

# PROPSEGMENT: A Large-Scale Corpus for Proposition-Level Segmentation and Entailment Recognition

Sihao Chen<sup>\*1,2</sup>, Senaka Buthpitiya<sup>1</sup>, Alex Fabrikant<sup>1</sup>, Dan Roth<sup>2</sup>, Tal Schuster<sup>1</sup>

<sup>1</sup>Google Research

<sup>2</sup>University of Pennsylvania

{sihaoc, senaka, fabrikant, talschuster}@google.com, {sihaoc, danroth}@cis.upenn.edu

## Abstract

The widely studied task of Natural Language Inference (NLI) requires a system to recognize whether one piece of text is textually entailed by another, i.e. whether the *entirety* of its meaning can be inferred from the other. In current NLI dataset and models, textual entailment relations are typically defined on the sentence- or paragraph-level. However, even a simple sentence often contains multiple *propositions*, i.e. distinct units of *meaning* conveyed by the sentence. As these propositions can carry different truth values in the context of a given premise, we argue for the need to segment a sentence into propositions, and individually infer their textual entailment relations.

We propose PROPSEGMENT, a corpus of over 35K propositions annotated by expert human raters. Our dataset structure resembles the tasks of (1) segmenting sentences within a document to the set of propositions, and (2) classifying the entailment relation of each proposition with respect to a different yet topically-aligned document, i.e. documents describing the same event or entity. We establish strong baselines for the segmentation and entailment tasks. Through case studies on summary hallucination detection and document-level NLI, we demonstrate that our conceptual framework is potentially useful for understanding and explaining the compositionality of NLI labels.

## 1 Introduction

Natural Language Inference (NLI), or Recognizing Textual Entailment (RTE), is the task of determining whether the meaning of one text expression can be inferred from another (Dagan and Glickman, 2004). Given two pieces of text ( $P, H$ ), we say that the premise  $P$  *entails* the hypothesis  $H$  if the *entirety* of  $H$ ’s meaning can be most likely inferred true after a human reads  $P$ . If some units of meaning in  $H$  is contradicted or cannot be determined by

### Premise Document

Andrew Warhola, known as Andy Warhol, is an American artist born August 6, 1928 in Pittsburgh, Pennsylvania and died February 22, 1987 in New York. He is one of the main representatives of pop art. Warhol is known the world over for his work as a painter, music producer, author, avant-garde films... (7 more sentences omitted)

### Hypothesis Sentence

(from another document of the same topic)

... The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art. ...

| Propositions   | Entailment Label |
|--|------------------|
| The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art. | Neutral          |
| The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art. | Entailment       |
| The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art. | Neutral          |

Table 1: An example instance from our proposed PROPSEGMENT dataset with propositions (marked as token subsets highlighted in blue) and their entailment labels.

by  $P$ , we describe the relation between the two as *contradiction* or *neutral* (de Marneffe et al., 2008) respectively. This fundamentally challenging natural language understanding task provides a general interface for semantic inference and comparison across different sources of textual information.

In reality, most naturally occurring text expressions are composed by a variable number of *propositions*, i.e. distinct units of meaning conveyed by the piece of text. Consider the sentence shown in Table 1: “The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art.”. Despite the sentence being relatively compact, it still contains (at least) three propositions, as listed in Table 1. While the entire hypothesis would be classified as *neutral* or *not-entailed* to the premise, one of its proposi-

\* Work done as an intern at Google

tions “*Andy Warhol’s hometown is in Pittsburgh, Pennsylvania*” is in fact entailed by the premise, while the premise provides no support for the other two propositions. This phenomenon, namely *partial entailment* (Levy et al., 2013), is a blind spot for existing sentence- or paragraph-level NLI formulations. When a hypothesis is *compositional*, NLI labels coarsely defined on the sentence/paragraph-level cannot express the difference between partial entailment from the non-entailment cases.

This work argues for the need to study and model textual entailment relations on the level of *propositions*. As NLI tasks and applications typically involve different genre of text with variable length and number of propositions (Yin et al., 2021), decomposing textual entailment relation to the propositional level provides a more fine-grained yet accurate description of textual entailment relation between two arbitrary text expressions. Modeling *propositional textual entailment* provides a more unified inference format across NLI tasks, and would potentially improve the generalization capabilities of NLI models, e.g. with respect to the variability in input lengths (Schuster et al., 2022).

To facilitate the study along this line, we propose PROPSEGMENT, a large-scale, multi-domain corpus with over 35K human-annotated propositions<sup>1</sup>. We define the tasks of proposition-level segmentation and entailment, where given a hypothesis sentence and a premise paragraph, a system is expected to segment the sentence into the set of all propositions within the sentence, and recognize whether each of the propositions can be directly inferred from the premise.

Interestingly, we observe in our study that existing syntactic or semantic-based notions of proposition (Baker et al., 1998; Kingsbury and Palmer, 2002; Meyers et al., 2004) often fail to account for the complete set of propositions in a sentence, partly due to the fact that propositions do not necessarily correspond to direct predicate-argument relations in the sentence (§ 2). For this reason, we adopt a more flexible and unified way of representing a proposition as a *subset of tokens from the input sentence*, without explicitly annotating the semantic role or predicate-argument structure within the proposition, as illustrated in Table 1. We offer a more detailed discussion on the motivation and design desiderata in § 2.

<sup>1</sup>The dataset is available at <https://github.com/google-research-datasets/propsegment>

We construct PROPSEGMENT by sampling clusters of topically-aligned documents, i.e. documents focusing on the same entity or event, from WIKIPEDIA (Schuster et al., 2022) and news domains (Gu et al., 2020). We train and instruct expert annotators to identify all propositions exhaustively in a document, and label the textual entailment relation of each proposition with respect to another document in the cluster, viewed as premise.

We discuss the modeling challenges, and establish strong baselines for the segmentation and entailment tasks. We demonstrate the utility of our dataset and models through downstream use case studies on summary hallucination detection (Maynez et al., 2020), and DocNLI (Yin et al., 2021), through which we show that recognizing and composing entailment relations at the proposition-level could provide fine-grained characterization and explanation for NLI-like tasks, especially with long and compositional hypotheses.

In summary, the main contributions in our paper include: (1) Motivating the need to recognize textual entailment relation on proposition level; (2) Introducing the first large-scale dataset for studying proposition-level segmentation and entailment recognition; and (3) Leveraging PROPSEGMENT to train Seq2Seq models as strong baselines for the tasks, and demonstrating their utility in document-level NLI and hallucination detection tasks.

## 2 Motivations & Design Challenges

Our study concerns the challenges of applying NLI/RTE formulations and systems to downstream applications in *real-world settings*. As textual entailment generally captures the relation between the meanings and implications of two pieces of text, one natural type of downstream use cases for NLI systems is to identify alignments and discrepancies between the semantic content presented in different documents/sources (Kryscinski et al., 2020; Zhang et al., 2020; Schuster et al., 2021, 2022).

For instance, our study is motivated by the task of comparing the information presented in two topically related documents, e.g. news documents covering the same event (Gu et al., 2020), or Wikipedia pages from different languages for similar entities (Schuster et al., 2022). As existing NLI datasets typically define the textual entailment relation at the sentence or paragraph level, NLI systems trained on such resources can only recognize whether or not the entirety of a

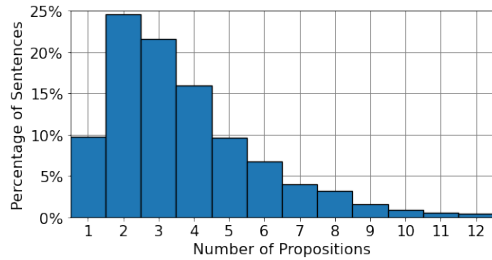


Figure 1: Distribution of proposition counts within sentences sampled from Wikipedia and news, as estimated in our PROPSEGMENT dataset.

