

LongtoNotes: OntoNotes with Longer Coreference Chains

Anonymous ACL submission

Abstract

001 Ontonotes has served as the most important
002 benchmark for coreference resolution. How-
003 ever, for ease of annotation, several long doc-
004 uments in Ontonotes were split into smaller
005 parts. In this work, we build a corpus of
006 coreference-annotated documents of signifi-
007 cantly longer length than what is currently
008 available. We do so by providing an accu-
009 rate, manually-curated, merging of annota-
010 tions from documents that were split into mul-
011 tiple parts in the original Ontonotes annota-
012 tion process (Pradhan et al., 2013). The result-
013 ing corpus, which we call LongtoNotes contains
014 documents in multiple genres of the English
015 language with varying lengths, the longest of
016 which are up to 8x the length of documents in
017 Ontonotes, and 2x those in Litbank. We evalu-
018 ate state-of-the-art neural coreference systems
019 on this new corpus, analyze the relationships
020 between model architectures/hyperparameters
021 and document length on performance and effi-
022 ciency of the models, and demonstrate areas
023 of improvement in long-document coreference
024 modelling revealed by our new corpus.

025 1 Introduction

026 Coreference resolution is an important prob-
027 lem in modelling discourse with applications in
028 knowledge-base construction (Luan et al., 2018),
029 question-answering (Reddy et al., 2019) and read-
030 ing assistants (Azab et al., 2013; Head et al., 2021).
031 In many such settings, the documents of interest,
032 are significantly longer and/or on wider varieties of
033 domain than the currently available corpora with
034 coreference annotations (Pradhan et al., 2013; Bam-
035 man et al., 2019; Mohan and Li, 2019; Cohen et al.,
036 2017).

037 The Ontonotes corpus (Pradhan et al., 2013) is
038 perhaps the most widely used benchmark for coref-
039 erence (Lee et al., 2013; Durrett and Klein, 2013;
040 Wiseman et al., 2016; Lee et al., 2017; Joshi et al.,
041 2020; Toshniwal et al., 2020b; Thirukovalluru et al.,

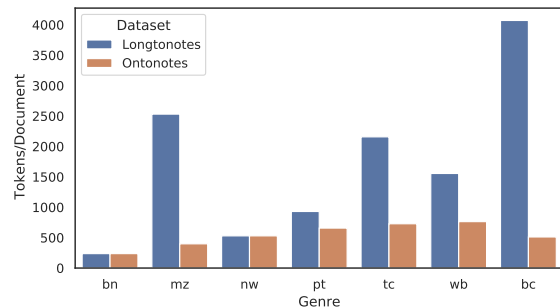


Figure 1: **Comparing Average Document Length.** Long documents in genres such as *broadcast conversations* (*bc*) were split into smaller parts in Ontonotes. Our proposed dataset, LongtoNotes, restores documents to their original form, revealing dramatic increases in length in certain genres.

2021; Kirstain et al., 2021). The construction process for Ontonotes, however, resulted in documents with artificially reduced length. For ease of annotation, longer documents were split into smaller parts and each part was annotated separately and treated as an independent document (Pradhan et al., 2013). The result is a corpus in which certain genres, such as *broadcast conversations* (*bc*), have greatly reduced length compared to their original form (Figure 1). As a result, the long, bursty spread of coreference chains in these documents is missing from the evaluation benchmark.

In this work, we present an extension to the Ontonotes corpus, called LongtoNotes. LongtoNotes combines coreference annotations in various parts of the same document, leading to a full document coreference annotation. This was done by our annotation team, which was carefully trained to follow the annotation guidelines laid out in the original Ontonotes corpus (Section 3). This led to a dataset where the average document length is over 0% longer than the standard OntoNotes benchmark and the average size of coreference chains increased by 25%. While other

066 datasets such as Litbank (Bamman et al., 2019) and
 067 CRAFT (Cohen et al., 2017) focus on long doc-
 068 uments in specialized domains, LongtoNotes
 069 comprises of documents in multiple genres (Ta-
 070 ble 1).

071 To illustrate the usefulness of LongtoNotes,
 072 we evaluate state-of-the-art coreference resolution
 073 models (Kirstain et al., 2021; Toshniwal et al.,
 074 2020b; Joshi et al., 2020) on the corpus and analyze
 075 the performance in terms of document length (§4.2).
 076 We illustrate how model architecture decisions and
 077 hyperparameters that support long-range dependen-
 078 cies have the greatest impact on coreference perfor-
 079 mance and importantly, these differences are only
 080 illustrated using LongtoNotes and are not seen
 081 in Ontonotes (§4.3). LongtoNotes also presents
 082 a challenge in scaling coreference models as pre-
 083 diction time and memory requirement increases
 084 substantially on the long documents (§4.4).

085 2 Our Contribution: LongtoNotes

086 We present LongtoNotes, a corpus that ex-
 087 tends the English coreference annotation in the
 088 OntoNotes Release 5.0 corpus¹ (Pradhan et al.,
 089 2013) to provide annotations for longer documents.
 090 In the original English OntoNotes corpus, the gen-
 091 res such as *broadcast conversations (bc)* and *tele-
 092 phone conversation (tc)* contain long documents
 093 that were divided into smaller parts to facilitate
 094 easier annotation. LongtoNotes is constructed
 095 by collecting annotations to combine within-part
 096 coreference chains into coreference chains over the
 097 entire long document. The annotation procedure,
 098 in which annotators merge coreference chains, is
 099 described and analyzed in Section 3.

100 The divided parts of a long document in
 101 Ontonotes are all assigned to the same partition
 102 (train/dev/test). This allows LongtoNotes to
 103 maintain the same train/dev/test partition, at the
 104 document level, as Ontonotes (Appendix, Table 11).
 105 The size of these partitions however does change as
 106 the divided parts are combined into a single anno-
 107 tated text in LongtoNotes. We will release the
 108 scripts to convert OntoNotes to LongtoNotes
 109 under Creative Commons 4.0 license and own-
 110 ing OntoNotes dataset is a prerequisite to run the
 111 scripts. We refer to LongtoNotes_s: Subset of
 112 LongtoNotes comprising only of long docu-
 113 ments (i.e. documents merged by the annotators).

¹The Arabic and Chinese parts of the Ontonotes dataset are not considered in our study.

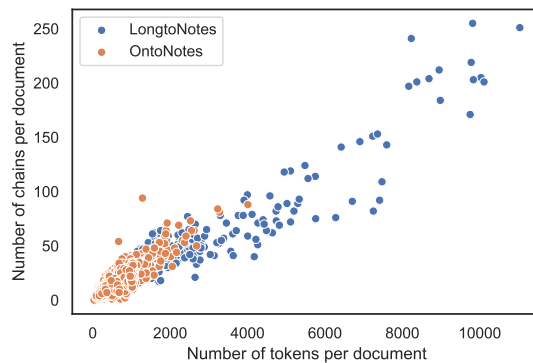


Figure 2: **Document and Coref Chain Length.** The number of coreference chains increases with the increase in token length in LongtoNotes.

