et al. (2023), which details a state-of-the-art detection model. However, we introduce three key distinctions: (i) we train a single classifier across all ten domains of the DEEPFAKETEXTDETECT dataset, (ii) we utilize approximately 56% of the original dataset to mitigate data imbalance issues, and (iii) we avoid adjusting the decision boundary across ten domainspecific classifiers. Our goal is to develop a general classifier capable of evaluating discourse motifs across various model architectures, rather than tailoring optimizations for individual domains.

We note that our approach can be seamlessly integrated into other baselines by concatenating discourse motif features with input representations. Additionally, it can enhance other detection strate-

gies by providing extra statistical features for zeroshot detectors or by aiding in the selection of more challenging and deceptive samples for adversarial learning setups.

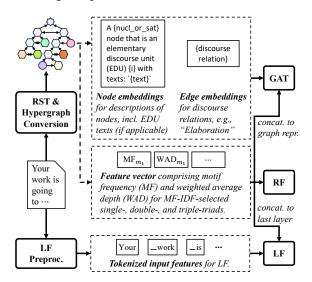


Figure 4: A data flow diagram illustrating how various types of input features are fed into the three baseline models: Graph Attention Network (GAT), Random Forest (RF), and Longformer (LF).

**Formality Scorer.** To identify a relationship between formality style of texts and discourse hierarchies (studied in §5.3), we make use of a publicly available sentence-level formality scorer (Babakov et al., 2023) that is finetuned on a dataset for formality style transfer task, GYAFC (Rao and Tetreault, 2018). The model is based on a ROBERTA-BASE model (Zhang et al., 2020).

#### 4.2 Datasets

Our experiments are conducted on two existing benchmark datasets and one new dataset (TENPAGESTORIES) that we created for more indepth analysis.

HC3-ENGLISH. Guo et al. (2023) published the first corpus on human writer vs. ChatGPT comparison. It comprises 24K *paired* responses from human authors and ChatGPT for mostly question-answering (QA) style prompts across 5 domains: Reddit-ELI5, medicine, finance, open QA, and Wikipedia-CS.

**DEEPFAKETEXTDETECT.** Li et al. (2023) proposed DEEPFAKETEXTDETECT dataset which consists of 448K *mostly unpaired* human-written and machine-generated texts from 10 diverse datasets including news article writing, story generation, opinion writing, etc. The machine-generated texts

 $<sup>^5 \</sup>text{ALLENAI/LONGFORMER-BASE-4096}, available on <math display="block"> \text{https://huggingface.co.}$ 

are from 27 mainstream LLMs from 7 sources such as OpenAI and LLaMA. Out of ten domains in the dataset, we present results on the following five considering the diversity and sizes of domains: review writing (YELP), story writing (WP), argument writing (CMV), news summarization (XSUM), and writing descriptions for scientific tables (SCI\_GEN).

**TENPAGESTORIES (Ours).** Existing datasets for the authorship detection task primarily contain short documents of 200 to 500 tokens, limiting the capture of long-term discourse relations; to overcome this limitation, we construct exceptionally lengthy generations based on 20 fictional stories on Project Gutenberg published in early November 2023.<sup>6</sup>

Our generated content is prepared in three distinct settings:

- 1. **Unconstrained**: We provide the first human-written paragraph of a story and instruct the LLM to iteratively generate continuations of the story. However, due to the limited input length of the existing discourse parser, we did not use the unconstrained generations in experiments.
- 2. "Fill-in-the-gap": We mask  $N \in {1,3,5}$  paragraphs between a preceding human-written paragraph and a subsequent human-written paragraph.
- 3. Constrained "fill-in-the-gap": Similar to (2), but the masked paragraph(s) now include the first and last sentences of the corresponding human-written paragraph(s), providing more guided contexts.

These settings are designed to enable the observation of long-term discourse patterns while varying the constraints applied to the original contexts. We generate content up to 8K tokens ( $\approx$ 10 A4 pages) in an iterative manner, continuing from the previously generated texts. We provide more details on the dataset construction in Appendix E.

## 5 Results

We present findings from three sets of experiments:

- 1. Evaluation of the utility of discourse motifs in the authorship detection task including the "paraphrase attack" scenario (§5.1)
- 2. Investigation of structural variations after paraphrasing (§5.2)

3. Identifying the relationship between formality and hyperedges (§5.3)

## 5.1 Human vs. Machine Authorship Detection

In this section, we report F1 scores for the binary classification task, following the experimental setup detailed in Section 4.1.

HC3 and DEEPFAKETEXTDETECT. From Table 3, we can see that incorporating motif information consistently enhances classification performance across various base encoders, underscoring the effectiveness of discourse motifs as auxiliary information for capturing deeper linguistic structures beyond surface lexicons. Specifically, when the GAT model leverages the overall hierarchical discourse structure, notable performance gains are observed (e.g.,  $0.67 \rightarrow 0.73$  on HC3). Additionally, given the LF model's proficiency in comprehending lengthy texts, the inclusion of discourse motifs offers explicit structural insights, further enhancing the performance. Even on the paraphrased outof-domain test set ("OOD-Para"), where samples come from "unseen domains" and "unseen models," augmenting the detector with discourse motifs yields significant improvements over the baselines.

Models	HC3	DeepfakeTextDetect		
Models	Test	Test	OOD	OOD-Para
RF (All Motifs)	0.55	0.57	0.48	0.60
RF (Motifs)	0.55	0.58	0.49	0.61
GAT	0.67	0.68	0.66	0.53
GAT+Motifs	0.73	0.72	0.67	0.54
LF	0.97	0.90	0.74	0.60
LF+Motifs	0.98	0.93	0.82	0.62

Table 3: Overall detection F1 scores of models on the benchmark datasets. "RF (All Motifs)" denotes Random Forest models with all of the found motifs (3,055 of them) considered as inputs. "RF (Motifs)" indicates Random Forest models where only the MF-IDF-selected motifs (207 of them) are taken as inputs.

**TENPAGESTORIES.** In a more open-ended exploration, we deploy the trained LF model augmented with discourse motifs to analyze our TENPAGESTORIES dataset.

First, we check how the different generation settings affect the detection performance (Table 4). As expected the detection model performs better with longer generated texts (i.e., 1, 3, and 5 paragraphs). Also, when the first and last sentences of the human-written texts are additionally provided

<sup>&</sup>lt;sup>6</sup>https://www.gutenberg.org

Models	FIG			Constrained FIG		
Models	1	3	5	1	3	5
LF LF+Motifs	0.30	0.40	0.45	0.59	0.50	0.55
LF+Motifs	0.43	0.65	0.71	0.71	0.69	0.70

Table 4: Overall detection F1 scores on TENPAGESTO-RIES where "FIG" refers to the "fill-in-the-gap" setting. The numbers 1, 3, and 5 indicate the number of paragraphs to be generated (or "filled in").

in the prompt (i.e., "constrained FIG"), the corresponding generations are easier to detect. We note that the models' f1 scores are low as they produce a lot of false negatives, i.e., predicting actual human-written texts to be machine-generated. However, we can observe that the addition of a small discourse vector significantly improves the performance.

