# Learning to Ignore: Long Document Coreference with Bounded Memory Neural Networks

**Shubham Toshniwal[1], Sam Wiseman[1], Allyson Ettinger[2], Karen Livescu[1], Kevin Gimpel[1]**
[1]Toyota Technological Institute at Chicago
[2]Department of Linguistics, University of Chicago

{shtoshni, swiseman, klivescu, kgimpel}@ttic.edu, aettinger@uchicago.edu

## Abstract

Long document coreference resolution remains a challenging task due to the large memory and runtime requirements of current models. Recent work doing incremental coreference resolution using just the global representation of entities shows practical benefits but requires keeping all entities in memory, which can be impractical for long documents. We argue that keeping all entities in memory is unnecessary, and we propose a memory-augmented neural network that tracks only a small bounded number of entities at a time, thus guaranteeing a linear runtime in length of document. We show that (a) the model remains competitive with models with high memory and computational requirements on OntoNotes and LitBank, and (b) the model learns an efficient memory management strategy easily outperforming a rule-based strategy.

## 1 Introduction

Long document coreference resolution poses runtime and memory challenges. Current best models for coreference resolution have large memory requirements and quadratic runtime in the document length (Joshi et al., 2019; Wu et al., 2020), making them impractical for long documents.

Recent work revisiting the entity-mention paradigm (Luo et al., 2004; Webster and Curran, 2014), which seeks to maintain explicit representations only of entities, rather than all their constituent mentions, has shown practical benefits for memory while being competitive with state-of-the-art models (Xia et al., 2020). In particular, unlike other approaches to coreference resolution which maintain representations of both mentions *and* their corresponding entity clusters (Rahman and Ng, 2011; Stoyanov and Eisner, 2012; Clark and Manning, 2015; Wiseman et al., 2016; Lee et al., 2018), the entity-mention paradigm stores representations only of the entity clusters, which are updated incrementally as coreference predictions are made. While such an approach requires less memory than those that additionally store mention representations, the number of entities can be impractically large when processing long documents, making the storing of all entity representations problematic.

Is it necessary to maintain an unbounded number of mentions or entities? Psycholinguistic evidence suggests it is not, as human language processing is incremental (Tanenhaus et al., 1995; Keller, 2010) and has limited working memory (Baddeley, 1986). In practice, we find that most entities have a small spread (number of tokens from first to last mention of an entity), and thus do not need to be kept persistently in memory. This observation suggests that tracking a limited, small number of entities at any time can resolve the computational issues, albeit at a potential accuracy tradeoff.

Previous work on bounded memory models for coreference resolution has shown potential, but has been tested only on short documents (Liu et al., 2019; Toshniwal et al., 2020). Moreover, this previous work makes token-level predictions while standard coreference datasets have span-level annotations. We propose a bounded memory model that performs quasi-online coreference resolution,[1] and test it on LitBank (Bamman et al., 2020) and OntoNotes (Pradhan et al., 2012). The model is trained to manage its limited memory by predicting whether to "forget" an entity already being tracked in exchange for a new (currently untracked) entity. Our empirical results show that: (a) the model is competitive with an unbounded memory version, and (b) the model's learned memory management outperforms a strong rule-based baseline.[2]

---

[1]"Quasi-online" because document encoding uses bidirectional transformers with access to future tokens.
[2]https://github.com/shtoshni92/long-doc-coref

Table 1: Max. Total Entity Count vs. Max. Active Entity Count.

|  | LitBank | OntoNotes |
|---|---|---|
| Max. Total Entity Count | 199 | 94 |
| Max. Active Entity Count | 18 | 24 |

## 2 Entity Spread and Active Entities

Given input document $\mathcal{D}$, let $(x_n)_{n=1}^N$ represent the $N$ mention spans corresponding to $M$ underlying entities $(e_m)_{m=1}^M$. Let $\text{START}(x_i)$ and $\text{END}(x_i)$ denote the start and end token indices of the mention span $x_i$ in document $\mathcal{D}$. Let $\text{ENT}(x_i)$ denote the entity of which $x_i$ is a mention. Given this notation we next define the following concepts.

**Entity Spread**   Entity spread denotes the interval of token indices from the first mention to the last mention of an entity. The entity spread $\text{ES}(e)$ of entity $e$ is given by:

$$\text{ES}(e) = [\min_{\text{ENT}(x)=e} \text{START}(x), \max_{\text{ENT}(x)=e} \text{END}(x)]$$

**Active Entity Count**   Active entity count $\text{AE}(t)$ at token index $t$ denotes the number of unique entities whose spread covers the token $t$, i.e., $\text{AE}(t) = |\{e \mid t \in \text{ES}(e)\}|$.