2.1 Length of Documents in LongtoNotes

114 The average number of tokens per document
 115 (rounded to the nearest integer) in LongtoNotes
 116 is 674, 44% higher than in Ontonotes (466). Ta-
 117 ble 1 breaks down the changes in document length
 118 by genre. We observe that the genre with the
 119 longest documents is *broadcast conversation* with
 120 4071 tokens per document, which is a dramatic
 121 increase from the length of the divided parts in
 122 Ontonotes which had 511 tokens per document in
 123 the same. The number of coreference chains and
 124 the number of mentions per chain grows as well.
 125 The long documents that were split into multiple
 126 parts during the original OntoNotes annotation are
 127 not evenly distributed among the genres of text
 128 present in the corpus. In particular, text categories
 129 *broadcast news (bn)* and *newswire (nw)* consist ex-
 130 clusively of short non-split documents, which were
 131 not affected by the LongtoNotes merging process.
 132 A detailed distribution of what documents are
 133 merged in LongtoNotes is provided in Table 10
 134 in the Appendix.

2.2 Number of Coreference Chains

136 As a consequence of the increase in document
 137 length, LongtoNotes presents a higher number
 138 of coreference chains per document (16), compared
 139 to OntoNotes (12). Figure 2 shows the length and
 140 number of coreference chains for each document in
 141 the two corpora. As expected, the number of chains
 142 in a document tends to get larger as the document
 143 size increases.

144 For genres with longer average document lengths
 145 like *broadcast conversation (bc)*, the increase in
 146 the number of chains is as high as 85%, while this
 147

increase is only 25% for *pivot (pt)* genre when the document length is comparatively shorter. It is worth noting that the majority of documents had number of chains in the range of 20 to 50 and only about 20 documents out of 3493 in the OntoNotes dataset had >50 chains per document. For LongtoNotes the number increases to 96 documents. A comparison of the number of chains per document between OntoNotes and LongtoNotes is shown in Figure 3.

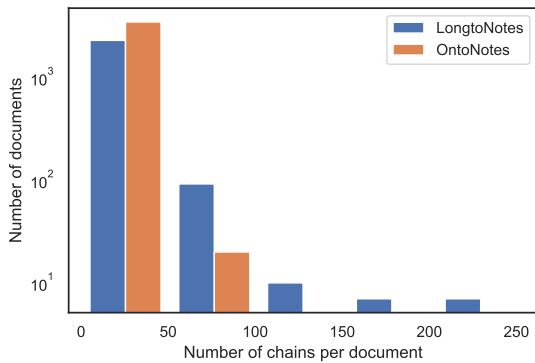


Figure 3: **Number of Chains per Document.** A histogram log plot reveals the long tailed distribution of the number of coreference chains present per document in LongtoNotes. Ontonotes contains more documents with fewer chains.

2.3 Number of Mentions per Chain

The number of mentions per coreference chain in LongtoNotes has gone up by over 30% compared to OntoNotes. This is primarily because of longer documents and an increase in the number of coreference chains per document. Mentions per chain increase with the increase in document length. For the *broadcast conversation (bc)* genre, the increase in the mentions per chain is highest with 87%, while for the *pivot (pt)* genre it is only 30% as it has shorter documents.

2.4 Distances to the Antecedents

For each coreference chain, we analyzed the distance between the mentions and their antecedents. The largest distance for a mention to its antecedent grew 3x for LongtoNotes dataset when compared to OntoNotes from 4885 to 11473 tokens. Figure 4 shows a detailed breakdown of the mention to antecedent distance. While there are no mentions that are more than 5K tokens distant from its antecedent in OntoNotes, there are 178 such mentions in LongtoNotes.

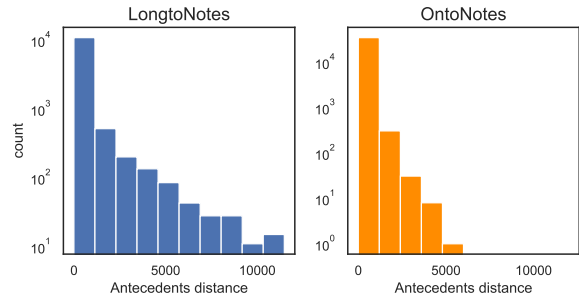


Figure 4: **Distance to Antecedent.** Histogram (log-scale) shows that the largest distance of mention to their antecedents per chain increases in LongtoNotes compared to OntoNotes.

2.5 Comparison with other Datasets

The literature contains multiple works proposing datasets for coreference resolution: Wiki coref (Ghaddar and Langlais, 2016), LitBank (Bamman et al., 2019), PreCo (Chen et al., 2018), Quiz Bowl Questions (Rodriguez et al., 2019; Guha et al., 2015), ACE corpus (Walker et al., 2006), MUC (Chinchor and Sundheim, 1995), MedMentions (Mohan and Li, 2019), inter alia. We compare LongtoNotes to these datasets in terms of number of documents, total number of tokens, and document length (Table 2).

Litbank, in particular, is a popular long document coreference dataset, presenting a high tokens/document ratio. However, the datasets consist of only 100 documents, rendering model development challenges. Moreover, it focuses only on the literary domain. Other datasets containing long documents (e.g., WikiCoref) are also very small in size. On the other hand, datasets consisting of a larger number of texts tend to contain shorter documents (e.g., PreCo). Thus, by proposing LongtoNotes, we address the scarcity of a multi-genre corpus with a collection of long documents containing long-range coreference dependencies.

3 Annotation Procedure & Quality Assurance

In this section, we describe the annotation procedure used to build LongtoNotes and assess the quality of the annotation.

3.1 Annotation Task

The annotations needed to build LongtoNotes are: antecedent labels for coreference chains in

Categories	# Docs		Tokens/Doc		# Chains		Ment./Chains	
	Ont.	Long.	Ont.	Long.	Ont.	Long.	Ont.	Long.
broadcast conversation (bc)	397	50	511	4071	14	85	65	519
broadcast news (bn)	947	947	237	237	8	8	29	29
magazine (mz)	494	78	398	2531	8	41	32	208
newswire (nw)	922	922	529	529	12	12	47	47
pivot (pt)	369	261	657	930	20	27	131	186
telephone conversation (tc)	142	48	728	2157	17	44	108	319
web data (wb)	222	109	763	1555	17	31	73	149
Overall	3493	2415	466	674	12	16	55	80

Table 1: **Genre Comparison.** Comparison of document and coreference chain statistics per genre in OntoNotes 5.0 and our proposed dataset, LongtoNotes.

Dataset	# Docs	Total Size	Tokens/Doc
WikiCoref	30	60K	2000
ACE-2007	599	300K	500
MUC-6	60	30K	500
MUC-7	50	25K	500
QuizBowl	400	50K	125
PreCo	37.6K	12.4M	330
LitBank	100	200K	2105
MedMentions	4392	1.1M	267
OntoNotes	3493	1.6M	466
LongtoNotes	2415	1.6M	674
LongtoNotes _s	283	740K	2615

Table 2: **Coreference Datasets.** A comparison of various coref datasets with our proposed dataset LongtoNotes.

part $i + 1$ of a document to chains present in parts $1, \dots, i$. We reformulate this annotation process as a question answering task where we ask annotators a series of questions using our own annotation tool designed for this task (Appendix, Figure 8). We display parts $1, \dots, i$ with color-coded mention spans. We then show a highlighted concept from part $i + 1$ and ask the question: *The highlighted concept below refers to which concept in the above paragraphs?* The annotators select one of the colour-coded chains from parts $1, \dots, i$ from a list of answers or the annotators can specify that the highlighted concept in part $i + 1$ does not refer to any concept in parts $1, \dots, i$, (i.e., a new concept emerging in part $i + 1$).

The annotation tool proceeds with a question for each coreference chain ordered (sorted by the first token offset of the first mention in the chain). After answering questions for all chains in one part of the document, the annotators are presented with a summary of their annotations and allowed to confirm/change their responses.