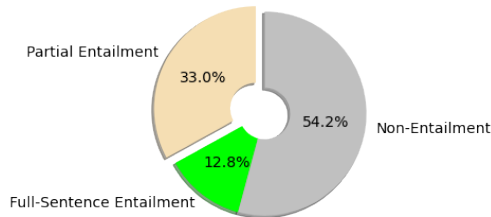


Figure 2: The percentage of sentences with *partial entailment* relation to another topically-related document (from Wikipedia or news) in the PROPSEGMENT dataset. Typically, NLI/RTE datasets do not distinguish partial entailment from the non-entailment categories.

hypothesis sentence/paragraph is entailed by a premise. However, we estimate that around 90% of the sentences from the two domains contain more than one propositions (Figure 1). In the presence of multiple propositions, *partial entailment* (Levy et al., 2013) describes the phenomenon where only a subset of propositions in the hypothesis is entailed by the premise.

**Partial entailment is 3× more common than full-sentence entailment.** In our corpus, we observe that, given two topically related documents from news or Wikipedia, 46% of sentences in one document have at least some information supported by the other document (Figure 2). But 74% of that slice is sentences that are *partially entailed*, with only some propositions supported by the other document. In this sense, a traditional system designed for sentence-level NLI labels will thus only detect a quarter of sentences that have meaningful entailment relations. In applications that seek a full understanding of cross-document semantic links, there is thus 4× headroom, a significant blind spot for sentence-level NLI models.

As we observe that most natural sentences are compositional, i.e. contain more than one proposition, we argue for the need to decompose and recognize textual entailment relation at the more granular level of propositions. In other words, in-

stead of assessing the entire hypothesis as one unit in the context of a premise, we propose to evaluate the truth value of each proposition individually, and aggregate for the truth value of the hypothesis.

### Current predicate-argument based methods often fail to extract all propositions in a sentence.

The linguistic notion of a proposition refers to a single, contextualized unit of meaning conveyed in a sentence. In the NLP community, propositions are usually represented by the predicate-argument structure of a sentence. For example, corpora resources such as FrameNet (Baker et al., 1998), PropBank (Palmer et al., 2005), NomBank (Meyers et al., 2004), among others, represent a proposition by a single predicate (verbal, nominal, etc.) along with its arguments with respect to their proto-thematic roles. Such resources facilitate the development of semantic role labeling (SRL) systems (Palmer et al., 2010) for proposition extraction, with a closed, predefined set of proto-roles.

To increase the coverage of propositions extracted, Open information extraction (OpenIE) formulations (Etzioni et al., 2008; Del Corro and Gemulla, 2013; Cui et al., 2018) were proposed to forgo the limits on a fixed set of semantic roles, account for both explicit and implicit predicates, and represent propositions in the unified form of (*subject, relation, object*) triples. However, we observe that OpenIE systems often fail to account for the complete set of propositions in a sentence. In many realistic cases, such as the *Andy Warhol’s hometown* example in Table 1, arguments in a proposition might not follow the same granularity as the ones in the sentence, e.g. *Andy Warhol* vs *Andy Warhol Museum*. Also, as OpenIE triples are still defined on direct predicate-argument relations, they often fail to produce a *decontextualized* (Choi et al., 2021) view of a proposition. For example, an OpenIE system would recognize the possessive relation “he has a hometown”, but fail to decontextualize the mentions of *he* → *Andy Warhol*, and *hometown* → *Pittsburgh*.

Furthermore, Gashteovski et al. (2020) and Fatahi Bayat et al. (2022) observe that existing neural OpenIE systems tend to extract long and over-specific arguments that potentially contain more compact propositions within them, leading to an incomplete set of propositions extracted. From the perspective of textual entailment, the reason that we want to extract the complete set of propositions in the most *compact* form possible is due to the

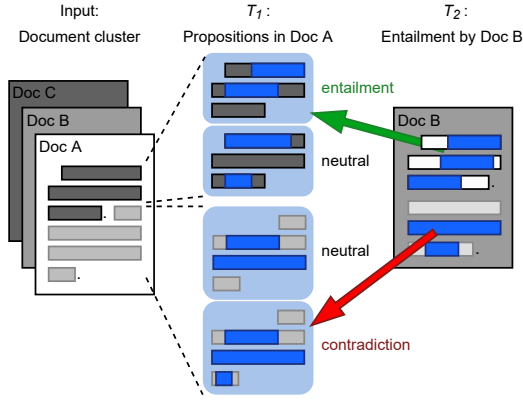


Figure 3: Given a cluster of related documents,  $T_1$  asks for each sentence of each document to be segmented into propositions, represented as subsets of a sentence’s tokens.  $T_2$  asks to classify the entailment relation  $\{entails, neutral, contradicts\}$  of each proposition in document A w.r.t. another document B from the same cluster; Our annotations also feature a single proposition in B that best supports each *entails* or *contradicts* label.

fact that their truth value could vary individually. Consider the simple example as follows - “Alice and Bob went to the zoo.” To recognize its truth value in the context of a premise, we need to individually evaluate each of the two propositions “Alice went to the zoo” and “Bob went to the zoo”. However, as the two propositions share the same verb predicate in the sentence, typically a neural SRL/OpenIE model that represents proposition in sequence labeling format would not be able split the conjunction “Alice and Bob”, and so would not manage to recognize the two propositions in separate forms (Kolluru et al., 2020).

To illustrate the difference between OpenIE and our approach, we offer a list of example propositions from our proposed PROPSEGMENT dataset, and compared them to extractions from rule-based and neural OpenIE systems, in Appendix C.

### 3 PROPSEGMENT Dataset

To facilitate research on recognizing propositional textual entailment, we propose PROPSEGMENT, a large-scale dataset with clusters of topically similar news and Wikipedia documents, with human annotated propositions and entailment labels.

#### 3.1 Task Definitions

We formulate the task of recognizing propositional textual entailment into two sub-tasks (Fig. 3). Given a hypothesis sentence and a premise document, a system is expected to (1) identify all the propositions within the hypothesis sentence, and

(2) classify the textual entailment relation of each proposition with respect to the premise document.

**$T_1$ : Propositional Segmentation** Given a sentence  $S$  with tokens  $[t_0, t_1, \dots, t_l]$  from a document  $D$ , a system is expected to identify the set of propositions  $\mathcal{P} \subseteq 2^S$ , where each proposition  $p \in \mathcal{P}$  is represented by a unique subset of tokens in sentence  $S$ . In other words, each proposition can be represented in sequence labeling format, per the example from Table 1. Each proposition is expected (1) to correspond to a distinct fact that a reader learns directly from reading the given sentence, (2) include all tokens within the sentence that are relevant to learning this fact, and (3) to not be equivalent to a conjunction of other propositions. We opt for this format as it does not require explicit annotation of the predicate-argument structure. This allows for more expressive power for propositions with implied or implicit predicates (Stern and Dagan, 2014). Also, representing each proposition as a separate sequence could effectively account for cases with shared predicate or arguments spans, and make evaluation more readily accessible.

Since the propositions, as we demonstrated earlier, do not necessarily have a unique and identifiable predicate word in the sentence, the typical inference strategy, e.g. in SRL or OpenIE, which first extracts the set of predicates, and then identifies the arguments with respect to each predicate would not work in this case. For this reason, given an input sentence, we expect a model on the task to directly output *all* propositions. In such *one-to-set* prediction setting, the output propositions of the model are evaluated as an unordered set.

**$T_2$ : Propositional Entailment** Given a hypothesis proposition  $p$  from document  $D_{hyp}$  and a whole premise document  $D_{prem}$ , a system is now expected to classify whether the premise entails the proposition, i.e. if the information conveyed by the proposition would be inferred true from the context provided by the premise.