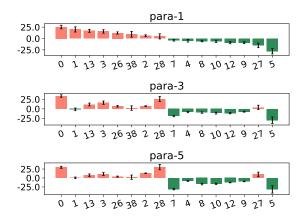


Figure 5: Difference distribution of motifs for TEN-PAGESTORIES under the *fill-in-the-gap* settings. The x-axis represents unique indices of single-triads while the y-axis shows the difference in motif frequency (scaled by 1e-3) of machine-generated and human-written texts for each motif. Discourse motifs indexed at 0 (Elaboration), 5 (Joint), 7 (Joint), and 28 (Temporal) seem to be useful in distinguishing the two groups.

Similarly, Figure 5 elucidates the distinct motif distributions within TENPAGESTORIES across various settings. It shows that as the generation length increases, motifs indexed at 0 (Elaboration) and 28 (Temporal) become increasingly prevalent in machine-generated texts. Conversely, motifs indexed at 7 (Joint) and 5 (Joint) appear more frequently in texts authored by humans. This may imply that human authors tend to construct narratives with a higher prevalence of Joint relations, indicative of more evenly branching structures, in contrast to the patterns found in texts generated by LLMs.

#### 5.2 Structural Variations after Paraphrasing

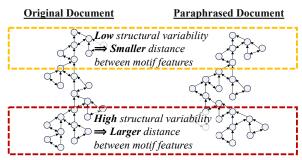


Figure 6: We posit that sentence-level paraphrasing typically leads to greater variability in lower segments of discourse structures compared to higher segments near the root node.

In this experiment, we look at the discourse structures of original and paraphrased documents regardless of their authorship. We aim to validate our hypothesis that sentence-level paraphrasing induces greater variability in the lower segments of discourse graphs, as these segments correspond to the individual EDUs and their aggregation into sentences. Given the assumption that the overall higher-level discourse structure remains relatively intact, we anticipate that the upper segments will display less structural variability (c.f. Figure 6). To this end, we calculate the discourse motif features (§3.2.1) for both upper and lower segments of the graphs corresponding to each pair of original and paraphrased documents. Here, the upper segment denotes the neighborhood graph of the root node within an edge distance of 3. Similarly, the lower segment encompasses the union of neighborhood graphs of individual EDU nodes within one edge distance. We average the motif features over the original and paraphrased documents and compute the absolute distance between the two groups.

Figure 7 depicts the results, delineated by motif frequency and weighted average depth. Notably, both features indicate that the lower segments of graphs display a greater absolute distance compared to the upper portions, thereby supporting our hypothesis.

## **5.3** Formality Scores and Hyperedges

Concerning text styles, we note that the most noticeable disparity among the various text domains appears to lie in the formality dimension. Thus, we computed the formality score for every text generated by humans and machines (c.f. §4.1). Also, for every corresponding document graph, we count the (normalized) frequency of hyperedges for all

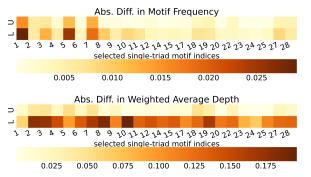


Figure 7: Two heatmaps illustrating the absolute differences of normalized motif counts (above) and weighted average depth (below) between motifs present in discourse graphs of original and paraphrased texts. A graph is further divided into two separate subgraphs (i.e., "upper (U)" and "lower (L)") depending on the distance from the root node.

motifs present in the graph. This is because the hyperedge relations represent higher-level abstractions of text and are thus linked to the discourse hierarchy. From there, we compute the Pearson correlation coefficient between the scores and frequencies. As shown in Table 5, the human-written texts over all domains exhibit weak to moderate correlation with the formality style while machinegenerated ones seem to show no correlation. This hints that when humans write more structured texts in terms of discourse, the contents are more formal.

Machine-Generated	Human-Written	
0.08	0.39	

Table 5: Pearson's R between formality scores and frequency of hyperedges.

#### **5.4** Other Experiments

We conduct preliminary analyses on discourse motifs and investigate the correlation between document graph shapes across different domains in Appendices D and H, respectively.

#### 6 Conclusion

As texts generated by LLMs become increasingly challenging to differentiate from those authored by humans, our approach involves identifying discernible, unique properties inherent in human-crafted texts. We posit that these distinctive signatures manifest at levels beyond mere surface structure and thus opt for looking into their hierarchical discourse structures. Viewing the structures as recursive hypergraphs and conducting motif analyses on them, we find that these motifs are useful in the

authorship detection task and can highlight the subtle structural differences in the two author groups depending on the domains of texts and their corresponding nature. Future plans include extending this approach to long documents by merging multiple document graphs and incorporating topological information beyond discourse.

#### 7 Limitations

In our study, the validity of our experiments and subsequent findings relies on the parsed discourse trees generated by an existing parser. It is important to acknowledge the potential for using alternative discourse frameworks like Segmented Discourse Representation Theory (Lascarides and Asher, 2007) though finding or training a robust parser is a challenge. Despite employing a stateof-the-art model for RST parsing, it is crucial to recognize that parsed results may still be suboptimal due to the inherent difficulty of hierarchical discourse parsing. This challenge is exacerbated by the limited scale of existing datasets used for model training and the inherent ambiguities present in discourse relations. In this aspect, an important future direction would be to build a more robust discourse parser using current LLMs.

In our experimental setup, motif analyses were carried out using single-, double-, and triple-triads. It is worth considering that longer and potentially significant patterns may exist beyond these sizes. However, expanding the motif sizes within our current computational approach is not feasible. Nevertheless, there is potential to capture the distributions of longer patterns through approximation techniques, such as subsampling or leveraging deep learning methodologies.

In the experiment investigating the correlation between formality scores and the number of hyperedges (Section 5.3), it is important to acknowledge a limitation: the formality scorer had been trained on short sentence-level inputs and may not be adequately suited for assessing the formality of longer documents spanning multiple paragraphs.

Furthermore, our study highlights the need for more comprehensive analysis within the TEN-PAGESTORIES dataset. Exploring the correlation between discourse structures and specific linguistic features tailored to creative writing, such as tagged events between named entities, presents an intriguing avenue for future research.

#### 8 Ethical Statement

While our research aims to distinguish linguistic features in human-written and machine-generated texts to improve authorship detection, it is important to acknowledge the potential for these features to be utilized in developing LLMs capable of generating texts that more closely resemble human-authored ones. To address this concern, we advocate for the implementation of transparent reporting practices, ethical review processes, responsible use policies, and collaboration among stakeholders to ensure the ethical development and deployment of LLMs.

Throughout the paper, we have referenced datasets and tools utilized in our experiments, ensuring they originate from open-source domains and do not pose any conflicts with the public release or usage of this paper. Our results are also consistent with the licensing terms of the open-source domain from which the datasets and tools were sourced. We also note that our constructed dataset, TENPAGESTORIES, stems from fictional ebooks available in the public domain and contains no information that names or uniquely identifies individual real people, nor does it include any offensive content.

We acknowledge the use of Grammarly and ChatGPT 3.5 for correcting any less fluent sentences but not for generating new content.

