

Entailment Progressions: A Framework for Identifying Author Usage of Logic

Anonymous ACL submission

Abstract

Textual entailment, or the ability to deduce whether a proposed hypothesis is logically supported by a given premise, has historically been applied to the evaluation of language modeling efficiency in tasks like question answering and text summarization. However, we believe that these zero-shot entailment evaluations can extend to a sequential evaluation of entailment on a sentence-by-sentence basis within a larger text. We refer to this approach as “entailment progressions”. Additionally, entailment progressions shed light on the intentional logical approaches authors typically employ to construct their arguments, illustrating the points at which authors choose to integrate contradiction and entailment. Our results suggest that entailment progressions can both identify consistency in logical structures and establish a connection between this consistency and how humans typically author texts, as opposed to more formulaic approaches.

1 Introduction

As Large Language Models (LLMs) expand and evolve to accommodate more complex language generation tasks, model developers and researchers have leveraged validation mechanisms to ensure the accuracy of text outputs. These mechanisms, ranging from active feedback mechanisms informed by human input (Christiano et al., 2017; MacGlashan et al., 2017) to passive benchmarking designed to test LLMs using metrics indicative of human linguistic capability (Wang et al., 2018; Lin, 2004; Papineni et al., 2002), work towards the primary goal of bridging the gap between model capability and human language.

Textual Entailment – the ability to deduce whether a proposed hypothesis is logically supported by a given premise (Bentivogli et al., 2009) – has helped modelers understand the inferential capabilities of a given language model. Its origins lie in the belief that for a model to conduct specific

Natural Language Processing (NLP) tasks, it must be capable of the elementary logical inference that underlies these tasks (Zaenen et al., 2005). However, we believe that entailment can potentially describe stylistic choices made by the author of the evaluated text through an examination of how the author chooses to introduce new information or support information they previously provided.

While RTE primarily focuses on the logical relationship between two statements, logical approaches require examining multiple statement-to-statement relationships for coherence. For example, statement (1) alone does not entail (2). However, if we introduce statement (3) in between the previous two statements, a new logical relationship emerges.

- (1) *Blue light has the shortest wavelength in the electromagnetic spectrum.*
- (2) *The sky is blue.*
- (3) *Gas and particles in the sky reflect light with the shortest wavelength.*

In the context of RTE, the former case would yield a single outcome, indicating neutral entailment. In contrast, the latter scenario would generate two outcomes: neutral entailment between (1) and (3); and positive entailment between the first two statements and the final premise. Extending RTE to encompass intermediary logical relationships between a text and its final hypothesis not only enhances our understanding of the employed logical approach, but also clarifies why such an approach was chosen by the author. This raises a pair of interesting and, as far as we are aware, unexplored questions: *is logical reasoning inherently linked to the traits of the author employing it?* And if it is, *can this relationship be identified?*

Our work builds upon previously established definitions and modeling approaches for RTE, demonstrating its applications beyond traditional use cases. In this paper, we:

081 **(1) Introduce entailment progressions**, a frame- 131
082 work in which a given piece of text can be repre- 132
083 sented as a series of values, with each value repre- 133
084 senting the level of textual entailment between two 134
085 consecutive sentences in a text. This entailment 135
086 progression describes the logical flow of the text, 136
087 identifying how new information is introduced in 137
088 relation to the preceding content (in support, in 138
089 rejection, or with no relation). 139

090 **(2) We show that analyzing the entailment pro-** 140
091 **gressions in a set of documents** written by a spe-
092 cific author **can unveil similarities in their logical**
093 **style, that can be attributed to the intentions of**
094 **the author.** 141

095 **(3) We show that entailment progressions are** 142
096 **capable of assessing whether automatically gen-** 143
097 **erated human-like text**, under specific prompted 144
098 conditions, **adhere to an underlying structure** 145
099 that modulates entailment at key points in the text. 146

100 The remainder of the paper is organized as fol- 147
101 lows. Section 2 presents the current state of the 148
102 art, Section 3 describes our methodology, while 149
103 Section 4 details the results of our experiment. We 150
104 conclude by providing some perspectives for future 151
105 work. 152

106 2 Related Work 153

107 The first notable definition of textual entailment 154
108 was formulated as follows: “*T textually entails H* 155
109 *if, typically, a human reading T would infer that H* 156
110 *is most probably true*” (Dagan et al., 2010). This 157
111 definition refined the RTE approach to specifically 158
112 focus on the logical relationship between T and H 159
113 on the basis of human evaluation rather than pre- 160
114 existing notions of implication, in which the plau- 161
115 sibility of T and H could potentially misconstrue 162
116 entailment (Korman et al., 2018). 163

117 Korman (Korman et al. 2018) later expanded 164
118 upon Dagan’s inferential approach to accommodate 165
119 for edge cases associated with human inference like 166
120 irrelevant trivialities, unexpressed conclusions, and 167
121 potential disagreements in human interpretation, 168
122 leading to an alternate understanding of RTE: “*a* 169
123 *text T textually entails a hypothesis H relative to a* 170
124 *group of end users G just in case, typically, a mem-* 171
125 *ber of G reading T would be justified in inferring* 172
126 *the proposition expressed by H from the proposition* 173
127 *expressed by T”. Korman’s definition differs from* 174
128 *Dagan’s in three important ways. First, RTE must* 175
129 *be restricted to a group G due to differing RTE ap-* 176
130 *proaches associated with variability in background* 177