As with the creation process for any textual entailment corpus or resources, the *indeterminacy* of entities and event coreference (Bowman et al., 2015), i.e. whether the entities and events described in a hypothesis can be assumed to refer to the same ones in the premise, has a large effect on the correctness of textual entailment label definition, and subsequently how a learned model would behave. Consider the following sentence re-



| Task/Setting   | WIKIPEDIA |       |       | NEWS  |       |       | FULL DATASET |       |       |
|----------------|-----------|-------|-------|-------|-------|-------|--------------|-------|-------|
|                | Train     | Dev   | Test  | Train | Dev   | Test  | Train        | Dev   | Test  |
| News Clusters  | 210       | 15    | 24    | 210   | 15    | 25    | 420          | 30    | 49    |
| Documents      | 630       | 45    | 72    | 630   | 45    | 75    | 1260         | 90    | 147   |
| Sentences      | 3305      | 259   | 1064  | 3293  | 236   | 1192  | 6598         | 495   | 2256  |
| Propositions   | 13879     | 1050  | 4679  | 11196 | 936   | 3964  | 25075        | 1986  | 8643  |
| ENTAIL Label % | 34.70     | 33.24 | 34.85 | 20.27 | 19.98 | 20.13 | 28.26        | 26.99 | 28.19 |

Table 2: Notable Statistics from the PROPSEGMENT dataset.

garding *Andy Warhol*: “A museum under his name is in Pittsburgh.”. This sentence alone cannot be verified true, unless enough context is provided to inform us that “his” refer to Andy Warhol here. In the case of propositions, this becomes critical, as each proposition presents an even narrower view of the whole context it appears in.

In previous RTE/NLI datasets (Dagan et al., 2005; Bowman et al., 2015; Williams et al., 2018), reference determinacy is guaranteed from the fact that hypotheses are created by annotators reading the premise. However, in our target application of comparing documents for example, reference determinacy cannot be safely assumed, as hypothesis and premise typically come from different documents, and their contexts thus differ. In PROPSEGMENT, to increase reference determinacy, the hypothesis and premise are sampled from different yet topically-aligned documents, i.e. about the same high-level event or entity.

### 3.2 Dataset Construction

We sample 250 document clusters from both the Wiki Clusters (Schuster et al., 2022) and NewSHead (Gu et al., 2020) datasets. Each cluster contains the first 10 sentences of three documents, either news articles on the same event, or Wikipedia pages in different languages (machine-translated into English) of the same entity. For each sentence, we train and instruct three human raters to annotate the set of propositions, each of which represented by a unique subset of tokens from the sentence. Conceptually, we instruct raters to include all the words that (1) pertain to the content of a proposition, and (2) are explicitly present in the sentence. For example, if there does not exist a predicate word for a proposition in the sentence, then only include the corresponding arguments. Referents present within the sentence are included in addition to pronominal and nominal references. We provide a more detailed description our rater guidelines and how propositions are defined with respect to various linguistic phenomena in Appendix B.

Given the three sets of propositions from the three raters for a sentence, we reconcile and select one of the three raters’ responses with the highest number of propositions that the other two sets also include as our ground truth set of propositions. Since the exact selection of tokens used to mark a proposition may vary slightly across different raters, we allow for fuzziness when measuring the match between two propositions. Following FitzGerald et al. (2018) and Roit et al. (2020), we use Jaccard similarity, i.e. intersection over union of the tokens included in the two propositions respectively, to measure the similarity between two propositions. We say two propositions match if their Jaccard similarity is greater or equal to a threshold  $\theta = 0.8$ , and align two raters’ responses using unweighted bipartite matching between propositions satisfying the Jaccard threshold.

Next, for all propositions in a document, we sample one other document from the document cluster as premise, and ask three raters to label the textual entailment relation between each proposition and the premise, i.e. one of  $\{Entailment, Neutral, Contradiction\}$ . We take the majority vote from the three as the gold entailment label. Interestingly, we observe that only 0.2% of all annotated labels from the rater are “contradictions”, suggesting a low presence of natural contradiction cases in the document clusters setup for the Wiki and news domain. For this reason, we choose to only consider two-way label ( $\{Entailment, Non-Entailment\}$ ) for the entailment task evaluation.

The statistics for the PROPSEGMENT dataset are shown in Table 2. We create the train/dev/test splits based on clusters, so that documents in each cluster exclusively belong to only one of the splits. Overall, the dataset features 1497 documents with  $\sim 35K$  propositions with entailment labels.

## 4 Baseline Methods

In this section, we provide and describe the baseline methods for the proposition extraction and propositional entailment tasks respectively.

| Task/Setting                                | Model                      | Jaccard $\theta = 0.8$     |                   |              | Exact Match                    |              |              |
|---|----------------------------|----------------------------|-------------------|--------------|--------------------------------|--------------|--------------|
|   |                            | Precision                  | Recall            | F1           | Precision                      | Recall       | F1           |
| T <sub>1</sub> : Propositional Segmentation | BERT-Base                  | 33.77                      | 33.53             | 33.65        | 14.33                          | 14.60        | 14.47        |
|   | BERT-Large                 | 34.97                      | 33.42             | 34.17        | 14.61                          | 14.16        | 14.38        |
|   | T5-Base                    | 54.96                      | 51.93             | 53.41        | 32.87                          | 31.54        | 32.19        |
|   | T5-Base w/ <i>Entail.</i>  | 53.54                      | 51.50             | 52.50        | 31.61                          | 30.67        | 31.13        |
|   | T5-Large                   | <b>55.95</b>               | <b>55.05</b>      | <b>55.50</b> | <b>32.40</b>                   | <b>32.16</b> | <b>32.28</b> |
|   | T5-Large w/ <i>Entail.</i> | <b>56.27</b>               | <b>55.50</b>      | <b>55.89</b> | <b>31.94</b>                   | <b>32.11</b> | <b>32.02</b> |
|   | Human Performance          | 69.63                      | 64.69             | 67.07        | 44.86                          | 42.93        | 43.87        |
| T <sub>2</sub> : Propositional Entailment   |                            | Performance (2-way Class.) |                   |              | Per-Label $F_1$ (3-way Class.) |              |              |
|   |                            | Accuracy                   | Balanced Accuracy |              | Entail.                        | Neutral      | Contra.      |
|   | <i>Always Entails.</i>     | 27.89                      | 50.00             |              | 43.62                          | 0.00         | 0.00         |
|   | <i>Always Neutral</i>      | 72.10                      | 50.00             |              | 0.00                           | 83.54        | 0.00         |
|   | T5-Base                    | 85.17                      | 81.44             |              | 73.32                          | 89.68        | 11.21        |
|   | T5-Large                   | <b>91.38</b>               | <b>89.75</b>      |              | <b>84.78</b>                   | <b>93.98</b> | <b>20.34</b> |
|   | Human Performance          | 90.20                      | 88.31             |              | -                              | -            | -            |

Table 3: Performance of the baseline models on the full (WIKI + NEWS) test set of PROPSEGMENT. Due to the low presence of contradiction examples in the test set ( $32/8643 = 0.4\%$  examples), the difference in  $F_1$  score w.r.t contradiction does not reflect statistically significant improvement.

#### 4.1 Propositional Segmentation Baselines

The key challenge with the proposition extraction task lies within the one-to-set structured prediction setting. Our one-to-set prediction format is similar to QA-driven semantic parsing such as QA-SRL (He et al., 2015; Klein et al., 2022), as both involve generating a variable number of units of semantic content under no particular order between them. As in propositions, there might not necessarily be a unique and identifiable predicate word associated with each proposition, extracting predicates first (e.g. as a sequence tagging task), and later individually produce one proposition for each predicate would not be a sufficient solution in this case.