The annotation of all parts of a document com-

prises an annotation task. That is, a single annotator is tasked with answering the multiple-choice question for each coreference chain in each part of a document.

From Annotations to Coreference Labels The annotations collected through this task are then converted into coreference labels for the merged parts of a document. The answers to the questions tell us the antecedent link between two coreference chains. These links are used to relabel all mentions in the two chains with the same coreference label, resulting in the LongtoNotes dataset.

3.2 Annotators and Training

We hired and trained a team of three annotators for the aforementioned task. The annotators were university-level English majors from India and were closely supervised by an expert with experience in similar annotation projects. The annotation team was paid a fair wage for the work. We had several hour-long training sessions outlining the annotation task, setup of the problem, and Ontonotes annotation guidelines. We reviewed example cases of difficult annotation decisions and collaboratively worked through example annotations. We then ran a pilot annotation study with a small number of documents. For these documents, the authors of this paper also provided annotation. We then reviewed the annotator’s work on these documents and discussed disagreements with the annotators and asked them to re-annotate the documents.

After the pilot annotation study, the tasks were assigned to the annotators in five batches of 60 tasks each. For 10% of the tasks, we had all three annotators provide annotations. For the remaining 90%, a single annotator was used. For the docu-

ments with multiple annotators, we used majority voting to settle disagreements. If all annotators disagreed on a specific case, we selected Annotator 1’s decision over the others (analysis in the Appendix).

3.3 Measuring Quality of Annotation

We would like to ensure that LongtoNotes meets high-quality standards. To do this, we define metrics of agreement between a pair of annotators. We consider (1) the question-answering agreement (i.e., how similar are the annotations made using the annotation tool), and (2) the coreference label agreement (i.e., at the level of the resulting coreference annotation).

We consider the following question answering metrics: *Each annotator receives a set of chains C_1, C_2, \dots, C_N . For each chain C_i , the annotator links it to a New chain or a chain from their (annotator specific) set of available chains. Let us call D_i this linking decision, which consists of a pair (C_i, A_i) , where A_i is the selected antecedent chain.*

- **Strict Decision Matching:** When two annotators agreed on merging two chains and there is an exact match between the merged chains. Calculated as $\frac{1}{N} \sum_i D_i^{(1)} = D_i^{(2)}$
- **Jaccard Decision Match:** Jaccard decision calculated as $\frac{1}{N} \sum_i \frac{(D_i^{(1)}.A_i^{(1)}) \cap (D_i^{(2)}.A_i^{(2)})}{(D_i^{(1)}.A_i^{(1)}) \cup (D_i^{(2)}.A_i^{(2)})}$
- **New Chain Agreement:** Number of times two annotators agreed on new chain choice divided by number of times at least one annotator labels New chain.
- **Not New Chain Agreement:** Pairwise agreement between annotators when the chain choice is not a New chain.

Table 3 presents the results for these metrics. We observed that on average annotators agreed with each other on over 90% of their decisions except when the No New chains were considered. Removing New chains reduces the total decisions to be made significantly, and hence a lower score on No New chains agreement. In general, Annotator 1 and 2 agree with each other more than with Annotator 3 (+1 – 2%). We found that Annotator 1 agreed most with the experts and hence Annotator 1’s decisions were preferred over the others in case of disagreement between all three annotators.

Metric	Comparison	Score
Strict Match	Authors	0.98
Strict Match	Each other	0.90
Jaccard Match	Authors	0.99
Jaccard Match	Each other	0.95
New Chain	Authors	0.96
New Chain	Each other	0.88
Not New Chain	Authors	0.92
Not New Chain	Each other	0.87

Table 3: **Quality Assessment of Annotation.** We report the average value of each metric over all pairs of annotators and the annotators and authors of this paper.

Where are disagreements found in annotation?

We would like to understand what kinds of mentions lead to the disagreement between annotators. To investigate this, we measure the part of speech of all the disagreed chain assignments between the annotators. We found that the 8% of the mentions within the disagreed chain assignments were pronouns, 8% were verbs, and 9% were common nouns. The number of proper nouns disagreements was lower with just 5%. When considering different genres, it was observed that genres with longer documents like *broadcast conversation (bc)* had more mentions that were pronouns when compared with genres with shorter documents *pivot (pt)*. As expected, the number of disagreements in general increased with the size of the documents. However, we found that the number of disagreements was manageably small even for long document genres such as *broadcast conversation (bc)*. A more comprehensive overlook is presented in the Appendix.

3.4 Time Taken per Annotation

We also recorded the time for each annotation. Figure 5 shows that the time taken per annotation increases with the increase in the document length. This is expected as more chains create more options to be chosen from and longer document length demands more reading and attention. In total, our annotation process took 400 hours.

4 Empirical Analysis with LongtoNotes

We hope to show that LongtoNotes can facilitate the empirical analysis of coreference models in ways that were not possible with the original OntoNotes. We are interested in the fol-

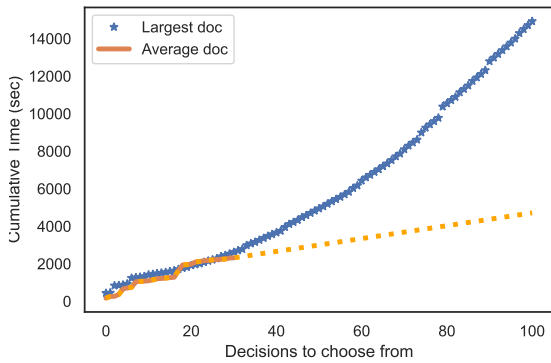


Figure 5: **Annotation Time and Document Length.** Annotation time (cumulative) increases exponentially with the increase in the number of decisions to choose from. A comparison is shown between the longest document in LongtoNotes vs an average document. The dotted lines represent the increase in annotation time if the growth was linear.

lowing empirical questions using the datasets—Ontonotes (Pradhan et al., 2013), and our proposed LongtoNotes and LongtoNotes_s:

- How does the length of documents play a role in the empirical performance of models?
- Does the empirical accuracy of models depend on different hyperparameters in LongtoNotes and Ontonotes?
- Does LongtoNotes reveal properties about the efficiency/scalability of models not present in Ontonotes?

4.1 Models

Much of the recent work on coreference can be organized into three categories: span based representations (Lee et al., 2017; Joshi et al., 2020), token-wise representations (Thirukovalluru et al., 2021; Kirstain et al., 2021) and memory networks / incremental models (Toshniwal et al., 2020b,a). We consider one approach from all three categories.

Span based representation We used the Joshi et al. (2020) implementation of the higher-order coref resolution model (Lee et al., 2018) with SpanBERT. Here, the documents were divided into a non-overlapping segment length of 384 tokens. We used SpanBERT Base as our model due to memory constraints. The number of training sentences was set to 3. We set the maximum top antecedents, $K = 50$. We used Adam (Kingma and Ba, 2014) as our optimiser with a learning rate of $2e^{-4}$.

# Tokens	Training	CoNLL F1
$\leq 2K$	Ontonotes	78.85
	LongtoNotes	78.25
$> 2K$	Ontonotes	65.11
	LongtoNotes	66.20

Table 4: **Performance and Document Length for Span-based Models.** F_1 score across different document length for SpanBERT Base trained model on OntoNotes and LongtoNotes dataset.

Token-wise representation We used the LongFormer Large (Beltagy et al., 2020) version of Kirstain et al. (2021) work, as this approach is less memory demanding and it is possible to fit this model in our memory. The maximum sequence length was set to 384 or 4096. Adam was used as an optimiser with a learning rate of $1e^{-5}$. A dropout (Srivastava et al., 2014) probability of 0.3 was used.