**Maximum Active Entity Count**   Maximum active entity count $\text{MAE}(\mathcal{D})$ for a document $\mathcal{D}$ denotes the maximum number of active entities at any token index in $\mathcal{D}$, i.e., $\text{MAE}(\mathcal{D}) = \max_{t \in [|\mathcal{D}|]} \text{AE}(t)$. This measure can be simply extended to a corpus $\mathcal{C}$ as: $\text{MAE}(\mathcal{C}) = \max_{\mathcal{D} \in \mathcal{C}} \text{MAE}(\mathcal{D})$.

Table 1 shows the MAE and the maximum total entity count in a single document, for LitBank and OntoNotes. For both datasets the maximum active entity count is much smaller than the maximum total entity count. Thus, rather than keeping all the entities in memory at all times, models can in principle simply focus on the far fewer active entities at any given time.

## 3 Model

Based on the preceding finding, we will next describe models that require tracking only a small, bounded number of entities at any time.

To make coreference predictions for a document, we first encode the document and propose candidate mentions. The proposed mentions are then processed sequentially and are either: (a) added to an existing entity cluster, (b) added to a new cluster, (c) ignored due to limited memory capacity (for bounded memory models), or (d) ignored as an invalid mention.

**Document Encoding** is done using the SpanBERT$_{\text{LARGE}}$ model finetuned for OntoNotes and released as part of the coreference model of Joshi et al. (2020). We don't further finetune the SpanBERT model. To encode long documents, we segment the document using the *independent* and *overlap* strategies described in Joshi et al. (2019).[3] In *overlap* segmentation, for a token present in overlapping BERT windows, the token's representation is taken from the BERT window with the most neighboring tokens of the concerned token. For both datasets we find that *overlap* slightly outperforms *independent*.

**Mention Proposal**   Given the encoded document, we next predict the top-scoring mentions which are to be clustered. The goal of this step is to have high recall, and we follow previous work to threshold the number of spans chosen (Lee et al., 2017). Given a document $\mathcal{D}$, we choose $0.3 \times |\mathcal{D}|$ top spans for LitBank, and $0.4 \times |\mathcal{D}|$ for OntoNotes.

We pretrain the mention proposal model before training the mention proposal and mention clustering pipeline end-to-end, as done by Wu et al. (2020). The reason is that without pretraining, most of the mentions proposed by the mention proposal model would be invalid mentions, i.e., spans that are not mentions, which would not provide any training signal to the mention clustering stage.

**Mention Clustering**   Let $(x_i)_{i=1}^K$ represent the top-$K$ candidate mention spans from the mention proposal step and let $s_m(x_i)$ represent the mention score for span $x_i$, which indicates how likely it is that a span constitutes a mention. Assume that the mentions are already ordered based on their position in the document and are processed sequentially in that order.[4] Let $E = (e_m)_{m=1}^M$ represent the $M$ entities currently being tracked by the model (initially $M = 0$). For ease of discussion, we will overload the terms $x_i$ and $e_j$ to also correspond to their respective representations.

In the *first* step, the model decides whether the

---

[3]We modify the *overlap* segmentation to respect sentence boundary or token boundary when possible.

[4]Specifically, they are ordered based on $\text{START}(\cdot)$ index with ties broken using $\text{END}(\cdot)$.

span $x_i$ refers to any of the entities in $E$ as follows:

$$s_c(x_i, e_j) = f_c([x_i; e_j; x_i \odot e_j; g(x_i, e_j)]) + s_m(x_i)$$
$$s_c^{top} = \max_{j=1...M} s_c(x_i, e_j)$$
$$e^{top} = \arg\max_{j=1...M} s_c(x_i, e_j)$$

where $\odot$ represents the element-wise product, and $f_c(\cdot)$ corresponds to a learned feedforward neural network. The term $g(x_i, e_j)$ correponds to a concatenation of feature embeddings that includes embeddings for (a) number of mentions in $e_j$, (b) number of mentions between $x_i$ and last mention of $e_j$, (c) last mention decision, and (d) document genre (only for OntoNotes).

Now if $s_c^{top} > 0$ then $x_i$ is considered to refer to $e^{top}$, and $e^{top}$ is updated accordingly.[5] Otherwise, $x_i$ does not refer to any entity in $E$ and a *second* step is executed, which will depend on the choice of memory architecture. We test three memory architectures, described below.