#### Acknowledgements

This research was primarily funded by a generous research gift from Grammarly. Kwang Hee Lee received funding from the National Research Foundation of Korea (NRF), which is supported by the Korean government (MSIT), under Grant No. RS-2024-00345567. We are grateful to the Minnesota NLP group members for their valuable feedback and constructive comments on our initial draft. We also acknowledge Jonghwan Hyeon and Chae-Gyun Lim from KAIST for their assistance with the design of the figures.

#### References

Kenza Amara, Rex Ying, Zitao Zhang, Zhihao Han, Yinan Shan, Ulrik Brandes, Sebastian Schemm, and Ce Zhang. 2022. Graphframex: Towards systematic evaluation of explainability methods for graph neural networks. *Preprint*, arXiv:2206.09677.

Wissam Antoun, Virginie Mouilleron, Benoît Sagot, and Djamé Seddah. 2023. Towards a robust detection of

language model-generated text: Is ChatGPT that easy to detect? In Actes de CORIA-TALN 2023. Actes de la 30e Conférence sur le Traitement Automatique des Langues Naturelles (TALN), volume 1 : travaux de recherche originaux – articles longs, pages 14–27, Paris, France. ATALA.

Nikolay Babakov, David Dale, Ilya Gusev, Irina Krotova, and Alexander Panchenko. 2023. Don't lose the message while paraphrasing: A study on content preserving style transfer. In *Natural Language Processing and Information Systems*, pages 47–61, Cham. Springer Nature Switzerland.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *Preprint*, arXiv:2004.05150.

Daria Beresneva. 2016. Computer-generated text detection using machine learning: A systematic review. In *International Conference on Applications of Natural Language to Data Bases*.

Amrita Bhattacharjee, Tharindu Kumarage, Raha Moraffah, and Huan Liu. 2023. Conda: Contrastive domain adaptation for ai-generated text detection. *Preprint*, arXiv:2309.03992.

Amrita Bhattacharjee and Huan Liu. 2023. Fighting fire with fire: Can chatgpt detect ai-generated text? *Preprint*, arXiv:2308.01284.

Leo Breiman. 2001. Random forests. *Machine Learning*, 45(1):5–32.

Megha Chakraborty, S.M Towhidul Islam Tonmoy, S M Mehedi Zaman, Shreya Gautam, Tanay Kumar, Krish Sharma, Niyar Barman, Chandan Gupta, Vinija Jain, Aman Chadha, Amit Sheth, and Amitava Das. 2023. Counter Turing test (CT2): AI-generated text detection is not as easy as you may think - introducing AI detectability index (ADI). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2206–2239, Singapore. Association for Computational Linguistics.

Simon Corston-Oliver, Michael Gamon, and Chris Brockett. 2001. A machine learning approach to the automatic evaluation of machine translation. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, pages 148–155, Toulouse, France. Association for Computational Linguistics.

Doraid Dalalah and Osama M.A. Dalalah. 2023. The false positives and false negatives of generative ai detection tools in education and academic research: The case of chatgpt. *The International Journal of Management Education*, 21(2):100822.

Yao Dou, Maxwell Forbes, Rik Koncel-Kedziorski, Noah A. Smith, and Yejin Choi. 2022. Is GPT-3 text indistinguishable from human text? scarecrow: A framework for scrutinizing machine text. In *Proceedings of the 60th Annual Meeting of the Association* 

- for Computational Linguistics (Volume 1: Long Papers), pages 7250–7274, Dublin, Ireland. Association for Computational Linguistics.
- Liam Dugan, Daphne Ippolito, Arun Kirubarajan, Sherry Shi, and Chris Callison-Burch. 2023. Real or fake text?: Investigating human ability to detect boundaries between human-written and machinegenerated text. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(11):12763–12771.
- Andy Extance. 2023. Chatgpt has entered the classroom: how llms could transform education. *Nature*, 623(7987):474–477.
- Edyta Gołąb-Andrzejak. 2023. The impact of generative ai and chatgpt on creating digital advertising campaigns. *Cybernetics and Systems*, page 1–15.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *Preprint*, arXiv:2301.07597.
- Xinlei He, Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2024. Mgtbench: Benchmarking machine-generated text detection. *Preprint*, arXiv:2303.14822.
- Patrick Huber and Giuseppe Carenini. 2022. Towards understanding large-scale discourse structures in pretrained and fine-tuned language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2376–2394, Seattle, United States. Association for Computational Linguistics.
- Cliff A Joslyn and Kathleen E. Nowak. 2017. Ubergraphs: A definition of a recursive hypergraph structure. *ArXiv*, abs/1704.05547.
- John Kirchenbauer, Jonas Geiping, Yuxin Wen,
   Jonathan Katz, Ian Miers, and Tom Goldstein. 2023.
   A watermark for large language models. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages 17061–17084. PMLR.
- Ryuto Koike, Masahiro Kaneko, and Naoaki Okazaki. 2023. Outfox: Llm-generated essay detection through in-context learning with adversarially generated examples. *ArXiv*, abs/2307.11729.
- Andreas Kopf, Yannic Kilcher, Dimitri von Rutte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Rich'ard Nagyfi, ES Shahul, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. Openassistant conversations democratizing large language model alignment. *ArXiv*, abs/2304.07327.

- Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Frederick Wieting, and Mohit Iyyer. 2023. Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Alex Lascarides and Nicholas Asher. 2007. Segmented discourse representation theory: Dynamic semantics with discourse structure. In *Computing meaning*, pages 87–124. Springer.
- Yafu Li, Qintong Li, Leyang Cui, Wei Bi, Longyue Wang, Linyi Yang, Shuming Shi, and Yue Zhang. 2023. Deepfake text detection in the wild. *Preprint*, arXiv:2305.13242.
- Xiaoming Liu, Zhaohan Zhang, Yichen Wang, Hang Pu, Yu Lan, and Chao Shen. 2023a. Coco: Coherence-enhanced machine-generated text detection under data limitation with contrastive learning. *Preprint*, arXiv:2212.10341.
- Yikang Liu, Ziyin Zhang, Wanyang Zhang, Shisen Yue, Xiaojing Zhao, Xinyuan Cheng, Yiwen Zhang, and Hai Hu. 2023b. Argugpt: evaluating, understanding and identifying argumentative essays generated by gpt models. *Preprint*, arXiv:2304.07666.
- Zeyan Liu, Zijun Yao, Fengjun Li, and Bo Luo. 2023c. Check me if you can: Detecting chatgpt-generated academic writing using checkgpt. *Preprint*, arXiv:2306.05524.
- Zhengyuan Liu, Ke Shi, and Nancy Chen. 2021. DMRST: A joint framework for document-level multilingual RST discourse segmentation and parsing. In *Proceedings of the 2nd Workshop on Computational Approaches to Discourse*, pages 154–164, Punta Cana, Dominican Republic and Online. Association for Computational Linguistics.
- William C Mann and Sandra A Thompson. 1987. *Rhetorical structure theory: A theory of text organization.* University of Southern California, Information Sciences Institute Los Angeles.
- R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. 2002. Network motifs: Simple building blocks of complex networks. *Science*, 298(5594):824–827.
- Eleni Miltsakaki, Rashmi Prasad, Aravind Joshi, and Bonnie Webber. 2004. The Penn Discourse Treebank. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal. European Language Resources Association (ELRA).
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff,