Memory networks We used SpanBERT Large with a sequence length of 512 tokens. As in their work, an endpoint-based mention detector was trained first and then was used for coreference resolution. The number of training sentences was set to 5, 10, and 20. The number of memory cells was selected from 20 or 40. All experiments were performed with AutoMemory models with learned memory type.

4.2 Length of Documents & Performance

Impact of Training Corpus We first investigate whether or not training on the longer documents in LongtoNotes are needed to achieve state-of-the-art results on the dataset. We compare the performance of models trained on Ontonotes to those trained on LongtoNotes. We find that by training on LongtoNotes, we can achieve higher CoNLL F1 measures on LongtoNotes than training with Ontonotes for each model architecture (Table 5). This suggests that the longer dependencies formed by merging annotations in various parts of documents in OntoNotes are difficult to model when training on short documents.

We find that to achieve accuracy with hyperparameters such as learning rate/warmup size, we need to maintain a number of steps per epoch consistent with Ontonotes when training with LongtoNotes. A detailed analysis is presented in the Appendix Section 8.

	Training	OntoNotes			LongtoNotes _s			LongtoNotes		
		P	R	F ₁	P	R	F ₁	P	R	F ₁
Span-based (Joshi et al., 2020)	OntoNotes	76.5	77.6	77.4	72.7	69.1	70.8	74.4	73.0	73.7
	LongtoNotes	75.9	77.7	76.8	72.4	70.7	71.5	73.9	74.1	74.0
Token-Level (Kirstain et al., 2021)	OntoNotes	81.2	79.5	80.4	79.6	80.0	79.8	79.7	77.2	78.5
	LongtoNotes	80.0	78.2	79.1	80.3	80.3	80.3	80.2	78.0	79.1
Memory-Model (Toshniwal et al., 2020b)	OntoNotes	73.5	79.3	76.4	63.4	73.8	68.2	67.9	76.6	72.0
	LongtoNotes	73.8	79.4	76.6	66.3	74.6	70.2	69.3	77.0	72.9

Table 5: **Performance Variation by Training Set.** Comparison of F_1 scores on various datasets using different models. **All experiments have been performed atleast 2 times and a variance of only ± 0.1 was observed.**

Length Analysis - Number of Tokens We break down the performance of the Span-based model by the number of tokens in each document. We compare the performance of the model depending on the training set. Figure 2 shows that the majority of the documents in the OntoNotes dataset falls within a token length of 2000 per document. We create two splits of LongtoNotes_s, one having a token length greater than 2000 tokens, the other having a number of tokens smaller than 2000. Table 4 shows that for smaller document length (less than 2000 tokens), the SpanBERT model trained on OntoNotes performed better but the trend reverses for longer documents (more than 2000 tokens), on which the model trained on LongtoNotes outperformed the model trained on OntoNotes by +1%.

Length Analysis - Number of Clusters Table 6 displays the change in F_1 score with the increase in the number of clusters per document. The SpanBERT Base model trained on LongtoNotes outperforms the same model trained on OntoNotes (+0.6%) when the number of clusters is more than 40. Note that, 40 is selected based on the cluster distribution shown in Table 1 with the majority documents in LongtoNotes lying in this range.

4.3 Hyperparameters & Document Length

Each model has a set of hyperparameters that would seemingly lead to variation in performance with respect to document length. We consider the performance of the models on LongtoNotes as a function of these hyperparameters.

Span-based model hyperparameters We consider two hyperparameters: the number of antecedents to use, K and the max number of sentences used in each training example. We found that upon varying K : 10, 25 and 50, there was

# Chains	Training	CoNLL F1
≤ 40	OntoNotes	73.60
	LongtoNotes	72.86
> 40	OntoNotes	68.44
	LongtoNotes	69.09

Table 6: **Performance and Number of Chains for Span-based Models.** F_1 score across different document length for SpanBERT Base trained model on OntoNotes and LongtoNotes dataset.

only a small difference observed in the results for both the models trained on OntoNotes and LongtoNotes (increasing K led to only minor increases). The result is summarized in Table 7. We could not go beyond $K = 50$ due to our GPU memory limitations. However, going beyond 50 might further help for longer documents. Furthermore, we found that the *number of sentences* parameter used to create training batches does not play a significant role in performance either (Figure 6).

K	OntoNotes	LongtoNotes	LongtoNotes _s
10	77.05	73.44	70.37
25	76.93	73.99	71.61
50	77.60	74.01	71.58

Table 7: **Number of Antecedents vs. Performance** SpanBERT Base model trained on LongtoNotes dataset with varying K value.

Token-wise model hyperparameters We experimented with reducing the sequence length when testing from 4096 to 384 and we observe a drop in performance. Figure 7 shows the effect on performance due to the change in the sequence length. We observed that longer sequence length (4096) helps more for LongtoNotes_s as there are longer sequences than for OntoNotes, which

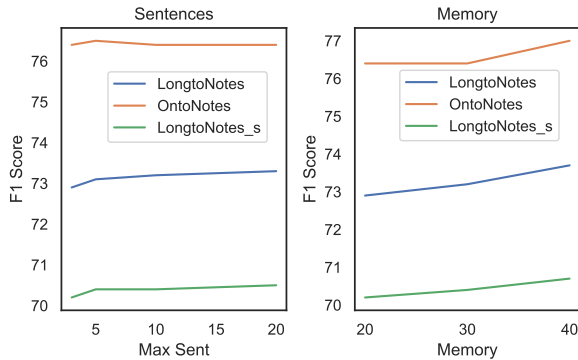


Figure 6: **Max Sentence Length.** Increasing max sentences from 3 to 20 has small effect on the performance of the SpanBERT large model. On the other hand, the increase is linear with the increase in the memory size alongside with the increase in max training sentences.

is evident in Figure 7. Furthermore, we analyzed the effect of sequence length on two genres: *magazine (mz)* having 6x longer sequences in LongtoNotes than OntoNotes vs *pivot (pt)* having just 1.4x longer documents. As observed in Figure 10, when the document is long as in *magazine (mz)*, there is a significant increase in performance with a longer sequence but the effect is negligible for *pivot (pt)* where the size of the document is almost the same. A detailed comparison is provided in the Appendix Table 15.

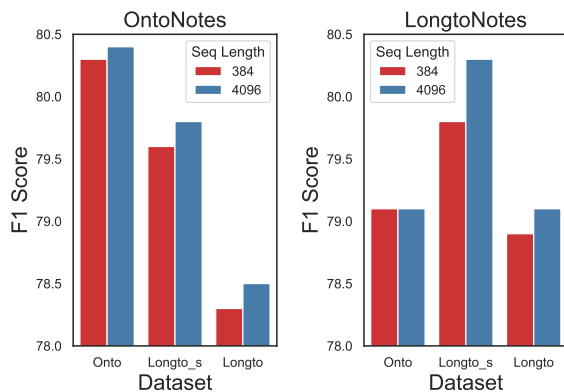


Figure 7: **Sequence Length vs. Performance.** LongFormer is significantly better on LongtoNotes with 4096 sequence length compared to 384. Two sequence lengths perform similarly on Ontonotes.

Memory model hyperparameters We consider two hyperparameters - the memory size which denotes the maximum active antecedents that can be considered and the max number of sentences used in training. We show that doubling the size of

the memory leads to an increase of 0.8 points of CoNLL F1 for LongtoNotes dataset. (Table 8). Figure 6 demonstrates that there is no significant improvement in the performance of the model with the increase in the number of training sentences.