For this particular problem setup, we introduce two classes of baseline architectures.

**Seq2Seq: T5** (Raffel et al., 2020) When formatting a output set as a sequence, Seq2Seq models, as they employ chain-rules to efficiently model the joint probability of output sequences, have been found to be a strong method for tasks with set outputs (Vinyals et al., 2016).

The obvious caveat for representing set outputs as sequences is that we need an ordering for the outputs. Having a consistent ordering helps seq2seq model learn to maintain the set structure in the output (Vinyals et al., 2016), and the best ordering scheme is often both model- and task-specific (Klein et al., 2022). In our preliminary experiments, we observe that sorting the propositions by the appearance order of the tokens in the sentence, i.e. positions of the foremost tokens of each proposi-

tion in the sentence, yields the best performance.

We start from the T5 1.1 checkpoints with the T5x library (Roberts et al., 2022). Given a sentence input, we finetune the T5 model to output the propositions in a single sequence. For each input sentence, we sort the output propositions using the aforementioned ordering scheme, and join them by a special token [TARGET]. The spans of tokens included in each proposition is surrounded by special tokens [M] and [/M]. For instance, “ [M]Alice [/M] and Bob [M]went to the Zoo [/M]. [TARGET] Alice and [M]Bob went to the Zoo. [/M] ”.

In addition, we evaluate the setting where the model is also given the premise document  $D_{prem}$ , and learns to output the entailment label along with each proposition (T5 w/ *Entail.* in Table 3).

**Encoder+Tagger: BERT** (Devlin et al., 2019)

For comparison, we provide a simpler baseline that does not model joint probability of the output propositions. On top of the last hidden layer output of an encoder mode, i.e. BERT, we add  $K$  individual linear layers that each correspond to one output proposition. Given an input sentence, the  $i^{th}$  linear layer produces a binary (0/1) label for each token, indicating whether the token should be inside the  $i^{th}$  proposition or not.  $K$  is set to be a sufficiently large number, e.g.  $K = 20$  in our experiments. We use the label of the [CLS] token of the  $i^{th}$  linear layer to indicate whether the  $i^{th}$  proposition should exist in the output. For such, we follow the same ordering of the output propositions

as in the seq2seq (T5) baseline setup.

## 4.2 Propositional Entailment Baselines

We formulate the propositional entailment task as a sequence labeling problem, and finetune T5 model as our baseline. The inputs consist of the hypothesis proposition  $p$  with its document context  $D_{hyp}$ , plus the premise document  $D_{prem}$ . The output is one of the three-way labels  $\{Entailment, Neutral, Contradiction\}$ . Due to the low (0.2%) presence of contradictions in our dataset, we evaluate the task as two-way classification, i.e.  $\{Entailment, Non-Entailment\}$ . Under such settings, we merge the *neutral* and *contradiction* outputs from the model as *non-entailments* during inference. To ensure that the model has access to the essential context information, our task input also include the document  $D_{hyp}$  of the hypothesis proposition  $p$ , so that model has a decontextualized view of  $p$  when inferring its textual entailment relation with  $D_{prem}$ .

## 5 Experiments and Results

### 5.1 Evaluation Metrics

**Propositional Segmentation** We measure the precision and recall between the set of predicted and gold propositions for a given sentence. As the set of gold propositions do not follow any particular ordering, we first produce a bipartite matching between them using the Hungarian algorithm (Kuhn, 1955). To measure whether two propositions match with each other, we adopt two different functions. First, we follow similar strategy as in our rater annotation reconciliation (§ 3.2), by measuring if the Jaccard similarity between the predicted and sets of tokens is larger than a threshold  $\theta$ . We use same threshold value  $\theta = 0.8$ . We also use exact match, an even more restrictive measure where two propositions match if and only if they have the exact same tokens. We report the macro-averaged precision and recall over sentences in the test set.

**Propositional Entailment** We report the baseline performance under two-way classification results in accuracy. As there exist label imbalance in our dataset (Table 2), we also report the balanced accuracy, i.e. average of true positive rate and true negative rate. To help us understand the per-label performance, we also report the  $F_1$  score w.r.t. each of the three-way label respectively.

| Train Domain | Test Domain<br>( $P/R/F_1$ w/ Jaccard $\theta = 0.8$ ) |                   |                   |
|--------------|--|-------------------|-------------------|
|              |  | WIKI              | NEWS              |
|              |  |                   |                   |
| WIKI         |  | 53.95/53.16/53.56 | 44.93/44.95/44.94 |
| NEWS         |  | 45.21/43.65/44.42 | 49.58/47.81/48.68 |

Table 4: Cross-domain (i.e. train on NEWS  $\rightarrow$  test on WIKI, and train on WIKI  $\rightarrow$  test on NEWS) generalization results of T5-large models on the Propositional Segmentation ( $T_1$ ) task.

### 5.2 Baseline Results

Table 3 shows the evaluation results for the segmentation ( $T_1$ ) and entailment task ( $T_2$ ) respectively.

For the segmentation task ( $T_1$ ), the seq2seq T5 model setup yields superior performance compared to the simpler encoder+tagger BERT setup. As the encoder+tagger setup predicts each proposition individually, and does not attend on other propositions during inference, we observe that the model predicts repeated/redundant propositions in  $> 20\%$  of the input sentences. On the other hand, in the seq2seq T5 setup, the repetition rate is less than 1%. For both setups, we manually remove the redundant outputs as a post processing step. We also evaluate the multi-task setup (i.e. T5 w/ *Entail.* in Table 3) where the model outputs the entailment label along with each proposition, and observe no significant difference in the performance.

For the entailment task ( $T_2$ ), we see that T5-Large yields the best overall performance. We observe that the performance with respect to the *entailment* label is lower compared to the *neutral* label, due to label imbalance. As the *contradiction* label has low presence both in the training and test split, the performance on the label do not offer statistically significant comparison.

For both tasks, we estimate the averaged human expert performance by comparing annotations from three of the authors to ground truth on 50 randomly sampled examples from the dataset. We observe that for the segmentation task  $T_1$ , there remains a sizable gap between the best model, T5-Large, and human performance, while on the entailment task  $T_2$ , T5-Large exceeds human performance, which is not uncommon among natural language understanding tasks of similar classification formats (Wang et al., 2019).

### 5.3 Cross-Domain Generalization

On the propositional segmentation ( $T_1$ ) task, we evaluate the how the best baseline model generalizes across the Wikipedia (Wiki) and News do-

mains. Table 4 shows the results of T5-Large models finetuned on data from each domain, and evaluated on the test split of both domains.

When applying a model trained on `Wiki`, we see a larger drop in performance when tested on `News`, as the `News` domain features more syntactic and style variations compared to the `Wiki` domain.

## 6 Analysis and Discussion

In this section, we exemplify the utilities of our propositional segmentation and entailment framework, which we refer to as PropNLI, through the lens of two downstream use cases, e.g. summary hallucination detection (§ 6.1), and document-level NLI w/ variable-length hypotheses (§ 6.2).

### 6.1 Application: Hallucination Detection

To demonstrate the utility of our dataset and the PropNLI framework, we look at the task of summarization hallucination detection, i.e. given a machine-generated summary of a source document, identify whether its content can be inferred from, or *faithful* to the document. Naturally the task can be represented as a NLI problem, and NLI systems have been shown effective on the task (Kryscinski et al., 2020; Chen et al., 2021). As the summaries can be long and compositional, recognizing partial entailment, and identifying which part(s) of a summary is hallucinated becomes important (Goyal and Durrett, 2020; Laban et al., 2022). The goal closely resembles our framework of recognizing propositional entailment.