1. **Unbounded Memory (U-MEM)**: If $s_m(x_i) > 0$ then we create a new entity $e_{M+1} = x_i$ and append it to $E$. Otherwise the mention is ignored as invalid, i.e., it doesn't correspond to an entity. Ignoring invalid mentions is important for datasets such as LitBank where singletons are explicitly annotated and used for evaluation. For OntoNotes, where singletons are not annotated and ignored for evaluation, we also consider a variant U-MEM* which appends all non-coreferent mentions, as done in Xia et al. (2020).

2. **Bounded Memory**: Suppose the model has a capacity of tracking $C$ entities at a time. If $C > M$, i.e., the memory capacity has not been fully utilized, then the model behaves like U-MEM. Otherwise, the bounded memory models must decide between: (a) evicting an entity already being tracked, (b) ignoring $x_i$ due to limited capacity, and (c) ignoring the mention as invalid. We test two bounded memory variants that are described below.

(a) **Learned Bounded Memory (LB-MEM)**: The proposed LB-MEM architecture tries to predict a score $f_r(.)$ corresponding to the anticipated number of remaining mentions for any entity or mention, and compares it against the mention score $s_m(x_i)$ as follows:

$$d = \arg\min[f_r(e_1), \ldots, f_r(e_M), f_r(x_i), s_m(x_i)]$$

Table 2: Results for LitBank (CoNLL F1).

| Model | Dev F1 | Test F1 |
|---|---|---|
| U-MEM | 77.1 | 76.5 |
| LB-MEM | | |
| 5 cells | 71.9 | 70.3 |
| 10 cells | 75.0 | 74.7 |
| 20 cells | 75.7 | 75.1 |
| RB-MEM | | |
| 5 cells | 58.5 | 57.8 |
| 10 cells | 69.9 | 69.0 |
| 20 cells | 75.3 | 74.4 |
| Bamman et al. (2020) | - | 68.1 |

where $f_r(\cdot)$ is a learned feedforward neural network. If $1 \leq d \leq M$ then then the model evicts the previous entity $e_d$ and reinitialize it to $x_i$. Otherwise if $d = M + 1$ then the model ignores $x_i$ due to limited capacity. Finally if $d = M + 2$ then the model predicts the mention to be invalid.

(b) **Rule-based Bounded Memory (RB-MEM)** The Least Recently Used (LRU) principle is a popular choice among memory models (Rae et al., 2016; Santoro et al., 2016). While LB-MEM considers all potential entities for eviction, with RB-MEM this choice is restricted to just the LRU entity, i.e., the entity whose mention was least recently seen. The rest of the steps are similar to the LB-MEM model.

**Training** All the models are trained using teacher forcing. The ground truth decisions for bounded memory models are chosen to maximize the number of mentions tracked by the model (details in Appendix A.3). Finally, the training loss is calculated via the addition of the cross-entropy losses for the two steps of mention clustering.

## 4 Experimental Setup

### 4.1 Datasets

**LitBank** is a recent coreference dataset for literary texts (Bamman et al., 2020). The dataset consists of prefixes of 100 novels with an average length of 2100 words. Singletons are marked and used for evaluation. Evaluation is done via 10-fold cross-validation over 80/10/10 splits.[6]

**OntoNotes** consists of 2802/343/348 documents in the train/development/test splits, respectively (Pradhan et al., 2012). The documents span 7 genres and have an average length of 463 words. Singletons are not marked in the dataset.

---

[5]We use weighted averaging where the weight for $e^{top}$ corresponds to the number of previous mentions seen for $e^{top}$.

[6]https://github.com/dbamman/lrec2020-coref/tree/master/data

Table 3: Results for OntoNotes (CoNLL F1).

| Model | Dev F1 | Test F1 |
|---|---|---|
| U-MEM | 78.4 | 78.1 |
| U-MEM* | 79.6 | 79.6 |
| LB-MEM | | |
| 5 cells | 74.0 | 73.3 |
| 10 cells | 77.1 | 76.8 |
| 20 cells | 78.1 | 78.2 |
| RB-MEM | | |
| 5 cells | 69.8 | 69.6 |
| 10 cells | 75.9 | 75.5 |
| 20 cells | 78.2 | 77.8 |
| U-MEM* (Xia et al., 2020) | 79.7 | 79.4 |
| Joshi et al. (2020) | 80.1 | 79.6 |
| Wu et al. (2020) | 83.4 | 83.1 |