Dataset	Memory Size	
	20	40
OntoNotes	76.6	77.0
LongtoNotes	72.9	73.7
LongtoNotes _s	70.2	70.7

Table 8: **Memory Size vs. Performance.** We compare two settings of the memory size parameter in memory model (Toshniwal et al., 2020b) and find that the larger memory version achieves better results on each dataset.

4.4 Model Efficiency

We compare the prediction time for the span-based model on the longest length and average length documents in LongtoNotes and Ontonotes in Table 9. We observe that there is a significant jump in running time and memory required to scale the model to long documents on LongtoNotes; this jump is much smaller on Ontonotes. This suggests that our proposed dataset is better suited for assessing the scaling properties of coreference methods.

Dataset	Type	Pred. Time	Pred. Mem
Ontonotes	Average	0.11 sec	1.50 GB
LongtoNotes	Average	0.47 sec	6.50 GB
Ontonotes	Longest	0.37 sec	5.84 GB
LongtoNotes	Longest	2.35 sec	42.68 GB

Table 9: **Model Efficiency of Span-based Models.** We find that LongtoNotes documents have extended length leading to greater variation of prediction time and prediction memory.

5 Conclusion

In this paper, we introduced LongtoNotes, a dataset that merges the coreference annotation of documents that in the original OntoNotes dataset were split into multiple independently-annotated parts. LongtoNotes has longer documents and coreference chains than the original OntoNotes dataset. Using LongtoNotes, we demonstrate that scaling current approaches to long documents has significant challenges both in terms of achieving a better performance as well as scalability. We demonstrate the merits of using LongtoNotes as an evaluation benchmark for coreference resolution and encourage future work to do so.

Ethical Considerations

Our dataset is comprised solely of English texts, and our analysis, therefore, applies uniquely to the English language. The annotation was performed with a data annotation service which ensured that the annotators were paid fair compensation. The annotation process did not solicit any sensitive information from the annotators. Finally, while our models are not tuned for any specific real-world application, the methods could be used in sensitive contexts such as legal or health-care settings, and any work must use our methods undertake extensive quality-assurance and robustness testing before using them in their setting.

Replicability: As part of our contributions, we will release the models trained on `LongtoNotes` discussed in this manuscript.

References

- Mahmoud Azab, Ahmed Salama, Kemal Oflazer, Hideki Shima, Jun Araki, and Teruko Mitamura. 2013. [An NLP-based reading tool for aiding non-native English readers](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*, pages 41–48, Hissar, Bulgaria. INCOMA Ltd. Shoumen, BULGARIA.
- Amit Bagga and Breck Baldwin. 1998. Algorithms for scoring coreference chains. In *The First International Conference on Language Resources and Evaluation Workshop on Linguistics Coreference*, pages 563–566.
- David Bamman, Sejal Popat, and Sheng Shen. 2019. [An annotated dataset of literary entities](#). 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 2138–2144, Minneapolis, Minnesota. Association for Computational Linguistics.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Hong Chen, Zhenhua Fan, Hao Lu, Alan L Yuille, and Shu Rong. 2018. [Preco: A large-scale dataset in preschool vocabulary for coreference resolution](#). *arXiv preprint arXiv:1810.09807*.
- Nancy A Chinchor and Beth Sundheim. 1995. Message understanding conference (muc) tests of discourse processing. In *Proc. AAAI Spring Symposium on Empirical Methods in Discourse Interpretation and Generation*, pages 21–26.

- K Bretonnel Cohen, Arrick Lanfranchi, Miji Joo-young Choi, Michael Bada, William A Baumgartner, Natalya Panteleyeva, Karin Verspoor, Martha Palmer, and Lawrence E Hunter. 2017. Coreference annotation and resolution in the colorado richly annotated full text (craft) corpus of biomedical journal articles. *BMC bioinformatics*, 18(1):1–14.
- Greg Durrett and Dan Klein. 2013. Easy victories and uphill battles in coreference resolution. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1971–1982.
- Abbas Ghaddar and Phillippe Langlais. 2016. [Wiki-Coref: An English coreference-annotated corpus of Wikipedia articles](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 136–142, Portorož, Slovenia. European Language Resources Association (ELRA).
- Anupam Guha, Mohit Iyyer, Danny Bouman, and Jordan Boyd-Graber. 2015. [Removing the training wheels: A coreference dataset that entertains humans and challenges computers](#). In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1108–1118, Denver, Colorado. Association for Computational Linguistics.
- Andrew Head, Kyle Lo, Dongyeop Kang, Raymond Fok, Sam Skjonsberg, Daniel S. Weld, and Marti A. Hearst. 2021. [Augmenting Scientific Papers with Just-in-Time, Position-Sensitive Definitions of Terms and Symbols](#). Association for Computing Machinery, New York, NY, USA.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. 2020. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8:64–77.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Yuval Kirstain, Ori Ram, and Omer Levy. 2021. Coreference resolution without span representations. In *ACL/IJCNLP*.
- Heeyoung Lee, Angel Chang, Yves Peirsman, Nathanael Chambers, Mihai Surdeanu, and Dan Jurafsky. 2013. Deterministic coreference resolution based on entity-centric, precision-ranked rules. *Computational linguistics*, 39(4):885–916.
- Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. [End-to-end neural coreference resolution](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 188–197, Copenhagen, Denmark. Association for Computational Linguistics.