To show that PropNLI can be used for hallucination detection, we experiment on the model generated summaries on the XSum dataset (Narayan et al., 2018), where Maynez et al. (2020) provide human annotations of the sets of hallucinated spans (if they exist) in the summaries. Table 5 illustrates our idea of using entailment labels of propositions to infer whether and where hallucination exists. If a proposition in a summary is *entailed* by the document, then all spans covered by the proposition are faithful to the source document. Otherwise if a proposition is not entailed, then some of its spans are hallucinated, which suggest that the summary contains *hallucinated* information.

Following such intuitions, we first evaluate the performance of our method in zero-shot settings as a hallucination classifier, i.e. binary classification for whether a summary is hallucinated or not. For baseline comparison, we use a T5-large model fine-

**Document:** The incident happened near Dr Gray’s Hospital shortly after 10:00. The man was taken to the hospital with what police said were serious but not life-threatening injuries. The A96 was closed in the area for several hours, but it has since reopened.

#### Summary w/ human labeled hallucinated spans:

A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire.

#### Predicted propositions (blue) and entailment labels

#1: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✓

#2: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✗

#3: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✗

#4: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✗

#### Predicted hallucinated spans (union of ✗- union of ✓)

A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire.

Table 5: An example model generated summary on the XSum dataset, with human-annotated hallucination spans from Maynez et al. (2020). We show that we can infer the hallucinated spans from the set of four propositions and their entailment labels (*entail*=✓, *not-entail*=✗), predicted by our T5-Large models. More examples can be found in Appendix D

| Method  | Hallu. Class. B. Acc. | Span Detection |     |                |               |     |                |
|---------|-----------------------|----------------|-----|----------------|---------------|-----|----------------|
|         |                       | Faith. Tokens  |     |                | Hallu. Tokens |     |                |
|         |                       | P              | R   | F <sub>1</sub> | P             | R   | F <sub>1</sub> |
| PropNLI | <b>.62</b>            | .78            | .50 | .61            | .64           | .71 | .67            |
| MNLI    | .59                   | .96            | .17 | .30            | .56           | .88 | .68            |

Table 6: Zero-shot performance of PropNLI vs. T5-Large MNLI model on hallucination identification and span detection tasks from Maynez et al. (2020).

tuned on MNLI (Williams et al., 2018) to classify a full summary as entailed ( $\rightarrow$  *faithful*) or not ( $\rightarrow$  *hallucinated*). As ~89% of the summaries annotated by Maynez et al. (2020) are hallucinated, we again adopt balanced accuracy (§ 5.1) as the metric. On 2500 examples, our method achieved 61.68% balanced accuracy, while MNLI achieved 58.79%.

Next, we study whether the entailment labels of propositions can be composed to detect hallucinated spans in a summary. As illustrated in Table 5, we take the union of the spans in *non-entailed* propositions, and exclude the spans that has appeared in *entailed* propositions. The intuition is that the hallucinated information likely only exists within the non-entailed propositions, but not the entailed ones.

We evaluate hallucinated span detection as a token classification task. For each summary, we evaluate the precision and recall of the *faithful* and



*hallucinated* set of predicted tokens respectively against the human-labeled ground truth set. We report the macro-averaged precision, recall and  $F_1$  score over all 2,500 summaries. We compare to the same T5-Large MNLI model, where we label all tokens as *faithful* if the summary is predicted to be *entailed*, and all tokens as *hallucinated* otherwise. We report the performance with respect to each of the two labels in Table 6. As the MNLI model don’t distinguish partial entailment from non-entailment cases, it predicts more tokens to be hallucinated, and thus having low precision and high recall on the hallucinated tokens, and vice versa. On the other hand, we observe our model can be used to detect the nuance between faithful and hallucinated tokens with good and more balanced performance for both cases. Table 5 shows one example summary and PropNLI’s predictions, and we include more examples in Appendix D.

## 6.2 Proposition-Level → Sentence/Paragraph-Level Entailment

As our study is motivated in part to decomposing sentence/paragraph-level NLI task to the proposition-level, we would like to see whether proposition-level entailment labels can potentially be *composed* to explain the sentence/paragraph-level NLI inference and predictions.

Given a hypothesis sentence/paragraph, and a premise, our PropNLI framework takes three steps. First we segment the hypothesis into propositions. For each proposition, we infer its entailment relation with the premise. And lastly, we compose the proposition-level entailment predictions into a overall hypothesis-level label, with an aggregation function. The first two steps resemble the two tasks of our dataset respectively. As for the aggregation function, we start with the simple strategy of using logical conjunction, i.e. the hypothesis is entailed by the premise if and only if all propositions within it are entailed. We hypothesize that in cases where multiple propositions exist in the hypothesis, the framework offers a more precise and finer-grained description of the textual entailment relation between the premise and hypothesis.

To demonstrate the utility of the idea, we conduct a case study on DocNLI (Yin et al., 2021), which features premise and hypothesis of different length, and so varying number and compositions of propositions. We take the baseline T5-Large segmentation and entailment models respectively, and use logical conjunction to aggregate the

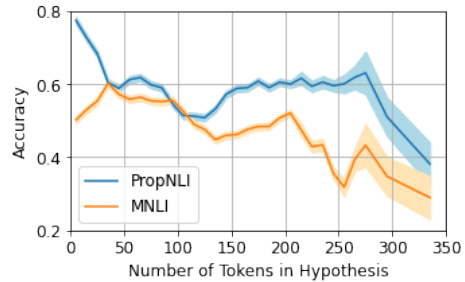


Figure 4: Zero-shot performance of T5-large MNLI model compared to our PropNLI T5-large models (i.e. proposition-level segmentation → entailment → aggregation) with respect to varying *hypothesis length* in DocNLI dev set. The shaded region shows 95% confidence interval.

proposition-level entailment prediction. We compare PropNLI in a zero-shot setting against the T5-Large MNLI model. The MNLI model takes the entire hypothesis and premise and input without any segmentation or decomposition.

The comparison results are shown in Figure 4. We take the development set of DocNLI and split examples into buckets according to number of tokens in the hypothesis. We examine the zero-shot performance of the PropNLI setup versus the finetuned MNLI model. We observe that with shorter hypotheses (< 100 tokens), the two setups demonstrated similar performance, as the hypothesis length is similar to the distribution of MNLI training set (avg. 21.73 tokens  $\pm$  30.70). As the length of the hypothesis increases, the performance of MNLI model starts to drop, while PropNLI’s performance remains relatively stable. Such observations suggest the potential of using the PROPSEGMENT and PropNLI framework to develop more generalizable NLI models, especially in the realistic case where input hypotheses are compositional.

## 7 Conclusion

In this paper, we presented PROPSEGMENT, the first large-scale dataset for studying proposition-level segmentation and entailment. We demonstrate that segmenting a text expression into propositions, i.e. atomic units of meanings, and assessing their truth values would provide a finer-grained characterization of the textual entailment relation between two pieces of text. Beyond NLI/RTE tasks, we hypothesize that proposition-level segmentation might be helpful in similar ways for other text classification tasks as well. We hope that PROPSEGMENT will serve as a starting point, and pave a path for research forward along the line.

## Acknowledgements

We thank Michael Collins, Corinna Cortes, Paul Haahr, Ilya Kornakov, Ivan Kuznetsov, Annie Louis, Don Metzler, Jeremiah Milbauer, Pavel Nalivayko, Fernando Pereira, Sandeep Tata, Yi Tay, Andrew Tomkins, and Victor Zaytsev for insightful discussions, suggestions, and support. We are grateful to the annotators for their work in creating PROPSEGMENT.