138 knowledge, linguistic proficiency, and other human 139
139 traits that can affect interpretation of logical rela- 140
140 tionships (Bos and Markert, 2005). Second, RTE 141
141 requires justifiable inference in order to allow read- 142
142 ers to assume logical transfer without believing the 143
143 plausibility of T or H (Feldman, 2003). Finally, 144
144 RTE should limit the scope of T and H to the literal 145
145 expression of T and H in order to condition for in- 146
146 ferential effects associated with differing grammar, 147
147 semantic, and syntactical choices (Braun, 2001). 148

149 While textual entailment could be philosophi- 149
150 cally outlined in the broader context of logical infer- 150
151 ence, the practical application of RTE approaches 151
152 was constrained by a limited understanding of the 152
153 linguistic underpinnings and specific criteria that 153
154 govern how expressions entail one another (Amoia, 154
155 2009). This became particularly evident in the 155
156 academic analyses of the Pascal RTE challenges 156
157 (Dagan et al., 2005; Giampiccolo et al., 2007, 2008; 157
158 Bentivogli et al., 2009, 2011), a series of competi- 158
159 tions in which participants tested RTE approaches 159
160 against datasets comprised of premise-hypothesis 160
161 pairs, along with “*support*” and “*reject*” labels 161
162 (Dagan et al., 2005). While these datasets at- 162
163 tempted to capture entailment through the binary 163
164 classification of these premise-hypothesis pairs, 164
165 this approach limits the broader scope of entail- 165
166 ment, which can vary depending on factors such 166
167 as world knowledge and negation (De Marneffe 167
168 et al., 2008). Additionally, researchers found that 168
169 sentence structure and general syntax played a sig- 169
170 nificant role in improving entailment predictions, 170
171 thus furthering the notion that RTE datasets should 171
172 encompass the requisite linguistic diversity to com- 172
173 prehensively map entailment (Vanderwende and 173
174 Dolan, 2005; Blake, 2007). Although the subse- 174
175 quent development of the Stanford Natural Lan- 175
176 guage Inference (SNLI) (Bowman et al., 2015) 176
177 and Multi Natural Language Inference (MNL) 177
178 (Williams et al., 2017) corpora significantly im- 178
179 proved the general recognition of entailment by 179
180 incorporating human annotations of entailment 180
181 across different genres and varying real-world con- 181
182 versations, more specialized RTE datasets, such 182
183 as the Diverse Natural Language Inference Collec- 183
184 tion (Poliak et al., 2018) and the “NLI Stress Tests” 184
185 (Naik et al., 2018), advocate for the recasting of 185
186 datasets pertaining to specific linguistic phenom- 186
187 ena such as event factuality, sentiment analysis, and 187
188 numerical reasoning into RTE challenges (White 188

et al., 2017).¹ Not only does this approach extend RTE approaches to intrinsically identify diverse logical structures expressed in various linguistic styles, but it also extends the practical application of these RTE methods to various NLP applications, such as question-answering (Khot et al., 2018; Harabagiu and Hickl, 2006), text summarization (Lloret et al., 2008; Naserasadi et al., 2019), and machine translation (Padó et al., 2009).

Current RTE modelling approaches require two main steps. First, the features of premise T and hypothesis H are extracted in order to represent the statements in accordance with relevant linguistic mechanisms associated with textual entailment (Li et al., 2020). These mechanism-oriented approaches include lexical approaches that leverage part-of-speech tagging, stopword removal, and named entity recognition to represent statements through word choice (Lan and Jiang, 2018), syntactic approaches that utilize parse trees and dependency graph representations to represent statements through sentence structure (Iftene and Moruz, 2009), and semantic approaches that use semantic converters like the Universal Natural Language and cross-referenced paraphrasing to represent statements in a fuller definitive context. These feature extraction processes can be hybridized accordingly and are often represented through word embeddings (Basak et al., 2018). Second, the statements are fed into a supervised multi-class classification model which predicts whether a premise-hypothesis pair possesses *positive* (the hypothesis can be inferred to be true if the premise is true), *negative* (the hypothesis can be inferred to be false if the premise is true), or *neutral* (the hypothesis' truth is not sufficiently conditional upon the premise being true) entailment. This step has been greatly facilitated by the development of robust RTE corpora like the previously mentioned datasets.

3 Methodology

3.1 Hypothesis

We incorporate Korman's RTE approach into our methodology due to its emphasis on refined human inference, which has several important implications (Korman et al., 2018). Consider a scenario in which an individual is tasked with crafting an argument in 10 sentences. Under these circumstances, accord-

¹For an in depth analysis of existing corpora, see (Poliak, 2020).

ing to the Korman approach, the opening sentence of the individual's argument cannot directly entail the concluding sentence. This is due to the presence of intermediate premises, which are necessary to make the argument convincing and logical, thus enabling readers to justifiably infer entailment within the argument.

In each successive sentence, new information is presented, which can be either affirmative, contradictory, or neutral in relation to the preceding premises. This incremental accumulation of information ultimately leads to the conclusion asserted by the final sentence. Thus, unless there are no immediate prerequisites for logical or textual coherence in presenting a set of claims, individuals should continuously incorporate textual entailment to ensure that the overarching message is effectively delivered. This is reinforced when taking into account Korman's stipulation of identifying where the individual stands in relationship to their group G . If another individual, not part of the original group, was tasked with crafting the same 10-sentence argument in their unique writing style, they might incorporate textual entailment by structuring intermediary propositions differently or using distinct modes of textual expression compared to the first individual. Ultimately, while these two individuals may hold identical viewpoints and employ the same set of evidence to support their stance, what sets them apart are the variations in the structure of their intermediary propositions and their respective modes of expression. We suggest that in such a scenario, assessing fluctuations in textual entailment on a sentence-by-sentence basis can quantitatively illustrate the differences in argumentative approaches between these two individuals.