623	Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018.	Marc Vilain, John Burger, John Aberdeen, Dennis Con-	678
624	Higher-order coreference resolution with coarse-to-	nolly, and Lynette Hirschman. 1995. A model-	679
625	fine inference. <i>arXiv preprint arXiv:1804.05392</i> .	theoretic coreference scoring scheme . In <i>Proceed-</i>	680
626	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh	<i>ings of the 6th Conference on Message Understand-</i>	681
627	Hajishirzi. 2018. Multi-task identification of enti-	<i>ing</i> , MUC6 '95, page 45–52, USA. Association for	682
628	tities, relations, and coreference for scientific knowl-	Computational Linguistics.	683
629	edge graph construction. In <i>Proceedings of the 2018</i>	Christopher Walker, Stephanie Strassel, Julie Medero,	684
630	<i>Conference on Empirical Methods in Natural Lan-</i>	and Kazuaki Maeda. 2006. Ace 2005 multilin-	685
631	<i>guage Processing</i> .	gual training corpus. <i>Linguistic Data Consortium,</i>	686
632	Xiaoqiang Luo. 2005. On coreference resolution per-	<i>Philadelphia</i> , 57:45.	687
633	formance metrics . HLT '05, page 25–32, USA. As-	Sam Wiseman, Alexander M. Rush, and Stuart M.	688
634	sociation for Computational Linguistics.	Shieber. 2016. Learning global features for coref-	689
635	Sunil Mohan and Donghui Li. 2019. Medmentions: a	erence resolution . In <i>Proceedings of the 2016 Con-</i>	690
636	large biomedical corpus annotated with umls con-	<i>ference of the North American Chapter of the Asso-</i>	691
637	cepts. <i>arXiv preprint arXiv:1902.09476</i> .	<i>ciation for Computational Linguistics: Human Lan-</i>	692
638	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue,	<i>guage Technologies</i> , pages 994–1004, San Diego,	693
639	Hwee Tou Ng, Anders Björkelund, Olga Uryupina,	California. Association for Computational Linguis-	694
640	Yuchen Zhang, and Zhi Zhong. 2013. Towards ro-	tics.	695
641	bust linguistic analysis using OntoNotes . In <i>Pro-</i>		
642	<i>ceedings of the Seventeenth Conference on Computa-</i>		
643	<i>tional Natural Language Learning</i> , pages 143–152,		
644	Sofia, Bulgaria. Association for Computational Lin-		
645	guistics.		
646	Siva Reddy, Danqi Chen, and Christopher D. Manning.		
647	2019. CoQA: A conversational question answering		
648	challenge . <i>Transactions of the Association for Com-</i>		
649	<i>putational Linguistics</i> , 7:249–266.		
650	Pedro Rodriguez, Shi Feng, Mohit Iyyer, He He, and		
651	Jordan Boyd-Graber. 2019. Quizbowl: The case		
652	for incremental question answering. <i>arXiv preprint</i>		
653	<i>arXiv:1904.04792</i> .		
654	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky,		
655	Ilya Sutskever, and Ruslan Salakhutdinov. 2014.		
656	Dropout: A simple way to prevent neural networks		
657	from overfitting . <i>Journal of Machine Learning Re-</i>		
658	<i>search</i> , 15(56):1929–1958.		
659	Raghuveer Thirukovalluru, Nicholas Monath, Kumar		
660	Shridhar, Manzil Zaheer, Mrinmaya Sachan, and An-		
661	drew McCallum. 2021. Scaling within document		
662	coreference to long texts . In <i>Findings of the Associ-</i>		
663	<i>ation for Computational Linguistics: ACL-IJCNLP</i>		
664	<i>2021</i> , pages 3921–3931, Online. Association for		
665	Computational Linguistics.		
666	Shubham Toshniwal, Allyson Ettinger, Kevin Gimpel,		
667	and Karen Livescu. 2020a. Petra: A sparsely su-		
668	pervised memory model for people tracking. <i>arXiv</i>		
669	<i>preprint arXiv:2005.02990</i> .		
670	Shubham Toshniwal, Sam Wiseman, Allyson Ettinger,		
671	Karen Livescu, and Kevin Gimpel. 2020b. Learn-		
672	ing to Ignore: Long Document Coreference with		
673	Bounded Memory Neural Networks . In <i>Proceed-</i>		
674	<i>ings of the 2020 Conference on Empirical Methods</i>		
675	<i>in Natural Language Processing (EMNLP)</i> , pages		
676	8519–8526, Online. Association for Computational		
677	Linguistics.		

6 Appendix

6.1 Annotation tool

Figure 8 shows the annotation tool built by us.

6.2 Comparison with OntoNotes

A detailed genre-wise comparison of the documents from OntoNotes dataset which were merged in LongtoNotes is presented in Table 10. It can be seen that categories like bn and nw are completely missing in LongtoNotes, while pt is partially missing.

Documents in Corpus comparison		
Category	Onto	Longto
bc/cctv	✓	✓
bc/cnn	✓	✓
bc/msnbc	✓	✓
bc/phoenix	✓	✓
bn/abc	✓	✗
bn/cnn	✓	✗
bn/mnb	✓	✗
bn/nbc	✓	✗
bn/pri	✓	✗
bn/voa	✓	✗
mz/sinorama	✓	✓
nw/wsj	✓	✗
nw/xinhua	✓	✗
pt/nt	✓	✓
pt/ot	✓	✗
tc/ch	✓	✓
wb/a2e	✓	✓
wb/c2e	✓	✓
wb/eng	✓	✓

Table 10: Comparison of documents from various sub-categories that exists in OntoNotes 5.0 and our proposed dataset LongtoNotes

7 Train test dev split

A comparison between the number of documents in the train-test-dev split between LongtoNotes and OntoNotes is provided in the Table 11.

Dataset	Train	Dev	Test
OntoNotes	2802	343	348
LongtoNotes	1959	234	222

Table 11: Comparison of train-test-dev split of documents between OntoNotes and LongtoNotes

7.1 Genre wise disagreement analysis

Table 12 presents the genre wise disagreement analysis for strict decision matching. Genres with longer documents like bc, mz have more disagreements compared to genres with smaller document length like tc, pt.

The trend is very similar for new chain assignments where genres with larger documents have more disagreements over new chain assignments. The numbers are presented in Table 14.

bc			
	Ann1	Ann2	Ann3
Ann1	1.0	0.91	0.87
Ann2	0.91	1.0	0.88
Ann3	0.87	0.88	1.0

mz			
	Ann1	Ann2	Ann3
Ann1	1.0	0.91	0.94
Ann2	0.91	1.0	0.93
Ann3	0.94	0.93	1.0

pt			
	Ann1	Ann2	Ann3
Ann1	1.0	0.97	0.98
Ann2	0.97	1.0	0.96
Ann3	0.98	0.96	1.0

tc			
	Ann1	Ann2	Ann3
Ann1	1.0	0.99	0.98
Ann2	0.99	1.0	0.98
Ann3	0.98	0.98	1.0

wb			
	Ann1	Ann2	Ann3
Ann1	1.0	0.93	0.90
Ann2	0.93	1.0	0.92
Ann3	0.90	0.92	1.0

Table 12: Genre wise strict decision based disagreement analysis between the annotators.

7.2 Annotators disagreements analysis

Figure 9 shows the cases (in black) when the annotators disagreed for each part of speech categories (shown in big colored bubbles). The size of the bubbles are representative of their occurrence in

Figure 8: The tool designed by us for the annotation task. Upper box represents all the previous paragraphs while the box on the bottom left is the current paragraph. The mentions of the current chain to be merged are shown in yellow. On the right side, the answers are presented which are chains from previous paragraphs and the annotator can select one of them or choose the None of the below option which creates a new chain.

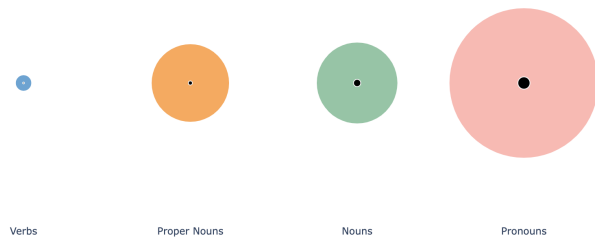
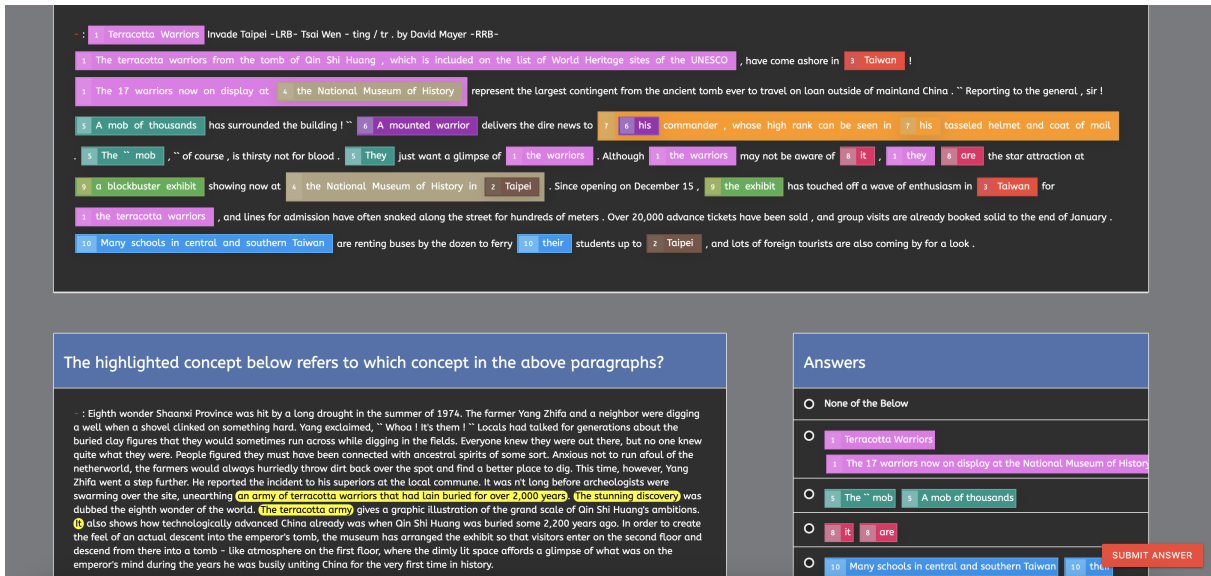