## References

- Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. TensorFlow: a system for Large-Scale machine learning. In *12th USENIX symposium on operating systems design and implementation (OSDI 16)*, pages 265–283.
- Collin F Baker, Charles J Fillmore, and John B Lowe. 1998. The Berkeley FrameNet project. In *COLING 1998 Volume 1: The 17th International Conference on Computational Linguistics*.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Sihao Chen, Fan Zhang, Kazuo Sone, and Dan Roth. 2021. [Improving faithfulness in abstractive summarization with contrast candidate generation and selection](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5935–5941, Online. Association for Computational Linguistics.
- Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. 2021. [Decontextualization: Making sentences stand-alone](#). *Transactions of the Association for Computational Linguistics*, 9:447–461.
- Lei Cui, Furu Wei, and Ming Zhou. 2018. [Neural open information extraction](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 407–413, Melbourne, Australia. Association for Computational Linguistics.
- I. Dagan and O. Glickman. 2004. Probabilistic textual entailment: Generic applied modeling of language variability. In *Learning Methods for Text Understanding and Mining*.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The PASCAL recognising textual entailment challenge. In *Machine learning challenges workshop*, pages 177–190. Springer.
- Marie-Catherine de Marneffe, Anna N. Rafferty, and Christopher D. Manning. 2008. Finding contradictions in text.
- Luciano Del Corro and Rainer Gemulla. 2013. ClausIE: clause-based open information extraction. In *Proceedings of the 22nd international conference on World Wide Web*, pages 355–366.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Oren Etzioni, Michele Banko, Stephen Soderland, and Daniel S Weld. 2008. Open information extraction from the web. *Communications of the ACM*, 51(12):68–74.
- Farima Fatahi Bayat, Nikita Bhutani, and H. Jagadish. 2022. [CompactIE: Compact facts in open information extraction](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 900–910, Seattle, United States. Association for Computational Linguistics.
- Nicholas FitzGerald, Julian Michael, Luheng He, and Luke Zettlemoyer. 2018. [Large-scale QA-SRL parsing](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2051–2060, Melbourne, Australia. Association for Computational Linguistics.
- Kiril Gashteovski, Rainer Gemulla, Bhushan Kotnis, Sven Hertling, and Christian Meilicke. 2020. [On aligning OpenIE extractions with knowledge bases: A case study](#). In *Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems*, pages 143–154, Online. Association for Computational Linguistics.
- Tanya Goyal and Greg Durrett. 2020. [Evaluating factuality in generation with dependency-level entailment](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3592–3603, Online. Association for Computational Linguistics.
- Xiaotao Gu, Yuning Mao, Jiawei Han, Jialu Liu, You Wu, Cong Yu, Daniel Finnie, Hongkun Yu, Jiaqi Zhai, and Nicholas Zukoski. 2020. Generating representative headlines for news stories. In *Proceedings of The Web Conference 2020*, pages 1773–1784.

- Luheng He, Mike Lewis, and Luke Zettlemoyer. 2015. [Question-answer driven semantic role labeling: Using natural language to annotate natural language](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 643–653, Lisbon, Portugal. Association for Computational Linguistics.
- Paul R Kingsbury and Martha Palmer. 2002. From TreeBank to PropBank. In *LREC*, pages 1989–1993.
- Ayal Klein, Eran Hirsch, Ron Eliav, Valentina Pyatkin, Avi Caciularu, and Ido Dagan. 2022. QASem parsing: Text-to-text modeling of QA-based semantics. *arXiv preprint arXiv:2205.11413*.
- Keshav Kolluru, Samarth Aggarwal, Vipul Rathore, Mausam, and Soumen Chakrabarti. 2020. [IMoJIE: Iterative memory-based joint open information extraction](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5871–5886, Online. Association for Computational Linguistics.
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. [Evaluating the factual consistency of abstractive text summarization](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346, Online. Association for Computational Linguistics.
- Harold W Kuhn. 1955. The Hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Omer Levy, Torsten Zesch, Ido Dagan, and Iryna Gurevych. 2013. [Recognizing partial textual entailment](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 451–455, Sofia, Bulgaria. Association for Computational Linguistics.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. [On faithfulness and factuality in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Adam Meyers, Ruth Reeves, Catherine Macleod, Rachel Szekely, Veronika Zielinska, Brian Young, and Ralph Grishman. 2004. The NomBank project: An interim report. In *Proceedings of the workshop frontiers in corpus annotation at hlt-naacl 2004*, pages 24–31.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. [Don’t give me the details, just the summary!](#) [topic-aware convolutional neural networks for extreme summarization](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The proposition bank: An annotated corpus of semantic roles. *Computational linguistics*, 31(1):71–106.
- Martha Palmer, Daniel Gildea, and Nianwen Xue. 2010. Semantic role labeling. *Synthesis Lectures on Human Language Technologies*, 3(1):1–103.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67.
- Adam Roberts, Hyung Won Chung, Anselm Levskaya, Gaurav Mishra, James Bradbury, Daniel Andor, Sharan Narang, Brian Lester, Colin Gaffney, Afroz Mohiuddin, et al. 2022. Scaling up models and data with t5x and seqio. *arXiv preprint arXiv:2203.17189*.
- Paul Roit, Ayale Klein, Daniela Stepanov, Jonathan Mamou, Julian Michael, Gabriel Stanovsky, Luke Zettlemoyer, and Ido Dagan. 2020. [Controlled crowdsourcing for high-quality QA-SRL annotation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7008–7013, Online. Association for Computational Linguistics.
- Tal Schuster, Sihao Chen, Senaka Buthpitiya, Alex Fabrikant, and Donald Metzler. 2022. Stretching Sentence-pair NLI Models to Reason over Long Documents and Clusters. In *Findings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. [Get your vitamin C! robust fact verification with contrastive evidence](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 624–643, Online. Association for Computational Linguistics.
- Gabriel Stanovsky, Julian Michael, Luke Zettlemoyer, and Ido Dagan. 2018. [Supervised open information extraction](#). 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 885–895, New Orleans, Louisiana. Association for Computational Linguistics.
- Asher Stern and Ido Dagan. 2014. [Recognizing implied predicate-argument relationships in textual inference](#). In *Proceedings of the 52nd Annual Meet-*

ing of the Association for Computational Linguistics (Volume 2: Short Papers), pages 739–744, Baltimore, Maryland. Association for Computational Linguistics.

Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2016. Order matters: Sequence to sequence for sets. In *Proceedings of the International Conference on Learning Representations*.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2019. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *7th International Conference on Learning Representations, ICLR 2019*.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. 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 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Wenpeng Yin, Dragomir Radev, and Caiming Xiong. 2021. DocNLI: A large-scale dataset for document-level natural language inference. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4913–4922, Online. Association for Computational Linguistics.

Yi Zhang, Zachary Ives, and Dan Roth. 2020. “who said it, and why?” provenance for natural language claims. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4416–4426, Online. Association for Computational Linguistics.

## A Model Implementation

**T5** We use T5 1.1 checkpoints from the T5x library (Roberts et al., 2022), with Flaxformer<sup>2</sup> implementation. For all sizes of T5 model and all tasks, we finetune the model for three epoch, with  $1e-3$  learning rate, 0.1 dropout rate, batch size of 128. We train the models on 16 TPU v3 slices.

**BERT** We use the BERT English uncased models from Tensorflow (Abadi et al., 2016), in large (24 layers, 16 attention heads, 1024 max sequence length) and base (12 layers, 12 attention heads, 768 max sequence length) sizes. For both sizes, we finetune the model for five epoch, with  $1e-5$  learning rate, 0.1 dropout rate, batch size of 16. We train the models on 8 TPU v3 slices.