The intentional fluctuation of textual entailment on behalf of the author adds a new dimension to Korman's RTE approach, one which examines RTE on behalf of the communicator rather than the audience. In the aforementioned scenario, if an evaluator were to assess the textual entailment of each communicator's argument sentence-by-sentence, disparities in RTE could be attributed to the variations in the evaluator's background knowledge or linguistic proficiency, which might hinder the evaluator from making valid inferences based on the modified expression. Additionally, these differences may also stem from the distinct approaches employed by the communicators in structuring their

messages in a logically coherent manner. If the former circumstance is adequately conditioned, then the evaluator’s RTE approach can provide quantifiable insight into the logical reasoning of the communicators. This becomes more evident when the evaluator assesses multiple instances from both communicators. The *entailment ‘progressions’* derived from the evaluator’s RTE approach can effectively illustrate key patterns in the communicators’ logical approach to crafting a cohesive argument or message.

The formal definition of the entailment progressions of a given text can be expressed as follows:

$$EP_{3 \times n} = \begin{bmatrix} c_1 & c_2 & \cdots & c_n \\ p_1 & p_2 & \cdots & p_n \\ n_1 & n_2 & \cdots & n_n \end{bmatrix}$$

where EP is an entailment progression matrix composed of c , p , n row vectors representing the contradiction, positive, and neutral entailment probabilities between two sentences in a given text. To compute these values at a given point in a text, we introduce the following equation:

$$EP_{c_i, p_i, n_i} = E(s_{i-1}, s_i)$$

where E represents the entailment model used for calculating entailment between two sentences, and s represents a sentence at a given point in the text.

Drawing from our analysis of the existing RTE literature, we propose the following. Given two texts composed of an equal number of sentences, denoted as T_1 and T_2 , which are equal in length and are composed of a series of similarly presented statements that serve to further the same logical premise, and an evaluator E , an individual tasked with recognizing the textual entailment on a statement-by-statement basis for T_1 and T_2 , if E identifies sufficient differences in the entailment progressions of T_1 and T_2 , then T_1 and T_2 can be distinguished based on the distinct logical approaches employed by their respective communicators.

3.2 Experimental Design

To ensure that our hypothesis is satisfied, we design an experimental setup that effectively accounts for potential confounding limitations that may arise during the evaluator’s analysis.

First, both C_1 and C_2 must employ similar linguistic mechanisms for presenting their premises. This is to avoid potential deterrents caused by

limited semantic or syntactical knowledge, which could hinder the evaluator’s ability to accurately assess the truth value of statements within C_1 and C_2 , a necessary prerequisite for RTE.

Second, both C_1 and C_1 must “*further the same logical premise*” by pertaining to the highly similar, if not identical, domains in which the evaluator possesses a sufficient and equal understanding. This is to ensure that the evaluator possesses the necessary background knowledge to proficiently implement their RTE approach with “*justifiable inference*”.

Third, both C_1 and C_2 must be equal in length, or possess an equal amount of statements. This final condition is in direct reference to the previously mentioned issue of inferential distance, in which the difficulty of reasoning from one statement to another in a piece of text is associated to the number of intermediary propositions required to effectively connect the statements. If C_1 and C_2 both advance the same logical premise, but C_1 is significantly longer than C_2 , with both covering an equal total inferential distance from their initial to final statements, then, on average, the statement-by-statement fluctuations in textual entailment for C_1 would be lesser than those of C_2 , simply by dividing the total inferential distance by the number of statements. While this explanation may not fully account for cases where C_1 offers a richer explanation of its relevant topic (and subsequently exhibits a higher degree of variability in the textual entailment it employs), it does hold true in a qualitative sense. If a communicator was tasked with crafting a 5-sentence argument on a subject that typically requires 15 statements for a cohesive exposition, the condensed length will force the communicator to emphasize specific points with greater urgency. In turn, this could lead to larger logical jumps and thus greater deviations in RTE.

When controlling for these conditions, we design an experimental setting in which the evaluator is capable of qualitatively and quantitatively distinguishing between the logical structures employed in C_1 and C_2 . Please note that this comparison can be made not only between texts from the same communicator but also between texts from different communicators, and the conclusions drawn from each approach may carry different implications. If both C_1 and C_2 are authored by the same communicator, analysing the textual entailment progressions of both texts can serve as a stress test. This stress test helps evaluate the robustness of the communi-

377 cator’s logical approach relative to minor topical
 378 differences that still adhere to the second condition.
 379 However, when C_1 and C_2 originate from different
 380 communicators, evaluating the textual entailment
 381 progressions of both texts allows for a comparison
 382 of intention in logical approaches between these
 383 communicators. By combining these “*intracom-*
 384 *municator*” (i.e., evaluating the logical approaches
 385 pertaining to a single author) and “*intercommuni-*
 386 *cator*” (i.e., comparing logical approaches across
 387 multiple authors) approaches, we can create a third
 388 type of comparison. This approach enables us to
 389 evaluate whether the differences in intended logical
 390 approaches between communicators remain distin-
 391 guishable and robust across various examples.

392 3.3 Data and Models

393 The data used to evaluate our hypothesis is sourced
 394 from mutual fund evaluations written in English
 395 (hereafter *narratives*) provided by a leading finan-
 396 cial analysis firm. These narratives supplement
 397 the firm’s assessments of mutual fund performance
 398 metrics by elaborating on key qualitative aspects
 399 of the fund’s performance.