- and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In *Proceedings* of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.
- Sudha Rao and Joel Tetreault. 2018. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 129–140, New Orleans, Louisiana. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi. 2023. Can ai-generated text be reliably detected? *Preprint*, arXiv:2303.11156.
- Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M. Borgwardt. 2011. Weisfeiler-lehman graph kernels. *J. Mach. Learn. Res.*, 12:2539–2561.
- Zhouxing Shi, Yihan Wang, Fan Yin, Xiangning Chen, Kai-Wei Chang, and Cho-Jui Hsieh. 2023. Red teaming language model detectors with language models. *Preprint*, arXiv:2305.19713.
- Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.
- Jinyan Su, Terry Zhuo, Di Wang, and Preslav Nakov. 2023. DetectLLM: Leveraging log rank information for zero-shot detection of machine-generated text. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12395–12412, Singapore. Association for Computational Linguistics.
- Frank W. Takes, Walter A. Kosters, Boyd Witte, and Eelke M. Heemskerk. 2018. Multiplex network motifs as building blocks of corporate networks. *Applied Network Science*, 3.

- Adaku Uchendu, Jooyoung Lee, Hua Shen, Thai Le, Ting-Hao 'Kenneth' Huang, and Dongwon Lee. 2023. Does human collaboration enhance the accuracy of identifying llm-generated deepfake texts? *Preprint*, arXiv:2304.01002.
- Christoforos Vasilatos, Manaar Alam, Talal Rahwan, Yasir Zaki, and Michail Maniatakos. 2023. Howkgpt: Investigating the detection of chatgpt-generated university student homework through context-aware perplexity analysis. *Preprint*, arXiv:2305.18226.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. *International Conference on Learning Representations*.
- Vivek Verma, Eve Fleisig, Nicholas Tomlin, and Dan Klein. 2023. Ghostbuster: Detecting text ghostwritten by large language models. *Preprint*, arXiv:2305.15047.
- Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan, Derek F. Wong, and Lidia S. Chao. 2023. A survey on llm-generated text detection: Necessity, methods, and future directions. *Preprint*, arXiv:2310.14724.
- Wen Xiao, Patrick Huber, and Giuseppe Carenini. 2021. Predicting discourse trees from transformer-based neural summarizers. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4139–4152, Online. Association for Computational Linguistics.
- Naganand Yadati. 2020. Neural message passing for multi-relational ordered and recursive hypergraphs. In *Advances in Neural Information Processing Systems*, volume 33, pages 3275–3289. Curran Associates, Inc.
- Lingyi Yang, Feng Jiang, and Haizhou Li. 2023. Is chatgpt involved in texts? measure the polish ratio to detect chatgpt-generated text. *Preprint*, arXiv:2307.11380.
- Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. 2019. Gnnexplainer: Generating explanations for graph neural networks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Xiao Yu, Yuang Qi, Kejiang Chen, Guoqiang Chen, Xi Yang, Pengyuan Zhu, Weiming Zhang, and Neng H. Yu. 2023. Gpt paternity test: Gpt generated text detection with gpt genetic inheritance. *ArXiv*, abs/2305.12519.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. *ArXiv*, abs/1905.12616.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2020. Pegasus: pre-training with extracted

gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org.

Xuandong Zhao, Lei Li, and Yu-Xiang Wang. 2022. Distillation-resistant watermarking for model protection in NLP. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5044–5055, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

#### A Comparison with prior work

With the advent of recent LLMs, the need, as well as the challenges in distinguishing machine-generated texts, has risen significantly. Various existing methodologies encompass fine-tuning LLMs (Guo et al., 2023; Liu et al., 2023b,c; Li et al., 2023) with supervised datasets, potentially utilizing techniques like contrastive learning (Liu et al., 2023a; Bhattacharjee et al., 2023), adversarial learning (Shi et al., 2023; Yang et al., 2023), or human-assisted learning (Uchendu et al., 2023; Dugan et al., 2023). Additionally, there are zero-shot methods (Beresneva, 2016; Vasilatos et al., 2023) that leverage statistical features of texts or intermediate values within LLMs; watermarking techniques (Zhao et al., 2022; Kirchenbauer et al., 2023) that introduce "green tokens" or embed "secret keys" as vectors to generated outputs; and prompt-based approaches that make use of (another) LLM as detectors (Zellers et al., 2019; Koike et al., 2023; Yu et al., 2023; Bhattacharjee and Liu, 2023). A comprehensive overview of these approaches can be found in Wu et al. (2023).

While the relevant recent studies on the detection task primarily adopt surface-level features, such as distributions of n-grams and part-of-speech (POS) tags, we approach the problem from another direction: using hierarchical discourse features within texts. This direction is reminiscent of the work by Corston-Oliver et al. (2001), who explored a comparable path in the pre-LLM era. Their study investigated the efficacy of the branching features observed in syntactic parse trees to determine the origin of texts, particularly whether they were generated by a machine translation model. Our work considers discourse-level features spanning multiple sentences and paragraphs (as opposed to sentence-level POS features).

The existing challenges in the field encompass several facets. First, there are pronounced out-of-distribution (OOD) challenges (Li et al., 2023; Antoun et al., 2023), wherein the detectors struggle when confronted with texts that fall outside the learned distribution. Second, the detectors are prone to potential attacks (Shi et al., 2023; He et al., 2024), such as paraphrasing attacks (Sadasivan et al., 2023; Krishna et al., 2023), where machine-generated texts undergo further paraphrasing to alter the distribution of lexical and syntactic features. This dynamic renders detectors reliant on surface-level features and watermarking technology ineffective. Finally, the inherent ambiguities present in the two groups of texts have become *more nuanced* over time, making it progressively challenging to distinguish them. In this paper, we touch upon these three challenges by showing that the addition of discourse features (i) improves the detection of OOD samples and (ii) trains a more robust classifier against paraphrased attacks. Furthermore, our analyses, utilizing discourse network motifs, shed light on the nuanced distinctions within the hierarchical structures of the two text categories.

### B Recursive hypergraphs to their standard traditional graph form

In this section, we introduce how recursive hypergraph representations of RST trees can be transformed into standard traditional graph forms. Suppose an RST tree consists of m EDUs  $\{EDU_1, \cdots, EDU_m\}$  with their relations such as  $(EDU_1, EDU_2)$  and  $((EDU_1, EDU_2), EDU_3)$  (see EDU1, EDU2, EDU3 in Figure 8 for example). Note that in this section, we ignore edge labels for better readability. A recursive hypergraph representation of the given RST in Figure 8 is a tuple (V, E) where a set of vertices  $V = \{EDU_1, \cdots, EDU_7\}$  and a set of edges  $E = \{e_{1-2} = (EDU_1, EDU_2), e_{1-3} = (e_{1-2}, EDU_3), \cdots, e_{1-7} = (e_{1-3}, e_{4-7})\}$ . To facilitate ease of analysis, the hypergraph can be transformed into a standard traditional graph form as follows.