Figure 9: Plot showing the part of speech distribution for the disagreed clusters between annotators.

the dataset, suggesting there are more pronominal mentions in the dataset than nouns or proper nouns.

7.2.1 Genre wise disagreement analysis

In general, annotators disagree more on pronouns than proper nouns and the trend is consistent for various genres as shown in Table 13.

PoS type	bc	pt
Pronouns	3.6	0.04
Nouns	3.2	0.05
Proper Nouns	1.9	0.03
Verbs	3.5	1.0

Table 13: Genre wise part of speech comparison for two genres: bc and pt. The numbers are normalized and presented in percentage.

8 Results

8.1 MUC, B^3 and CEAFE scores

Tables 16, 17 and 18 present the MUC (Vilain et al., 1995), B^3 (Bagga and Baldwin, 1998) and CEAFE (Luo, 2005) scores for SpanBERT Base (Lee et al., 2017) and LongDocCoref Models (Toshniwal et al., 2020b). On all three metrics, both models trained on LongtoNotes dataset outperforms the models trained on OntoNotes dataset. For SpanBERT base model, we compare three version of the LongtoNotes dataset: LongtoNotes_s and LongtoNotes dataset as mentioned in the paper and LongtoNotes_{eq} where LongtoNotes dataset is reweighted to create the total number of documents equal to the number of documents in OntoNotes dataset. For LongDocCoref model, n represents the maximum number of training sentences, while m refers to the memory used.

8.2 Genre wise F_1 scores vs sequence length

Table 15 shows that LongFormer Large model with larger sequence length (4096) outperforms the one with shorter sequence length (384) for all models. The difference is higher when the documents are longer (as seen in mz genre) than when the documents are shorter (as seen in pt).

bc			
	Ann1	Ann2	Ann3
Ann1	1.0	0.91	0.85
Ann2	0.91	1.0	0.86
Ann3	0.85	0.86	1.0

mz			
	Ann1	Ann2	Ann3
Ann1	1.0	0.89	0.91
Ann2	0.89	1.0	0.90
Ann3	0.91	0.90	1.0

pt			
	Ann1	Ann2	Ann3
Ann1	1.0	0.94	0.95
Ann2	0.94	1.0	0.91
Ann3	0.95	0.91	1.0

tc			
	Ann1	Ann2	Ann3
Ann1	1.0	0.98	0.98
Ann2	0.98	1.0	0.98
Ann3	0.98	0.98	1.0

wb			
	Ann1	Ann2	Ann3
Ann1	1.0	0.92	0.90
Ann2	0.92	1.0	0.91
Ann3	0.90	0.91	1.0

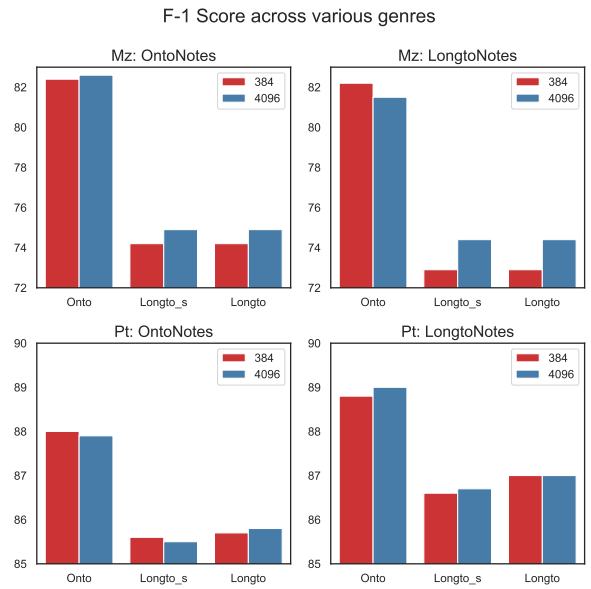


Figure 10: Plot comparing the sequence length effect on performance for two genres: *magazine (mz)* and *pivot (pt)*.

Table 14: Genre wise disagreement analysis between the annotators for new chain assignment.

	OntoNotes						LongtoNotes _s						LongtoNotes					
	Mention			Coref			Mention			Coref			Mention			Coref		
	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁
LongFormer Large (mz)																		
+ OntoNotes (384)	88.0	87.9	88.0	82.4	82.4	82.4	84.3	86.1	85.2	73.8	75.0	74.2	84.3	86.1	85.2	73.8	75.0	74.2
+ OntoNotes (4096)	87.9	88.3	88.1	82.4	82.9	82.6	84.4	86.7	85.5	74.1	75.9	74.9	84.4	86.7	85.5	74.1	75.9	74.9
+ LongtoNotes (384)	87.0	88.4	87.7	81.4	83.0	82.2	84.4	86.9	85.6	72.4	73.6	72.9	84.4	86.9	85.6	72.4	73.6	72.9
+ LongtoNotes (4096)	86.9	87.8	87.4	80.9	82.0	81.5	85.0	86.7	85.8	74.1	74.8	74.4	85.0	86.7	85.8	74.1	74.8	74.4
LongFormer Large (pt)																		
+ OntoNotes (384)	95.5	94.4	95.0	88.6	87.4	88.0	94.3	95.3	94.8	84.6	86.9	85.7	94.9	94.4	94.7	85.5	85.8	85.6
+ OntoNotes (4096)	95.6	94.2	94.9	88.9	86.9	87.9	94.4	94.8	94.6	84.8	86.8	85.8	94.9	94.0	94.5	85.5	85.2	85.5
+ LongtoNotes (384)	95.1	94.3	94.7	89.2	88.3	88.8	94.2	95.1	94.6	86.0	88.0	87.0	94.6	94.2	94.4	86.5	86.7	86.6
+ LongtoNotes (4096)	95.3	94.2	94.8	89.7	88.2	89.0	94.5	94.5	94.5	86.4	87.4	86.9	94.8	93.7	94.3	87.0	86.4	86.7

Table 15: Comparison of F_1 scores for mz and pt genres.