400 The narratives are primarily categorized based
 401 on authorship or the specific method used for
 402 their creation. The narratives for the top-
 403 performing 25% of the examined funds are au-
 404 thored by human analysts (referred to as *Analyst*
 405 *narratives*), while the narratives for the re-
 406 maining 75% are generated using a proprietary al-
 407 gorithm called Smart Text (referred to as *Quant*
 408 *narratives*). Smart Text was first launched in
 409 the spring of 2021 and is generated through a de-
 410 terministic rules-based process. The firm employs
 411 a rule-based approach to group funds into various
 412 mental models, where each mental model has an
 413 associated template structure that helps to assign
 414 text branches to. These branches contain embedded
 415 data points (such as fund name and expense ratio)
 416 and also a synonym bank of words to ensure text
 417 variation. It is important to note that this process is
 418 entirely rules-based and doesn’t utilize any recent
 419 generative AI techniques.²

420 While *Analyst narratives* are cate-
 421 gorized by the authoring analyst, *Quant*
 422 *narratives* are categorized by subject-specific
 423 templates used to evaluate a fund in relation to a
 424 specific area in investment management research

²The reason for employing a rule-based model is the firm’s focus on generating precise and reliable narratives.

425 that is most relevant to its performance. Our
 426 dataset contains 26 different templates, each cov-
 427 ering a specific aspect of a fund’s performance
 428 relative to its core characteristics. For example, the
 429 *Active Allocation* template is structured to
 430 evaluate funds that are both actively managed by
 431 a portfolio manager and are diversified across dif-
 432 ferent asset classes, and is primarily focused on
 433 the fund’s active allocation strategy rather than
 434 other aspects that the fund may possess. To ex-
 435 amine the relationship between authorship and log-
 436 ical approach, we divided the sets by analyst (for
 437 *Analyst narratives*) and by template type
 438 (for *Quant narratives*).

439 Several constraints were applied when im-
 440 plementing the proposed experimental approach.
 441 Firstly, as cosine similarity served as the metric
 442 used for evaluating the similarity of entailment pro-
 443 gressions within evaluated groups, and given that
 444 entailment progressions often varied in length due
 445 to differences in the number of sentences across
 446 sets of *Quant* and *Analyst narratives*, we
 447 grouped each set into subsets based on the number
 448 of sentences. Additionally, since cosine similarity
 449 is inflated in lower dimensional settings, we only
 450 considered subsets with 10 or more sentences in
 451 our final analysis. Furthermore, to ensure the ac-
 452 curacy of our analysis and prevent potential issues
 453 that could arise due to a low sample size, we only
 454 considered subsets with 10 or more narratives. The
 455 final cosine similarity scores were computed as a
 456 weighted average of the cosine similarity scores
 457 within each subset, with weights determined by the
 458 number of narratives within the subset belonging
 459 to the respective set.

460 Table 1 provides an overview of the dataset used
 461 in our experiments categorized by authorship (i.e.,
 462 whether the narrative originates from an analyst or
 463 the Smart Text model) as well as the number of
 464 categories within each type of narrative.

NARRATIVE TYPE	# OF AUTHOR CATEGORIES	TOTAL
Analyst	5	1000
Quant	26	4000

Table 1: Total count of narratives and categories within each narrative type.

465 To calculate the textual entailment on a sentence-
 466 by-sentence basis for the narratives in our dataset,
 467 we rely on **RoBERTa-base** (Liu et al., 2019), an
 468 optimized BERT-like (Devlin et al., 2019) encoder

Transformer. We specifically leverage a version of this model adapted to the domain of RTE by fine-tuning on datasets designed for Natural Language Inference (Reimers and Gurevych, 2019). For performing the experiments, we relied on the Hugging-Face transformers library (Wolf et al., 2020).

4 Results and Discussion

Table 2 presents the cosine similarity scores for the entailment progressions of the `Quant` and `Analyst` narratives. These results highlight that `Quant` narratives exhibit, on average, higher cosine similarity scores within sets of entailment progressions compared to `Analyst` narratives. These results align with our “*intra-communicator*” experimental design, where higher cosine similarity scores indicate greater robustness in the stylized logical approach characterizing the narrative set. Since the `Quant` narratives are generated through the Smart Text algorithm, this consistency within the narratives can be attributed to the (more) formulaic communication style embedded into the narratives’ delivery method. In contrast, `Analyst` narratives are prone to display greater variability within the narrative sets, reflecting the less structured approach often taken by human analysts when crafting their narratives.

Table 3 presents the cosine similarity scores for the entailment progressions of the `Quant` and `Analyst` narratives between different author groups, an example of the previously noted “*intercommunicator*” approach. While we anticipate the cosine similarity across different authors to be lower than the “*intracommunicator*” comparison, our analysis suggests that this is dependent upon the specific contextual factors that shape the writing styles of authors within a group. For example, the cosine similarities of different analyst pairs do not differ from the results presented in Table 2. This can be attributed to the fact that all the evaluated analysts are expected to adhere to the writing standards set by the financial services firm employing them. This trend is not upheld when examining the cosine similarities between different template pairs, where we observe not only a significant decrease in the cosine similarities as outlined in Table 2, but also a greater variance, ranging from -0.148 to 0.370. This aligns with the manner in which the templates are generated (i.e., on the basis of specific rules associated with the subject matter they address). Therefore, a higher cosine similarity is

observed when evaluating these templates in comparison to those that utilize a very similar logical structure. The most notable instance of this is phenomenon is evident when comparing a template to itself (cf. Table 2)).