**Definition 1** A standard traditional graph form (V', E') of a recursive hypergraph representation (V, E) of an RST tree is defined by  $V' = V \cup E$  and  $E' = \{(u, w) | (u, w) \in E \text{ or } u \in w\}$ .

This graph is thought of as an "extended" Levi graph of a recursive hypergraph (Joslyn and Nowak, 2017). We call it simply the transformed graph of an RST tree if there is no ambiguity.

In the definition of the transformed graph, each relation in an RST tree constitutes a triangle graph. A triangle graph is a graph with three vertices where there is at least one edge between two vertices. For

<sup>&</sup>lt;sup>7</sup>We can consider bi-direction by adding a reverse order of an edge.

<sup>&</sup>lt;sup>8</sup>Different from the definition of Levi graphs in (Joslyn and Nowak, 2017), we allow cycles in the extended Levi graphs.

#### **Quote Example**

Your work is going to fill a large part of your life, and the only way to be truly satisfied is to do what you believe is great work. And the only way to do great work is to love what you do. If you haven't found it yet, keep looking. Don't settle. As with all matters of the heart, you'll know when you find it. — Steve Jobs

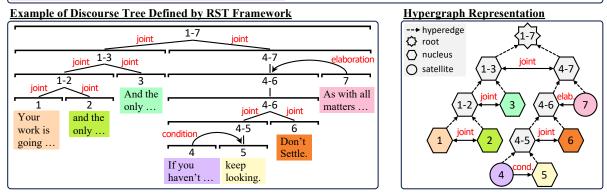


Figure 8: Steve Job's motivational quote and its corresponding RST tree.

example, a relation r = (x, y) (x and y are either an EDU or an edge) is transformed to the triangle graph,  $V' = \{r, x, y\}, E' = \{(x, y), (x, r), (y, r)\}$ . Now, we generalize this observation as follows.

**Theorem 1** The transformed graph of an RST tree only consists of the union of triangle graphs.

**Proof sketch by induction**. Before we prove, assume that every finite RST can be represented by a huge, single relation in a recursive manner. For example, the RST in Figure 8 can be represented by

$$(((EDU_1, EDU_2), EDU_3), (EDU_7, ((EDU_4, EDU_5), EDU_6))).$$
 (6)

Consider an RST tree (EDU<sub>1</sub>, EDU<sub>2</sub>) as an initial case. Then, the transformed graph is a triangle graph. Suppose x that is an RST tree with n-1 EDUs is the union of triangle graphs. Then, an RST tree  $(x, \text{EDU}_n)$  or  $(\text{EDU}_n, x)$  is also a graph with the union of triangle graphs by the definition.

More generally, suppose that an RST tree with less than n EDUs forms the union of triangle graphs. For an RST tree z containing n EDUs, this RST can be divided into an RST tree x with n-k EDUs and an RST tree y with k EDUs, i.e., z=(x,y). Since both transformed graphs of x and y are unions of triangle graphs by the assumption, and a relation (x,y) can be transformed into the triangle graph as discussed above, the RST z can be also transformed into a new graph with unions of triangle graphs.

#### C Calculation of motif features

Figure 9 illustrates the calculation of the two features, motif frequency and weighted average depth, of a discourse motif with index 18. This motif is one of the selected single-triads that exhibit MF-IDF scores (§3.2.2) surpassing at least one standard deviation.

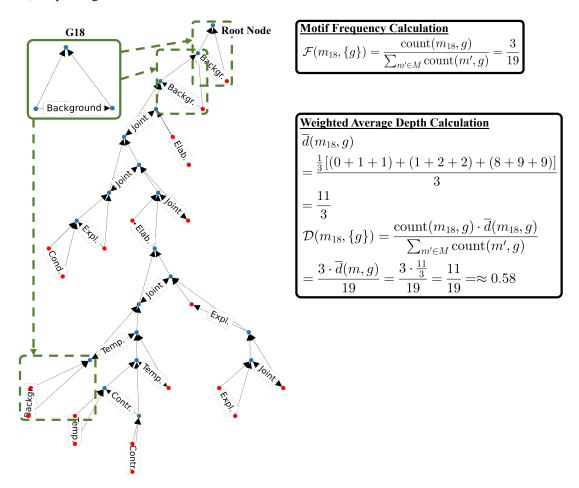


Figure 9: An example of calculating motif frequency and weighted average depth for the single-triad  $m_{18}$ .

#### D Preliminary analyses on discourse motifs

We begin our preliminary analyses by examining the relationship between the proposed discourse motifs and influential edges identified by the GAT explainer (§D.1). Subsequently, we calculate the difference distributions of motifs across five domains to emphasize their distinctive features (§D.2).

## D.1 Correlation between Motifs and Influential Edges for Explanation

As a way to evaluate the appropriateness of using discourse motifs, we explore how they are related to influential edges in GATs computed by an explanation method. Specifically, we utilize explainability techniques such as "GNNExplainer" (Ying et al., 2019) and "GraphFramEx" (Amara et al., 2022) for graph classification, aiming to assess their correlation with our proposed motif features. These methods elucidate the predictions of GNNs by identifying the nodes or edges whose masking has the most significant impact on the predicted outcome. Since our graphs consist of directed and labeled edges, we require an explanation method tailored to edge-based features for GAT. Thus, we employ "AttentionExplainer" from GraphFramEx, which computes edge masks using the attention coefficients generated by GATs trained on the binary classification task.

Note that the influential edges refer to edges with high masking values<sup>9</sup> generated by a GAT explainer as described in (§4). We then calculate a correlation between frequencies of selected discourse motifs based

<sup>&</sup>lt;sup>9</sup>The most influential edges whose masking values are larger than 0.99 are considered in [0, 1].

on MF-IDF scores and those of influential edges. To this end, we define two random variables as follows. Let X be the number of single-triads and let Y be the number of influential edges for each graph. Samples of X and Y are collected from each graph in the datasets (HC3-ENGLISH, DEEPFAKETEXTDETECT), and the Pearson correlation coefficient  $r_{XY}$  is calculated.

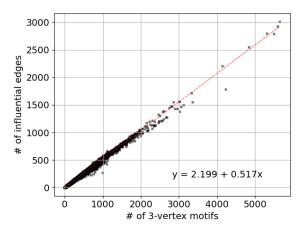


Figure 10: Correlation between selected single-triads with high MF-IDF scores and influential edges computed by a GAT explainer on DEEPFAKETEXTDETECT dataset.

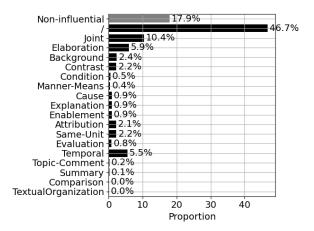


Figure 11: Proportions of labels of influential edges. Note that "Non-influential" refers to edges whose mask values are under the threshold. Note also that the value "0.0%" in this figure does not indicate the zero, but a very small positive number like "1e-5%."

Figure 10 shows that frequencies of selected single-triads and influential edges are highly correlated with  $r_{XY} \approx 0.99$ , indicating that the discourse motifs may be useful features in determining the authorship of texts. Figure 11 shows the proportion of edge labels (i.e., discourse relations) of influential edges. We note that the label "/" in the figure denotes the hyperedge relation ("hyp.").