	OntoNotes						LongtoNotes _s						LongtoNotes					
	Mention			Coref			Mention			Coref			Mention			Coref		
	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁
SpanBERT Base (Lee et al., 2017)																		
+ OntoNotes	86.6	87.5	87.0	83.1	83.6	83.4	88.4	85.0	86.7	84.2	80.8	82.4	86.7	85.4	86.1	83.0	81.3	82.1
+ LongtoNotes _s	73.3	91.0	81.2	70.0	85.7	77.1	78.3	90.5	84.0	73.8	85.5	79.2	73.2	90.4	80.9	69.4	85.1	76.5
+ LongtoNotes	86.6	87.1	86.8	83.0	82.9	86.8	88.1	84.6	86.3	83.3	80.1	81.7	86.6	85.5	86.0	82.4	81.0	81.7
+ LongtoNotes _{eq}	86.1	87.8	87.0	82.8	83.5	83.2	87.7	86.2	87.0	83.4	81.9	82.6	86.1	86.3	86.2	82.3	81.9	82.1
LongDocCoref (Toshniwal et al., 2020b)																		
+ OntoNotes	95.3	85.6	86.4	81.2	85.4	83.2	95.3	85.6	86.4	77.8	86.2	81.8	95.3	85.6	86.4	78.2	85.2	81.6
+ LongtoNotes _s	95.3	85.6	86.4	22.3	66.9	33.5	95.3	85.6	86.4	17.5	65.7	27.6	95.3	85.6	86.4	21.7	66.9	32.8
+ LongtoNotes	95.3	85.6	86.4	81.4	85.0	83.2	95.3	85.6	86.4	79.3	85.8	82.4	95.3	85.6	86.4	79.1	85.0	81.9
+ LongtoNotes _{eq} (n=3)	95.3	85.6	86.4	81.6	85.2	83.4	95.3	85.6	86.4	79.7	86.2	82.8	95.3	85.6	86.4	79.3	85.2	82.2
+ LongtoNotes _{eq} (n=5)	95.3	85.6	86.4	81.4	85.3	83.3	95.3	85.6	86.4	79.7	86.2	82.8	95.3	85.6	86.4	79.2	85.3	82.1
+ LongtoNotes _{eq} (n=10)	95.3	85.6	86.4	81.5	85.1	83.3	95.3	85.6	86.4	79.7	86.2	82.8	95.3	85.6	86.4	79.6	84.8	82.1
+ LongtoNotes _{eq} (n=10, m=40)	95.3	85.6	86.4	81.6	85.6	83.6	95.3	85.6	86.4	79.8	85.9	82.7	95.3	85.6	86.4	79.5	85.2	82.3

Table 16: Comparison of MUC scores

	OntoNotes						LongtoNotes _s						LongtoNotes					
	Mention			Coref			Mention			Coref			Mention			Coref		
	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁
SpanBERT Base (Lee et al., 2017)																		
+ OntoNotes	86.6	87.5	87.0	75.0	75.5	75.3	88.4	85.0	86.7	70.7	65.1	67.8	86.7	85.4	86.1	72.3	69.5	70.9
+ LongtoNotes _s	73.3	91.0	81.2	57.0	76.8	65.4	78.3	90.5	84	54.8	69.7	61.3	73.2	90.4	80.9	53.3	72.8	61.5
+ LongtoNotes	86.6	87.1	86.8	74.6	74.0	74.3	88.1	84.6	86.3	67.5	62.7	65.0	86.6	85.5	86.0	70.6	68.2	69.4
+ LongtoNotes _{eq}	86.1	87.8	87.0	74.9	75.2	75.0	87.7	86.2	87.0	69.7	67.0	68.3	86.1	86.3	86.2	71.7	70.6	71.2
LongDocCoref (Toshniwal et al., 2020b)																		
+ OntoNotes	95.3	85.6	86.4	72.2	77.9	74.9	95.3	85.6	86.4	57.9	71.7	64.0	95.3	85.6	86.4	63.9	74.7	68.9
+ LongtoNotes _s	95.3	85.6	86.4	18.3	61.7	28.2	95.3	85.6	86.4	10.7	53.6	17.9	95.3	85.6	86.4	16.1	58.7	25.2
+ LongtoNotes	95.3	85.6	86.4	73.3	76.7	75.0	95.3	85.6	86.4	61.0	70.1	65.2	95.3	85.6	86.4	65.5	73.7	69.4
+ LongtoNotes _{eq} (n=3)	95.3	85.6	86.4	73.7	76.9	75.2	95.3	85.6	86.4	64.4	70.4	67.3	95.3	85.6	86.4	67.5	73.7	70.5
+ LongtoNotes _{eq} (n=5)	95.3	85.6	86.4	73.4	77.3	75.3	95.3	85.6	86.4	64.5	70.9	67.6	95.3	85.6	86.4	67.5	74.2	70.7
+ LongtoNotes _{eq} (n=10)	95.3	85.6	86.4	73.6	77.0	75.3	95.3	85.6	86.4	64.5	70.9	67.6	95.3	85.6	86.4	68.3	73.5	70.8
+ LongtoNotes _{eq} (n=10, m=40)	95.3	85.6	86.4	73.5	78.1	75.7	95.3	85.6	86.4	65.0	70.5	67.6	95.3	85.6	86.4	67.9	74.4	71.0

Table 17: Comparison of BCUB scores

	OntoNotes						LongtoNotes _s						LongtoNotes					
	Mention			Coref			Mention			Coref			Mention			Coref		
	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁	P	R	F ₁
SpanBERT Base (Lee et al., 2017)																		
+ OntoNotes	86.6	87.5	87.0	71.5	73.7	72.1	88.4	85.0	86.7	63.3	61.6	62.4	86.7	85.4	86.1	68.1	68.4	68.2
+ LongtoNotes _s	73.3	91.0	81.2	53.2	69.5	60.3	78.3	90.5	84.0	51.5	59.2	55.1	73.2	90.4	80.9	50.4	64.2	56.5
+ LongtoNotes	86.6	87.1	86.8	70.8	73.1	71.9	88.1	84.6	86.3	63.4	60.5	61.9	86.6	85.5	86.0	67.7	68.2	67.9
+ LongtoNotes _{eq}	86.1	87.8	87.0	70.2	74.2	72.1	87.7	86.2	87.0	64.0	63.1	63.5	86.1	86.3	86.2	67.5	69.6	68.5
LongDocCoref (Toshniwal et al., 2020b)																		
+ OntoNotes	95.3	85.6	86.4	67.0	74.5	70.5	95.3	85.6	86.4	54.5	63.4	58.6	95.3	85.6	86.4	61.6	69.8	65.4
+ LongtoNotes _s	95.3	85.6	86.4	25.7	60.0	35.9	95.3	85.6	86.4	16.8	47.8	24.8	95.3	85.6	86.4	23.5	57.2	33.3
+ LongtoNotes	95.3	85.6	86.4	65.8	75.3	70.2	95.3	85.6	86.4	53.7	65.9	59.2	95.3	85.6	86.4	60.5	71.7	65.6
+ LongtoNotes _{eq} (n=3)	95.3	85.6	86.4	66.1	76.2	70.8	95.3	85.6	86.4	54.9	67.4	60.5	95.3	85.6	86.4	61.2	72.2	66.2
+ LongtoNotes _{eq} (n=5)	95.3	85.6	86.4	66.7	76.0	71.1	95.3	85.6	86.4	56.0	66.6	60.9	95.3	85.6	86.4	61.9	71.8	66.5
+ LongtoNotes _{eq} (n=10)	95.3	85.6	86.4	66.2	75.9	70.7	95.3	85.6	86.4	56.0	66.6	60.9	95.3	85.6	86.4	61.7	72.2	66.6
+ LongtoNotes _{eq} (n=10, m=40)	95.3	85.6	86.4	68.0	75.9	71.7	95.3	85.6	86.4	56.1	68.9	61.9	95.3	85.6	86.4	62.9	72.9	67.5

Table 18: Comparison of CEAFF scores