## B Annotation Guidelines

### B.1 Segmentation annotation guidelines

There is no unequivocally unique definition for precisely how to segmenting an English sentence in the context of a document into propositions defined as token subsets, due to a variety of complex language phenomena. Our raters were instructed to follow the following overall guidelines for the segmentation task:

1. Each proposition is expected to correspond to a distinct fact that a reader learns directly from reading the given sentence.
  - (a) The raters are instructed to focus on the text’s most literal *denotation*, rather than drawing further inferences from the text based on world knowledge, external knowledge, or common sense.
  - (b) The raters are instructed to consider *factivity*, marking only propositions that, in their judgement, the author intends the reader to take as factual from reading the sentence.
  - (c) With regard to quotes, raters are asked to estimate the author’s intent, including the proposition quoted when the reader is expected to take it as factual, and/or the proposition of the quote itself having been uttered if the reader is expected to learn that a speaker uttered that quote.
  - (d) The raters are instructed to omit text that are clearly non-factual, such as rhetorical

<sup>2</sup><https://github.com/google/flaxformer>



flourishes or first-person account of an article author’s emotional response to the topic. This rule is specific to the news and Wikipedia domains, since in other domains of prose, first-person emotions may well be part of the intended informational payload.

2. Each proposition should include all tokens within the sentence that are relevant to learning this fact.
  - (a) Specifically, the raters are asked to include any tokens in the same sentence that are antecedents of pronouns or other endophora in the proposition, or relevant bridging references.
  - (b) Raters are asked to ignore punctuation, spacing, and word inflections when selecting tokens, though a number of other minutiae, such as whether to include articles, are left unspecified in the rater instructions.
3. Choose the simplest possible propositions, so that no proposition is equivalent to a conjunction of the other propositions, and so that the union of all of the sentence’s proposition gives us all the information a reader learns from the sentence.

The raters are also asked to omit propositions from any text that doesn’t constitute well-formed sentences, typically arising from parsing errors or from colloquialisms.

Note that the resulting subsets of tokens do not, generally, constitute well-formed English sentences when concatenated directly, but can, in our ad hoc trials, easily be reconstituted into stand-alone sentences by a human reader.

## B.2 Entailment annotation guidelines

For the propositional entailment task, our instructions are somewhat similar to the RTE task (Dagan and Glickman, 2004), but specialized to the proposition level.

The raters are asked to read the premise document and decide whether a specific hypothesis proposition is entailed by it, contradicted, or neither. In the first two cases, the raters are asked to mark a proposition in the premise document that most closely supports the hypothesis proposition, using the same definition of proposition as above.

The interface nudges the raters to select one of the propositions marked by the segmentation rater, but allows the entailment rater to create a new proposition as well. Note that the choice of a specific supporting proposition is sometimes not well defined.

To judge entailment, the raters are asked “from reading just the premise document, do we learn that the hypothesis proposition is true, learn that it’s false, or neither?” More specifically, the raters are asked:

1. To consider the full document of the hypothesis as the context of the hypothesis proposition, and the full premise document.
2. To allow straightforward entailment based on “common sense or widely-held world knowledge”, but otherwise avoid entailment labels whenever “significant analysis” (any complex reasoning, specialized knowledge, or subjective judgement) is required to align the two texts.
3. To assume that the two documents were written in the same coarse spatiotemporal context — same geographical area, and the same week.

Raters have the option of marking that they don’t understand the premise and/or the hypothesis and skipping the question.

## C Example Propositions From OpenIE vs. PROPSEGMENT

To illustrate the difference between how we define propositions in PROPSEGMENT, versus OpenIE formulations, we include a few example sentences with propositions in PROPSEGMENT in Table 7 and 8, and compare propositions extracted with ClausIE, a rule-based OpenIE model (Del Corro and Gemulla, 2013), and a neural Bi-LSTM model from Stanovsky et al. (2018).

## D XSum Hallucination Detection - Examples

Table 9 and 10 show two example documents, with propositions and the inferred hallucinated spans in model-generated and gold summaries by our PropNLI model. We compare the predictions to the annotations of hallucinated span provided by Maynez et al. (2020).

---

**Sentence:** The 82nd NFL Draft took place from April 27-29, 2017 in Philadelphia.

**PROPSEGMENT**

#1: [The 82nd NFL Draft took place from April 27-29, 2017](#) in Philadelphia.

#2: [The 82nd NFL Draft took place from April 27-29, 2017 in Philadelphia.](#)

**ClausIE**

#1: (The 82nd NFL Draft, took place, from April 27-29, 2017 in Philadelphia)

#2: (The 82nd NFL Draft, took place, from April 27-29, 2017)

**Neural Bi-LSTM OIE** (*Splitting each modifier, i.e. ARGM*)

#1: (The 82nd NFL Draft, took, place, from April 27-29, 2017)

#2: (The 82nd NFL Draft, took, place, in Philadelphia)

---

**Sentence:** She has also appeared in films such as Little Women (1994), The Hours (2002), Self Defense (1997), Les Misérables (1998) and Orson Welles y yo (2009).

**PROPSEGMENT**

#1: [She has also appeared in films such as Little Women](#) (1994), The Hours (2002), Self Defense (1997), Les Misérables (1998) and Orson Welles y yo (2009).

#2: [She has also appeared in films such as](#) Little Women (1994), [The Hours](#) (2002), Self Defense (1997), Les Misérables (1998) and Orson Welles y yo (2009).

#3: [She has also appeared in films such as](#) Little Women (1994), The Hours (2002), [Self Defense](#) (1997), Les Misérables (1998) and Orson Welles y yo (2009).

#4: [She has also appeared in films such as](#) Little Women (1994), The Hours (2002), Self Defense (1997), [Les Misérables](#) (1998) and Orson Welles y yo (2009).

#5: [She has also appeared in films such as](#) Little Women (1994), The Hours (2002), Self Defense (1997), Les Misérables (1998) and [Orson Welles y yo](#) (2009).

#6: She has also appeared in films such as [Little Women](#) (1994), The Hours (2002), Self Defense (1997), Les Misérables (1998) and Orson Welles y yo (2009).

#7: She has also appeared in films such as Little Women (1994), [The Hours](#) (2002), Self Defense (1997), Les Misérables (1998) and Orson Welles y yo (2009).

#8: She has also appeared in films such as Little Women (1994), The Hours (2002), [Self Defense](#) (1997), Les Misérables (1998) and Orson Welles y yo (2009).

#9: She has also appeared in films such as Little Women (1994), The Hours (2002), Self Defense (1997), [Les Misérables](#) (1998) and Orson Welles y yo (2009).

#10: She has also appeared in films such as Little Women (1994), The Hours (2002), Self Defense (1997), Les Misérables (1998) and [Orson Welles y yo](#) (2009).

**ClausIE**

#1: (She, has appeared, in films such as Little Women also)

#2: (She, has appeared, in films such as The Hours also)

#3: (She, has appeared, in films such as Self Defense also)

#4: (She, has appeared, in films such as Les Misérables also)

#5: (She, has appeared, in films such as Orson Welles y yo also)

#6: (She, has appeared, in films such as Little Women)

#7: (She, has appeared, in films such as The Hours)

#8: (She, has appeared, in films such as Self Defense)

#9: (She, has appeared, in films such as Les Misérables)

#10: (She, has appeared, in films such as Orson Welles y yo)

#11: (Little Women, is, 1994)

#12: (The Hours, is, 1994)

#13: (Self Defense, is, 1994)

#14: (Les Misérables, is, 1994)

#15: (Orson Welles y yo, is, 1994)

#16: (The Hours, is, 2002)

#17: (Self Defense, is, 1997)

#18: (Les Misérables, is, 1998)

#19: (Orson Welles y yo, is, 2009)

**Neural Bi-LSTM OIE**

#1: (She, appeared, in films such as Little Women (1994), The Hours (2002), Self Defense (1997), Les Misérables (1998) and Orson Welles y yo (2009))

---

Table 7: Comparison of propositions in PROPSEGMENT with extractions with ClausIE (Del Corro and Gemulla, 2013), and the neural Bi-LSTM OIE model from Stanovsky et al. (2018).