The results presented in Tables 2 and 3 suggest that the entailment progressions of a given author’s work are not only influenced by the authors themselves, but also by the stylistic constraints imposed upon them. The first notable constraint observed in our analysis pertains to *structure*. Both the `Quant` and `Analyst` narratives adhere to the structural constraint in terms of length (i.e., the analysts and the Smart Text algorithm are required to condense their analysis of a given mutual fund into an informative and concise format). Were the narratives written without this constraint, and thus allowing for more sentences to be used for covering a given fund, then their entailment progressions would change to accommodate more intermediary propositions. This would reduce the inferential distance between the initial and final claims (Korman et al., 2018). The second notable constraint pertains to the subject matter of the text. A broader, more complex subject matter can create variability in the logical approaches employed to effectively support the claim at hand. This can be observed in both the `Quant` and `Analyst` narratives; while the former specializes in narrower subject matters (e.g., Active Equity), the latter focuses on multiple aspects of a fund that can be used to evaluate performance. While Table 2 shows how this subject constraint renders the cosine similarity within narrative groups higher for `Quant` narratives than for `Analyst` narratives, Table 3 highlights how the cosine similarity between `Quant` narratives groups experiences a more pronounced decrease compared to `Analyst` narratives groups by virtue of differing subject matters altogether.

This becomes even more evident when we visualize the entailment progressions associated with specific sets of `Quant` and `Analyst` narratives (cf. Figure 1). While a high cosine similarity generally indicates a similarity in trends between two vectors, this distinction is particularly pronounced when assessing the points at which significant fluctuations (or lack thereof) occur in the entailment progressions of the narrative set.

Depending on the template under examination,

ANALYST	COSINE SIMILARITY	TEMPLATE	COSINE SIMILARITY
Analyst A	0.285	Active Allocation	0.601
Analyst B	0.237	Passive Allocation	0.584
Analyst C	0.234	Active Equity	0.579
Analyst D	0.230	Long Term	0.554
Analyst E	0.201	Active Fixed Income	0.530

Table 2: Cosine similarity scores of entailment progressions within the evaluated narrative sets. A higher score indicates a greater similarity between entailment progressions and suggests a closer alignment in the underlying logical structure employed by the respective authors.

ANALYST PAIRS	COSINE SIMILARITY	TEMPLATE PAIRS	COSINE SIMILARITY
Analyst A, Analyst B	0.213	Active Fixed Income, Long Term	0.370
Analyst A, Analyst C	0.237	Active Fixed Income, Total Risk	0.112
Analyst A, Analyst D	0.189	Active Fixed Income, Passive Allocation	-0.148
Analyst B, Analyst C	0.244	Passive Allocation, Active Equity	0.192
Analyst B, Analyst D	0.188	Passive Allocation, Total Risk	0.199

Table 3: Cosine similarity scores of entailment progressions between different evaluated narrative sets. Analyst and Quant narratives are directly compared in pairs by their respective analyst or template.

we observe that the Quant narratives exhibit fluctuations in entailment at key points. Particularly, in the case of Active Allocation Quant narratives, we can observe spikes of contradictory entailment occurring around the 40% and 90% marks. In contrast, Analyst narratives are less uniform due to their non-formulaic nature – human analysts may not purposefully introduce contradictory entailment at consistent points across all their narratives. Additionally, while entailment progressions in Quant narratives exhibit a clearer dominance of neutral entailment (with occasional spikes in positive or negative entailment), Analyst narratives integrate both positive and negative entailment throughout the narrative, a distinction that can be attributed to the difference in communicative quality between the method used for generating the Quant narratives and human authorship.

We can link our analysis results to our initial hypothesis in three key ways. First, we identify the similarity in entailment progressions for template-based sets of Quant narratives as a result of the rigidity in the logical approaches employed within templates. If one were to read such narratives within a specific template, they would be able to recognize the similar logical structure and shared logical processes used for their generation. Second, we note the lower similarity in entailment progres-

sions for the Analyst narratives sets, indicating a less rigid logical structure in their crafting. The variability in logical approaches used by the analysts when rating mutual funds can make it challenging for a reader to identify the analyst solely based on logical flow. However, this does not preclude the reader from identifying the analyst on the basis of word choice, tone, and other stylistic indicators that are separate from RTE. Third, the varying levels of similarity between sets of Quant and Analyst narratives in their entailment progression subtly highlight the distinction between human and automatically-generated communication. Quant narratives tend to follow a more structured and explicit logical structure compared to Analyst narratives. Consequently, if a reader were to read a combined set of Quant narratives (from a specific template) and Analyst narratives (written by a specific author), they could discern whether a narrative was automatically generated or crafted by a human analyst based on its resemblance to narratives they previously encountered. While this may be more implicit than the previous two points, it underlies the sentiment expressed by the firm’s clients, who find Quant narratives to be “too robotic” and lacking the qualitative aspects of an analyst’s narrative. In all three situations, our hypothesis holds, confirming the benefits of leveraging entailment progressions.