It is noteworthy that while the motifs are motivated by the discourse framework rooted in linguistics, the notion of influential edges stems from a machine learning-driven, data-centric approach. The strong correlation observed between the two is particularly interesting.

#### **D.2** Difference Distributions of Motifs across Domains

Figure 12 shows bar plots for the difference distribution of single-triad motifs across the five domains of texts. We note that large-scaled plots can be found in Figure 16.

The mean and standard deviation of differences ( $\mathcal{M}(\mathcal{D}_{diff})$ ) are is 0.00012 and 0.0015, respectively. Notably, machine-generated texts tend to feature more motifs with Background relations, while human-written texts exhibit an inclination towards Temporal discourse relations in domains such as review writing (YELP) and story generation (WP), and contain more Joint relations in domains that demand formality or deeper logical depths (i.e., argument writing (CMV) and news summarization (XSUM). We note that

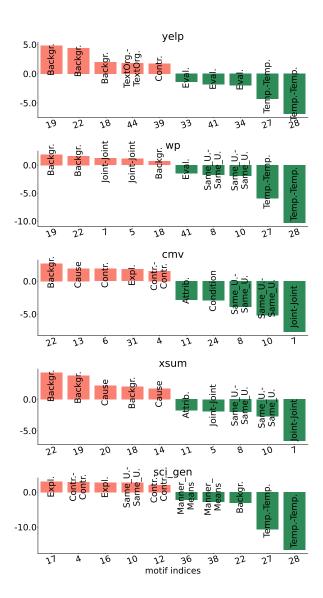


Figure 12: Difference in single-triad motif distribution of machine-generated and human-written texts for five different domains in DEEPFAKETEXTDETECT dataset. Y-axis represent the difference in MF-IDF scores scaled by 1e-3. Full-scaled versions including a larger number of motifs and plots for double- and triple-triads can be found in Appendix J.

Joint discourse relation signifies cases where the corresponding child trees are of equal importance (i.e., "nuclei") and thus branched evenly. The texts in table description generation (SCI\_GEN) are interesting in that the human-written descriptions (taken mostly from academic papers in NLP and machine learning) mainly highlight the changes in performance and this is captured as Temporal, Manner, or Means relations. Table 6 shows some grounded examples of texts for prevalent motifs in each domain.

## E Constructing "TenPageStories" dataset

The "TenPageStories" dataset consists of stories with approximate lengths of 8,000 tokens or 10 A4 pages. It was constructed using an iterative method of calling the OpenAI's GPT-4-1106-PREVIEW model. In total, we conducted 552 request calls, resulting in a total expenditure of 55 USD.

The construction process works by continuously adding the text generated by the model in one call to the prompt of the next call. The iterative method was used for all three generation settings: (1) unconstrained generation, (2) "fill-in-the-gap", and (3) constrained "fill-in-the-gap". Pseudo-code and examples for the three generation settings are provided below. Note that for the constrained fill-in-the-gap example, three different completion calls are made for an iteration: the first completes paragraph 1, the second generates

paragraph 2...n-1, the third completes the last paragraph n, where n is the number of paragraphs to be masked per iteration. The exception to this process is when n=1, then one completion call is made to complete the paragraph based on its first and last sentences.

## **Algorithm 1** Unconstrained Generation

```
1: function UNCONSTRAINEDGENERATION(promptInstruction)
2: prompt ← promptInstruction
3: while length(prompt) < 8000 tokens do
4: generatedText ← generate text using LLM based on the prompt
5: prompt ← prompt + generatedText
6: end while
7: return prompt
8: end function
```

## Algorithm 2 "Fill-in-the-gap"

```
1: function FILLINTHEGAP(promptInstruction, numMasked)
```

- 2: prompt ← promptInstruction + firstParagraph
- 3: **for** each group of numMasked+1 paragraphs from the second paragraph onwards **do**
- 4: generatedText  $\leftarrow$  generate text using the LLM to fill the first numMasked paragraphs in the group
- 5: prompt  $\leftarrow$  prompt + generatedText + lastParagraphInTheGroup
- 6: **end for**
- 7: **return** prompt
- 8: end function

## Algorithm 3 Constrained "fill-in-the-gap"

- 1: **function** CONSTRAINEDFILL(promptInstruction, numMasked)
- 2:  $prompt \leftarrow promptInstruction$
- 3: **for** each group of numMasked paragraphs **do**
- 4: firstParagraph ← generate text using the LLM to fill the first paragraph in the group based on its first sentence
- 5: prompt ← prompt + firstParagraph
- 6:  $middleParagraphs \leftarrow generate text using the LLM to fill all paragraphs in the group besides$ the first and last
- 7: prompt ← prompt + middleParagraphs
- 8: lastParagraph  $\leftarrow$  generate text using the LLM to fill the last paragraph based on its last sentence
- 9: prompt ← prompt + lastParagraph
- 10: end for
- 11: **return** prompt
- 12: end function

#### 1. Free Generation Example Prompt:

Take a look at the story and generate text according to the instructions in brackets []. Only return the generated parts. Here is the story:

Text from previous generation

Text from previous generation

Text from previous generation

[Generate text here]

## 2. "Fill-in-the-gap" Example Prompt (Fill 1 paragraphs):

Take a look at the story and generate text according to the instructions in brackets []. Only return the generated parts. Here is the story:

Original paragraph text

Text from previous generation

Original paragraph text

[Generate paragraph of approximately 94 tokens here]

Original paragraph text

#### 3. Constrained "Fill-in-the-gap" Example Prompt (Fill 3 paragraphs):

Take a look at the story and generate text according to the instructions in brackets []. Only return the generated parts. Here is the story:

Text from previous generation

Text from previous generation

Text from previous generation

Original first sentence of paragraph 1

[Complete preceding paragraph using approximately 67 tokens]

[Generate paragraph of approximately 129 tokens here]

[Complete following paragraph using approximately 87 tokens]

Original last sentence of paragraph 3

## F Details on experimental setup

Our experiments were conducted on a workstation equipped with 4 NVIDIA RTX A5500 GPUs, an AMD Ryzen Threadripper PRO 5975WX 32-core CPU, and 1TB of RAM. When running experiments, we primarily adhered to the default hyperparameter configurations of the baseline models and model optimizer, with precise details available in the accompanying code repository: [TBA].

We report that RF, GAT, and LF models underwent training for 1 epochs ( $\approx 10$  minutes), 10 epochs ( $\approx 7.5$  hours), and 4 epochs ( $\approx 14.5$  hours), respectively, with the best-performing checkpoint chosen based on validation dataset performance.

We also report that parsing RST structures and extracting discourse features across the entire dataset required approximately 9.5 hours, utilizing the multiprocessing capabilities of CPU cores to their fullest extent. We note that these features only need to be computed once offline.

## **G** More examples of grounded motifs in texts

Table 6 shows more examples of texts featuring discourse motifs with high MF-IDF scores.