---

**Sentence:** The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art.

**PROPSegment**

#1: The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art.

#2: The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art.

#3: The Andy Warhol Museum in his hometown, Pittsburgh, Pennsylvania, contains an extensive permanent collection of art.

**ClausIE**

#1: (his, has, hometown)

#2: (his hometown, is, Pittsburgh Pennsylvania)

#3: (The Andy Warhol Museum in his hometown, contains, an extensive permanent collection of art)

**Neural Bi-LSTM OIE**

#1: (The Andy Warhol Museum in his hometown Pittsburgh Pennsylvania, contains, an extensive permanent collection of art)

---

**Sentence:** The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

**PROPSegment**

#1: The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

#2: The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

#3: The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

#4: The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

#5: The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

#6: The Cleveland Cavaliers got the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.

**ClausIE**

#1: (The Cleveland Cavaliers, got, the first choice in the lottery)

#2: (the lottery, was used, on 20-year-old forward Anthony Bennett)

#3: (Anthony Bennett, is, a freshman from the University of Nevada)

**Neural Bi-LSTM OIE**

#1: (The Cleveland Cavaliers, got, the first choice in the lottery, which was used on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.)

#2: (the lottery, was used, on 20-year-old forward Anthony Bennett, a freshman from the University of Nevada.)

---

Table 8: (Cont.) Comparison of propositions in PROPSegment with extractions with ClausIE (Del Corro and Gemulla, 2013), and the neural Bi-LSTM OIE model from Stanovsky et al. (2018).

---

**Document:** The incident happened near Dr Gray’s Hospital shortly after 10:00. The man was taken to the hospital with what police said were serious but not life-threatening injuries. The A96 was closed in the area for several hours, but it has since reopened.

---

**Summary from BertS2S**

A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire.

**Predicted propositions (blue) and entailment labels**

#1: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✓

#2: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✗

#3: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✗

#4: A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

A man has been taken to hospital following a one-vehicle crash on the A96 in Aberdeenshire.

---

**Summary from TConvS2S**

a man has been taken to hospital after being hit by a car in Moray.

**Predicted propositions (blue) and entailment labels**

#1: a man has been taken to hospital after being hit by a car in Moray. ✓

#2: a man has been taken to hospital after being hit by a car in Moray. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

a man has been taken to hospital after being hit by a car in Moray.

---

**Gold Summary from the XSum dataset**

A cyclist has suffered serious head injuries after a collision with a car in Elgin.

**Predicted propositions (blue) and entailment labels**

#1: A cyclist has suffered serious head injuries after a collision with a car in Elgin. ✗

#2: A cyclist has suffered serious head injuries after a collision with a car in Elgin. ✗

#3: A cyclist has suffered serious head injuries after a collision with a car in Elgin. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

A cyclist has suffered serious head injuries after a collision with a car in Elgin.

---

**Summary from PTGen**

A man has been taken to hospital after being hit by a car in the A96 area of Glasgow.

**Predicted propositions (blue) and entailment labels**

#1: A man has been taken to hospital after being hit by a car in the A96 area of Glasgow. ✓

#2: A man has been taken to hospital after being hit by a car in the A96 area of Glasgow. ✗

#3: A man has been taken to hospital after being hit by a car in the A96 area of Glasgow. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

A man has been taken to hospital after being hit by a car in the A96 area of Glasgow

---

**Summary from TransS2S**

A man has been taken to hospital after a two-vehicle crash on the A96 in County Antrim.

**Predicted propositions (blue) and entailment labels**

#1: A man has been taken to hospital after a two-vehicle crash on the A96 in County Antrim. ✓

#2: A man has been taken to hospital after a two-vehicle crash on the A96 in County Antrim. ✗

#3: A man has been taken to hospital after a two-vehicle crash on the A96 in County Antrim. ✗

#4: A man has been taken to hospital after a two-vehicle crash on the A96 in County Antrim. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

A man has been taken to hospital after a two-vehicle crash on the A96 in County Antrim.

---

Table 9: More example of model generated summaries on the XSum dataset, with human-annotated hallucination spans from Maynez et al. (2020). For each document, Maynez et al. (2020) provide summaries and hallucination annotations from 5 different summarization systems. We randomly sample documents and show our model’s predictions for all 5 summaries here.



---

**Document:** Dervite, 28, made 14 appearances last season to help Wanderers finish second in League One and secure promotion. The French centre-back joined Bolton from Charlton in 2014 and has made 83 appearances in all competitions. "Dorian was a bit of a forgotten man last year but came in and made an excellent contribution towards the end of the campaign," manager Phil Parkinson told the club website. Dervite follows David Wheater, Gary Madine and Jem Karacan in signing new contracts with Bolton, following their promotion to the Championship.

---

**Summary from BertS2S**

Bolton defender Dorian Dervite has signed a new two-year contract with the championship club.

**Predicted propositions (blue) and entailment labels**

#1: Bolton defender Dorian Dervite has signed a new two-year contract with the championship club. ✓

#2: Bolton defender Dorian Dervite has signed a new two-year contract with the championship club. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

Bolton defender Dorian Dervite has signed a new two-year contract with the championship club.

---

**Summary from TConvS2S**

Bolton Wanderers have signed defender Dorian Dervite from bolton wanderers for an undisclosed fee.

**Predicted propositions (blue) and entailment labels**

#1: Bolton Wanderers have signed defender Dorian Dervite from bolton wanderers for an undisclosed fee. ✗

#2: Bolton Wanderers have signed defender Dorian Dervite from bolton wanderers for an undisclosed fee. ✗

#3: Bolton Wanderers have signed defender Dorian Dervite from bolton wanderers for an undisclosed fee. ✗

#4: Bolton Wanderers have signed defender Dorian Dervite from bolton wanderers for an undisclosed fee. ✓

**Predicted hallucinated spans** (union of ✗- union of ✓)

Bolton Wanderers have signed defender Dorian Dervite from bolton wanderers for an undisclosed fee.

---

**Gold Summary from the XSum dataset**

Defender Dorian Dervite has signed a new one-year contract with Bolton.

**Predicted propositions (blue) and entailment labels**

#1: Defender Dorian Dervite has signed a new one-year contract with Bolton ✓

#2: Defender Dorian Dervite has signed a new one-year contract with Bolton. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

Defender Dorian Dervite has signed a new one-year contract with Bolton.

---

**Summary from PTGen**

Bolton Wanderers defender Dorian Dervite has signed a new three-and-a-half-year contract with the league one club until the end of the 2018-19 season.

**Predicted propositions (blue) and entailment labels**

#1: Bolton Wanderers defender Dorian Dervite has signed a new three-and-a-half-year contract with the league one club until the end of the 2018-19 season. ✓

#2: Bolton Wanderers defender Dorian Dervite has signed a new three-and-a-half-year contract with the league one club until the end of the 2018-19 season. ✗

#3: Bolton Wanderers defender Dorian Dervite has signed a new three-and-a-half-year contract with the league one club until the end of the 2018-19 season. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

Bolton Wanderers defender Dorian Dervite has signed a new three-and-a-half-year contract with the league one club until the end of the 2018-19 season.

---

**Summary from TransS2S**

Bolton Wanderers midfielder Gary Wheat has signed a new one-year contract with the championship side.

**Predicted propositions (blue) and entailment labels**

#1: Bolton Wanderers midfielder Gary Wheat has signed a new one-year contract with the championship side. ✗

#2: Bolton Wanderers midfielder Gary Wheat has signed a new one-year contract with the championship side. ✗

**Predicted hallucinated spans** (union of ✗- union of ✓)

Bolton Wanderers midfielder Gary Wheat has signed a new one-year contract with the championship side.

---

Table 10: (Cont.) More example of model generated summaries on the XSum dataset, with human-annotated hallucination spans from Maynez et al. (2020).