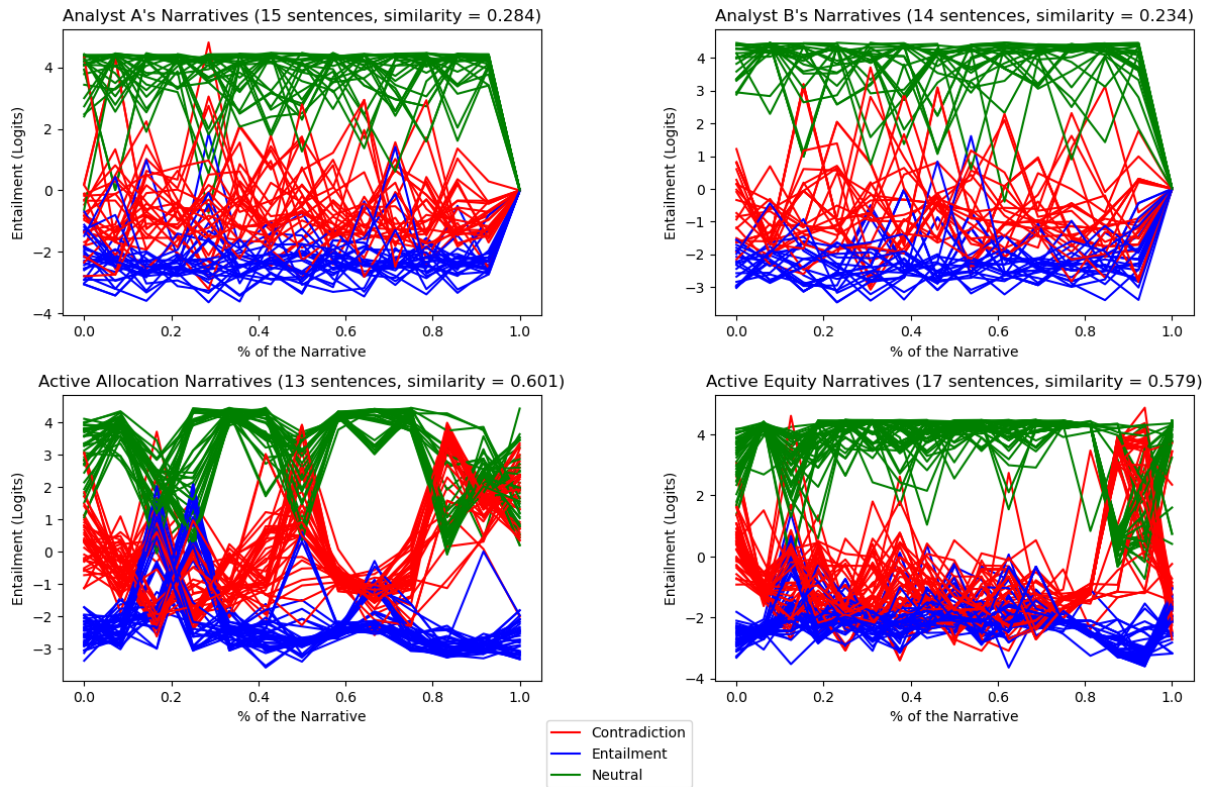


Figure 1: Line plots displaying the entailment progressions of narratives within selected sets. Entailment is measured in logits calculated by the RoBERTa model used to calculate the entailment scores for the progression.

5 Conclusion

In this paper we introduce entailment progressions, a framework which serves in identifying the logical flow of a text and highlights the way in which different authors choose to integrate nuance and affirmation on a sentence-by-sentence basis. This deviation from traditional RTE literature seeks to characterize textual entailment similarly to author style, which is often linked to lexical and semantic choices, as opposed to being seen as a purely objective or impersonal linguistic phenomenon. As demonstrated through the difference in cosine similarity between the `Quant` and `Analyst narrative` entailment progressions, the author’s stylistic rigidity influences the robustness of the underlying logical structures. These findings suggest that entailment progressions can potentially distinguish between human and “non-human” logical approaches. In future work, we will explore the logical approaches taken by LLMs and assess to what degree they align with the humanistic logical approaches. This exploration can play a crucial role in identifying whether these models possess an inherent logical structure to adhere to, which, in turn, can potentially contribute

to the larger area of LLM interpretability. Second, we will explore whether entailment progressions can serve as a benchmark for defining the similarity between two logical approaches, extending the analysis presented in Table 3 to encompass other authorship styles. If two sets of texts’ entailment progressions, both of which are conditioned in accordance with our defined methodology, were compared on the basis of cosine similarity, a high cosine similarity could indicate that both texts adhere to similar logical structures and that the authors of the texts have similarly integrated logic into how they have devised their overall texts. This concept can apply to both human and non-human approaches, where comparisons between human authors, human and model-based authors, and model-based authors can be analysed using entailment progressions as a heuristic to assess whether similar logical processes are at play.

Ethics Statement

The data that was used for conducting the experiments was provided by a leading financial analysis firm and is not publicly available. Both the `Quant` and `Analyst narrative` narratives were col-