Domains	Machine-Generated	Human-Written
Daviau Waitina	[Motif #19 (Background)] I recommend reservations. [[We were there at 8:30pm, food was excellent with very friendly staff]] << th>that made us feel welcome right away.>> The wine list had something	[Motif #28 (Temporal)] When Lagasse's Stadium first opened we were excited by the idea of watching the big game [[while at the same time enjoying some of those great dishes]] < <we -="" emeril="" make="" on="" tv="" watched="">&gt; perfect combination, right? Well, as it turns out</we>
Review Writing (YELP)	[Motif #22 (Background)] Buffets began popping up along the Strip late last decade < <when a="" casinos="" decided="" number="" of="">&gt;&gt; [[they wanted to boost traffic by offering some form of cheap dining experience without offending guests with lower cost options like pizza slices or pasta salads.]] While most buffets offer</when>	[Motif #34 (Evaluation)] [[The stir fry bowl was very fresh, the one thing I noticed was that they probably don't clean their stir fry grill that often between bowls as some of my beef tasted like fish (not good for those with food allergies), but I would assume you could ask them to clean it first.]] < <could a="" beef="" been="" bit="" bowl="" for="" have="" in="" more="" price.="" the="">&gt; Not sure if I'd go again.</could>
Story Generation (WP)	[Motif #19 (Background)] 10-year-old Lily eagerly awaited the arrival of hers. She had always been a good girl, always sharing and playing nicely with others. As she sat in her room, [[she wondered]] < <what animal="" assigned="" be="" her.="" kind="" of="" to="" would="">&gt; Would it be a furry puppy? A fluffy kitten? Or something more exotic</what>	[Motif #28 (Temporal)] I back to the side of the stage and the stage lights come up. < <suddenly audience="" backdrop.="" can="" i="" illuminated="" in="" of="" parts="" see="" the="">&gt; [[Pretty much everybody - a hundred or so, that is - got quiet after I started arguing with Jon about payment. Now there's clapping, screaming, and I think some lady took off her shirt way back there.]] So apparently the new guy is popular</suddenly>
	[Motif #7 (Joint-Joint)] Elizabeth went on to describe her life, her hopes, and her dreams. < <she about="" and="" family,="" friends,="" her="" home.="" wrote="">&gt; [[She even included a few sketches of her hometown, which had long since faded away.]] As I read, I felt a strange connection to Elizabeth. Though she lived in</she>	[Motif #10 (Same_Unit-Same_Unit)] It's not about literal wealth but instead "personal wealth." <there's a="" in="" man="" middle="" nowhere="" of="" the="">&gt; [[living with deer]] that uses up a lot of our Happiness supply. There's a janitor that spends his life cleaning up after others who gets a fresh supply every day.</there's>
Argument Writing (CMV)	[Motif #13 (Cause)] For example if you want to know something about your relative or friend but they don't answer your text messages < because he doesn't use his cellphone>> <<(because there would no longer be any need)>> then what else can you expect him to do except send an email? Well, unfortunately	[Motif #7 (Joint-Joint)] < <it's already="" been="" etc.,="" factory="" happening="" online="" shopping,="" with="" workers,="">&gt; [[and with self-driving cars on the horizon and computing technology getting better and better, the long-term job security for many professions is shaky at best.]] This surge in unemployment may even cause a great economic crisis.</it's>
	[Motif #6 (Contrast)] To start off, yes, I am Mormon <and am="" book="" by="" i="" mormon="" musical.="" no="" not="" of="" offended="" the="">&gt; <in a="" devout="" fact,="" i="" know="" mormons="" of="" ton="" very="">&gt; who I saw The Book of Mormon musical in Chicago a couple of weeks ago and it was incredible!</in></and>	[Motif #24 (Condition)] but what I'm saying is that they shouldn't be impeached or forced to resign, or face criticism on their ability to lead or shape policy. << If they are a good policy maker fine.>> << If they are making the world a better place>> [[and fighting for things that other politicians are afraid to]] then fine.
News Summarization (XSUM)	[Motif #22 (Background)] Hamilton went on to win and now trails Rosberg, who was given a five second penalty in Austria < <following an="" happened="" into="" investigation="" what="">&gt; [[after he returned ahead at Turn One but made contact with his Mercedes teammate as they headed towards Turns Three and Four respectively.]] Rosberg apologized immediately afterwards</following>	[Motif #7 (Joint-Joint)] the country's biggest artificial lake. The sub-sea cable will connect to the UK network at Blyth in Northumberland. < <crucially, able="" at="" be="" call="" notice.="" on="" power="" short="" the="" to="" uk="" will="">&gt; [[The energy will be used to manage the growing levels of intermittent wind power on the network.]] It will also be a two-way link.</crucially,>
	[Motif #20 (Cause)] The television personality said she was proud to be Welsh but recognized that many people felt differently < because there are different views across the United Kingdom>> [[about whether their vote could lead to independence for Wales from Britain or Scotland leaving Britain along with Northern Ireland into an independent EU bloe.]] She noted	Motif #11 (Attribution) getting rid of fees would help more young people into university, including the disadvantaged? Universities are worried that such a switch to direct funding, dependent on government finances, [would put a limit on places and a brake on expansion.] One of the quiet revolutions of recent years has been the complete removal of limits on student numbers
Table Descriptions (SCI_GEN)	[Motif #17 (Explanation)] For mobile robots navigating on sidewalks, it is essential to be able to safely cross street intersections. [Most existing approaches rely on the recognition of the traffic light signal]] < <to a="" can="" determine="" midcrossing.="" or="" proceed="" robot="" stop="" whether="">&gt; However in our settings</to>	[Motif #28 (Temporal)] In this context, all the polarimetric information seems irretrievably mixed. << A direct model, derived from a simple but original extension of the widespread "multiple scattering model" leads to a high dimensional linear inverse problem.>> [[It is solved by a fast dedicated imaging algorithm that performs to determine at a time three huge 3-D scatterer maps which correspond to HH, VV and HV polarizations at emission and reception.]] It is applied successfully to
	[Motif #4 (Contrast-Contrast)] < <recent (ace)="" automatic="" based="" chord="" extraction="" focused="" has="" improvement="" learning.="" machine="" models="" of="" on="" research="" the="">&gt; &lt;<however, account="" fail="" into="" knowledge.="" models="" most="" prior="" still="" take="" the="" to="">&gt; Unlike the current approaches, prior knowledge is largely ignored in the automatic chord extraction (ACE) problem,</however,></recent>	[Motif #38 (Manner_Means)] The accumulation of litter and plastic debris at the seafloor and the bottom of rivers are extremely harmful for the aquatic life. We propose a solution for monitoring this problem < <using (auvs)="" a="" autonomous="" of="" team="" underwater="" vehicles="">&gt; [[to exchange the recorded video in order to reconstruct the map of regions of interest.]] However, underwater video transmission is a challenge</using>

Table 6: Grounded examples of texts for some prevalent motifs. Exact shapes of the motifs can be found in Fig. 18.

#### H Additional experiment: branched vs. chain-like structures

In this experiment, we differentiate texts by domain and assume that the formality of texts is consistent within the domain. We can then look at the varying levels of graph structures per domain. As illustrated in the Introduction (§1), we assume that the linguistic structure of discourse can change depending on factors such as formality, spontaneity, and depth of reasoning in the texts.