679	lected internally by members of the firm’s team,	729
680	and are exclusively available to clients who have	730
681	subscribed to service offerings that grant access	731
682	to these narratives. Although the firm has given	732
683	explicit approval to our methodology and use of	733
684	their data, the proprietary nature of this information	734
685	warrants selective discretion when describing and	735
686	disseminating this data.	736
687	Limitations	737
688	Entailment progressions serve to identify the logi-	738
689	cal approach employed by a given author, to which	739
690	we characterize it as an aspect of the author’s style	740
691	that can be leveraged for tasks that involve author	741
692	classification or author style transfer. Since our	742
693	analysis extends to assessing the internal logical	743
694	structures employed by non-human approaches in	
695	relation to typical human logic approaches, the	
696	former task can aid tools seeking to distinguish	
697	between model-generated and human-generated	
698	texts when necessary, a problem that is newly aris-	
699	ing in areas concerning intellectual property and	
700	fraud detection (Yu et al., 2023; Májovský et al.,	
701	2023). Although entailment progressions can help	
702	identify stylistic differences between model and	
703	human-generated outputs, they can simultaneously	
704	improve these models in their closeness in inferen-	
705	tial capability to humans, potentially rendering the	
706	task of differentiating between human and model	
707	outputs more difficult.	
708	In this work, we acknowledge a number of limi-	
709	tations and discuss their implications.	
710	(1) Our analysis leverages cosine similarity to com-	
711	pare entailment progressions, but this is limited by	
712	specific metric requirements. First, cosine simi-	
713	larity requires that the vectors being compared be	
714	of the same length, thus limiting a comprehensive	
715	analysis. Although we address this issue within our	
716	methodology, this limitation significantly hinders	
717	our ability to perform an in-depth analysis of entail-	
718	ment progressions that fall outside this length con-	
719	straint. Second, cosine similarity is inflated in low-	
720	dimensional settings, which can hinder the compar-	
721	ison of entailment progressions derived from	
722	shorter texts. These limitations can obscure entail-	
723	ment progression analysis to the extent that it	
724	artificially inflates the similarity scores between	
725	entailment progressions in every instance.	
726	(2) As entailment progressions are generated using	
727	a pre-existing language model from the literature,	
728	the quality of these entailment progressions is di-	
	rectly dependent on the performance of the model.	
	In scenarios like ours, where the true entailment is	
	unknown, assessing the performance of the model	
	can be challenging. Additionally, it is important to	
	note that the model is trained on corpora containing	
	general language premise-hypothesis pairs, which	
	can limit its performance in more specialized do-	
	main, such as investment management research.	
	Given that our methodology relies on the model’s	
	capability to accurately infer textual entailment,	
	this requirement may not always be satisfied when	
	applied to areas outside the scope of general lan-	
	guage understanding. For example, a model might	
	construe (4) and (5) as negatives, when in fact they	
	are positives.	
	(4) The expensive ratio for this fund is low.	744
	(5) Manager turnover for this firm is below aver-	745
	age.	746
	This can be mitigated through ensemble modelling,	747
	where multiple entailment models are deployed,	748
	and the corresponding agreement/disagreement can	749
	be leveraged to determine the final entailment prob-	750
	abilities at a given point in the progression.	751
	(3) Given that entailment progressions adhere to the	752
	Markov property, wherein entailment between two	753
	statements in an entailment progression is sequen-	754
	tial, entailment progression approaches should evalu-	755
	ate whether sufficient context is captured by exam-	756
	ining only the immediately preceding sentence,	757
	rather than a larger set of previous sentences. Exam-	758
	ining only the preceding sentence when calculating	759
	entailment may lead to an inflation of contradic-	760
	tion and neutral probability scores, as it overlooks	761
	prior sentences which may sufficiently entail the	762
	sentence under evaluation. This limitation is more	763
	procedural than the other limitations, as entailment	764
	progression generation can be adapted to account	765
	for different Markov assumptions.	766
	References	767
	Marilisa Amoia. 2009. Linguistic-based computational	768
	treatment of textual entailment recognition.	769
	Rohini Basak, Sudip Kumar Naskar, and Alexander Gel-	770
	bukh. 2018. A simple hybrid approach to recognizing	771
	textual entailment. <i>Journal of Intelligent & Fuzzy</i>	772
	<i>Systems</i> , 34(5):2873–2885.	773
	Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo	774
	Giampiccolo. 2009. The fifth pascal recognizing	775
	textual entailment challenge. <i>TAC</i> , 7:8.	776

777	Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. 2011. The seventh pascal recognizing textual entailment challenge. In <i>TAC</i> .	830
778		831
779		832
780	Catherine Blake. 2007. The role of sentence structure in recognizing textual entailment. In <i>Proceedings of the ACL-PASCAL Workshop on Textual Entailment and Paraphrasing</i> , pages 101–106.	833
781		834
782		835
783		
784	Johan Bos and Katja Markert. 2005. Recognising textual entailment with logical inference. In <i>Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing</i> , pages 628–635.	
785		
786		
787		
788		
789	Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. <i>arXiv preprint arXiv:1508.05326</i> .	836
790		837
791		
792		
793	David Braun. 2001. Indexicals.	838
794		839
795	Paul F Christiano, Jan Leike, Tom Brown, Miljan Martić, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. <i>Advances in neural information processing systems</i> , 30.	840
796		841
797		
798	Ido Dagan, Bill Dolan, Bernardo Magnini, and Dan Roth. 2010. Recognizing textual entailment: Rational, evaluation and approaches—erratum. <i>Natural Language Engineering</i> , 16(1):105–105.	842
799		843
800		844
801		845
802	Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In <i>Machine learning challenges workshop</i> , pages 177–190. Springer.	846
803		847
804		
805		
806	Marie-Catherine De Marneffe, Anna N Rafferty, and Christopher D Manning. 2008. Finding contradictions in text. In <i>Proceedings of acl-08: Hlt</i> , pages 1039–1047.	848
807		849
808		850
809		851
810	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	852
811		853
812		854
813		855
814		856
815		857
816		858
817		859
818		
819	Richard Feldman. 2003. Epistemology. <i>Tijdschrift Voor Filosofie</i> , 68(2).	860
820		861
821	Danilo Giampiccolo, Hoa Trang Dang, Bernardo Magnini, Ido Dagan, Elena Cabrio, and Bill Dolan. 2008. The fourth pascal recognizing textual entailment challenge. In <i>TAC</i> .	862
822		863
823		
824		
825	Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and William B Dolan. 2007. The third pascal recognizing textual entailment challenge. In <i>Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing</i> , pages 1–9.	864
826		865
827		866
828		867
829		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
	Sanda Harabagiu and Andrew Hickl. 2006. Methods for using textual entailment in open-domain question answering. In <i>Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics</i> , pages 905–912.	
	Adrian Iftene and Mihai Alex Moruz. 2009. Uaic participation at rte5. In <i>TAC</i> .	
	Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. Scitail: A textual entailment dataset from science question answering. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 32.	
	Daniel Z Korman, Eric Mack, Jacob Jett, and Allen H Renear. 2018. Defining textual entailment. <i>Journal of the Association for Information Science and Technology</i> , 69(6):763–772.	
	Yunshi Lan and Jing Jiang. 2018. Embedding wordnet knowledge for textual entailment. <i>ACL</i> .	
	Peiguang Li, Hongfeng Yu, Wenkai Zhang, Guangluan Xu, and Xian Sun. 2020. Sa-nli: A supervised attention based framework for natural language inference. <i>Neurocomputing</i> , 407:72–82.	
	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach . <i>arXiv:1907.11692 [cs]</i> .	
	Elena Lloret, Oscar Ferrández, Rafael Munoz, and Manuel Palomar. 2008. A text summarization approach under the influence of textual entailment. In <i>NLPCS</i> , pages 22–31.	
	James MacGlashan, Mark K Ho, Robert Loftin, Bei Peng, Guan Wang, David L Roberts, Matthew E Taylor, and Michael L Littman. 2017. Interactive learning from policy-dependent human feedback. In <i>International conference on machine learning</i> , pages 2285–2294. PMLR.	
	Martin Májovský, Martin Černý, Matěj Kasal, Martin Komarc, and David Netuka. 2023. Artificial intelligence can generate fraudulent but authentic-looking scientific medical articles: Pandora’s box has been opened. <i>Journal of Medical Internet Research</i> , 25:e46924.	
	Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. <i>arXiv preprint arXiv:1806.00692</i> .	
	Ali Naserasadi, Hamid Khosravi, and Faramarz Sadeghi. 2019. Extractive multi-document summarization based on textual entailment and sentence compression via knapsack problem. <i>Natural Language Engineering</i> , 25(1):121–146.	

885	Sebastian Padó, Michel Galley, Dan Jurafsky, and	Lhoest, and Alexander Rush. 2020. Transformers:	941
886	Christopher D Manning. 2009. Robust machine trans-	State-of-the-Art Natural Language Processing . In	942
887	lation evaluation with entailment features. In <i>Pro-</i>	<i>ceedings of the 2020 Conference on Empirical</i>	943
888	<i>ceedings of the Joint Conference of the 47th Annual</i>	<i>Methods in Natural Language Processing: System</i>	944
889	<i>Meeting of the ACL and the 4th International Joint</i>	<i>Demonstrations</i> , pages 38–45, Online. Association	945
890	<i>Conference on Natural Language Processing of the</i>	for Computational Linguistics.	946
891	<i>AFNLP</i> , pages 297–305.		
892	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Zhiyuan Yu, Yuhao Wu, Ning Zhang, Chenguang Wang,	947
893	Jing Zhu. 2002. Bleu: a method for automatic evalu-	Yevgeniy Vorobeychik, and Chaowei Xiao. 2023.	948
894	ation of machine translation. In <i>Proceedings of the</i>	Codeiprompt: Intellectual property infringement	949
895	<i>40th annual meeting of the Association for Computa-</i>	assessment of code language models.	950
896	<i>tional Linguistics</i> , pages 311–318.		
897	Adam Poliak. 2020. A survey on recognizing textual	Annie Zaenen, Lauri Karttunen, and Richard Crouch.	951
898	entailment as an NLP evaluation . In <i>Proceedings of</i>	2005. Local textual inference: can it be defined or	952
899	<i>the First Workshop on Evaluation and Comparison</i>	circumscribed? In <i>Proceedings of the ACL workshop</i>	953
900	<i>of NLP Systems</i> , pages 92–109, Online. Association	<i>on empirical modeling of semantic equivalence and</i>	954
901	for Computational Linguistics.	<i>entailment</i> , pages 31–36.	955
902	Adam Poliak, Aparajita Haldar, Rachel Rudinger, J Ed-		
903	ward Hu, Ellie Pavlick, Aaron Steven White, and Ben-		
904	jamin Van Durme. 2018. Collecting diverse natural		
905	language inference problems for sentence representa-		
906	tion evaluation. <i>arXiv preprint arXiv:1804.08207</i> .		
907	Nils Reimers and Iryna Gurevych. 2019. Sentence-		
908	BERT: Sentence embeddings using Siamese BERT-		
909	networks . In <i>Proceedings of the 2019 Conference on</i>		
910	<i>Empirical Methods in Natural Language Processing</i>		
911	<i>and the 9th International Joint Conference on Natu-</i>		
912	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages		
913	3982–3992, Hong Kong, China. Association for Com-		
914	putational Linguistics.		
915	Lucy Vanderwende and William B Dolan. 2005. What		
916	syntax can contribute in the entailment task. In <i>Ma-</i>		
917	<i>chine Learning Challenges Workshop</i> , pages 205–		
918	216. Springer.		
919	Alex Wang, Amanpreet Singh, Julian Michael, Felix		
920	Hill, Omer Levy, and Samuel R Bowman. 2018.		
921	Glue: A multi-task benchmark and analysis platform		
922	for natural language understanding. <i>arXiv preprint</i>		
923	<i>arXiv:1804.07461</i> .		
924	Aaron Steven White, Pushpendre Rastogi, Kevin Duh,		
925	and Benjamin Van Durme. 2017. Inference is every-		
926	thing: Recasting semantic resources into a unified		
927	evaluation framework. In <i>Proceedings of the Eighth</i>		
928	<i>International Joint Conference on Natural Language</i>		
929	<i>Processing (Volume 1: Long Papers)</i> , pages 996–		
930	1005.		
931	Adina Williams, Nikita Nangia, and Samuel R Bow-		
932	man. 2017. A broad-coverage challenge corpus for		
933	sentence understanding through inference. <i>arXiv</i>		
934	<i>preprint arXiv:1704.05426</i> .		
935	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien		
936	Chaumond, Clement Delangue, Anthony Moi, Pier-		
937	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,		
938	Joe Davison, Sam Shleifer, Patrick von Platen, Clara		
939	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven		
940	Le Scao, Sylvain Gugger, Mariama Drame, Quentin		