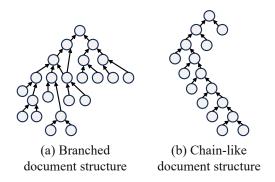


Figure 13: Illustration of document graph structures with different discourse hierarchies. For simplicity, the horizontal edge linking two child nodes is omitted.

One possible dimension of characterizing this difference in structures is by checking whether the structures are evenly branched or follow a more sequential pattern as illustrated in Figure 13.

To quantify this notion, we calculate the average shortest path length (ASPL) of a document graph. For linear or chain-like graphs, the ASPL tends to be relatively short, as nodes are connected in a linear fashion. In contrast, more spread-out graphs will likely have a longer ASPL due to increased distances between nodes. Another closely related metric is the diameter of a graph which is the longest shortest path between any two nodes in the graph. However, as human-written and machine-generated texts are not necessarily paired and could vary a lot in length, we opt for ASPL.

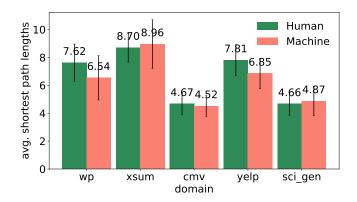
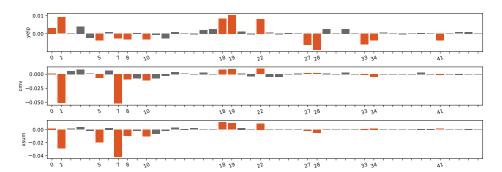


Figure 14: Average shortest path lengths per domain for document graphs. We

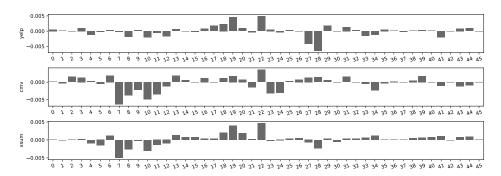
We observe that the news summarization task (XSUM) yields texts with the longest ASPL across both groups. Human-generated texts exhibit longer ASPL in the WP, CMV, and YELP domains. Notably, within these three domains, Joint discourse relations emerge as the most prevalent relation on the human-authored side, affirming our hypothesis regarding the "equal branching" characteristic of this relation (c.f., Fig. 16).

## I Difference distributions computed by motif frequency vs. MF-IDF scores

Figure 15 presents two bar plots, illustrating (a) motif frequency and (b) MF-IDF scores of motifs across three domains. Notably, due to the IDF scaling, the latter plot exhibits a slightly less skewed pattern compared to the former plot.



(a) Difference distribution of motifs computed by using their motif frequencies.



(b) Difference distribution of motifs computed by using their MF-IDF scores.

Figure 15: Two ways of computing the difference distributions.

## J Large-scale difference distributions

Figure 16 displays two bar plots for single- and double-motifs, depicting MF-IDF difference distributions. Triple-triads are excluded due to label clutter.

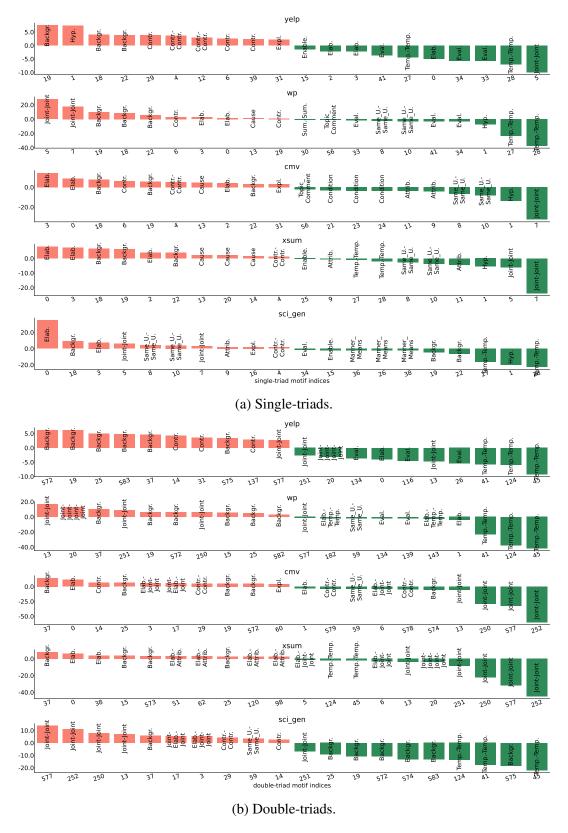


Figure 16: MF-IDF difference distributions of single- and double-triad motifs. To improve readability, the hyperedge relations have been excluded.

## K Graphical examples of identified motifs

Figure 17 illustrates examples of the different types of triads. Similarly, Figure 18 shows all 69 single-triad motifs.

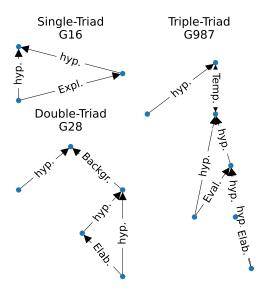


Figure 17: Examples of three types of motifs. More examples of single-triads are drawn in Figure 18.

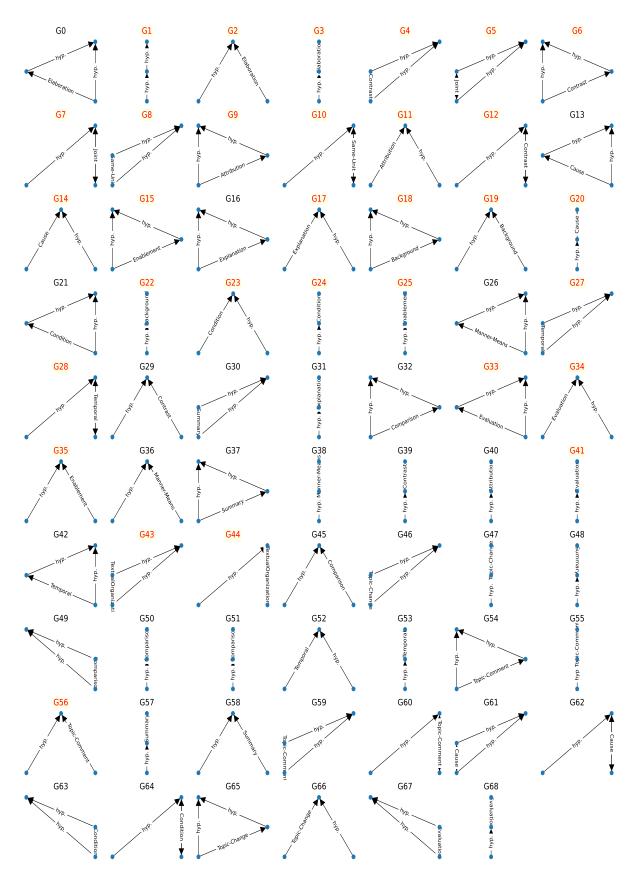


Figure 18: All 69 motifs of size three. The 31 marked motifs are the selected ones that exhibit MF-IDF scores surpassing at least one standard deviation.