Table 4: Peak memory and inference time statistics for the LitBank cross-validation split 0. Note that the training memory statistics depend on document truncation and sampling probability of invalid mentions. The models in this table are trained without document truncation and sample 20% of invalid mentions during training.

| Model | Peak training mem. (in GB) | Peak inference mem. (in GB) | Inference time (in s) |
|---|---|---|---|
| U-MEM | 11.6 | 3.1 | 29.25 |
| LB-MEM | | | |
| 5 cells | 8.0 | 3.2 | 27.31 |
| 10 cells | 8.4 | 3.2 | 27.44 |
| 20 cells | 9.1 | 3.2 | 27.86 |
| RB-MEM | | | |
| 5 cells | 8.0 | 3.2 | 26.19 |
| 10 cells | 8.3 | 3.2 | 26.50 |
| 20 cells | 8.9 | 3.2 | 26.19 |

## 4.2 Hyperparameters

Document encoding is done using the SpanBERT$_{LARGE}$ model of Joshi et al. (2020) which was finetuned for OntoNotes.[7] The SpanBERT model is not further finetuned. The other model parameters are trained using the Adam optimizer (Kingma and Ba, 2014) with an initial learning rate of $2 \times 10^{-4}$ which is linearly decayed. For span representation, we use the embedding function described in Lee et al. (2017). For OntoNotes we follow the setup of Xia et al. (2020). We differ, however, in training all the model parameters, except SpanBERT, from scratch. The models are trained for a maximum of 15 epochs for OntoNotes, and 25 epochs for LitBank. For both the datasets, the training stops if dev performance doesn't improve for 5 epochs. For more details see Appendix A.2.

## 5 Results

Tables 2 and 3 show results of all the proposed models for LitBank and OntoNotes respectively. Detailed results with the performance on different coreference metrics is presented in Table 11 for Lit-Bank, and Table 12 for OntoNotes (Appendix A.4).

As expected, the bounded memory models improve with increase in memory. For both the datasets, the LB-MEM model with 20 memory cells is competitive with the U-MEM model though the gap between them for LitBank is non-trivial. Among the bounded memory models, the LB-MEM model is significantly better than RB-MEM for lower numbers of memory cells. We analyze the

---

[7]From the original models, we stripped out just the Span-BERT part which is available at https://huggingface.co/shtoshni/spanbert_coreference_large

reasons for this in the next section.

For OntoNotes, the U-MEM* model easily outperforms the U-MEM model which is trained to ignore all non-gold mentions. These non-gold mentions also include singletons in case of OntoNotes. Thus, the U-MEM model essentially has to predict if a non-coreferent mention will be coreferent with future mentions or not. U-MEM* avoids this difficult problem by adding all non-coreferent mentions to memory. Since in OntoNotes singletons are removed during evaluation, the U-MEM* model is not penalized for predicting singletons corresponding to invalid mentions, and otherwise. Note that the U-MEM* model doesn't make sense for LitBank where singletons are used for evaluation. The initial empirical results also confirmed that decisively.

Between the two datasets, we see that the increase in memory results in larger improvement for LitBank. We also establish a new state-of-the-art for LitBank with the U-MEM model. For OntoNotes, the U-MEM* model matches the performance of a similar model by Xia et al. (2020). Remarkably, the two U-MEM* models almost match the performance of the computationally and memory intensive span-ranking model of Joshi et al. (2020) whose finetuned SpanBERT document encoder is used by these two models. We expect gains by further finetuning the SpanBERT model and learning a parameterized global entity representation, but we leave them for future work.

## 6 Analysis

In this section we analyze the behavior of the three memory models on LitBank and OntoNotes.

Table 5: Comparison of number of entities in memory.

| Model | LitBank | | OntoNotes | |
| --- | --- | --- | --- | --- |
| | Avg | Max | Avg | Max |
| U-MEM | 80.8 | 160 | 16.0 | 83 |
| U-MEM* | - | - | 173 | 962 |
| LB-MEM | | | | |
| 5 cells | 5.0 | 5 | 4.6 | 5 |
| 10 cells | 10.0 | 10 | 7.8 | 10 |
| 20 cells | 20.0 | 20 | 12.1 | 20 |
| RB-MEM | | | | |
| 5 cells | 5.0 | 5 | 4.6 | 5 |
| 10 cells | 10.0 | 10 | 7.9 | 10 |
| 20 cells | 20.0 | 20 | 11.9 | 20 |

Table 6: Average number of mentions ignored by the two bounded memory models.

| Memory size | LitBank | | OntoNotes | |
| --- | --- | --- | --- | --- |
| | LB-MEM | RB-MEM | LB-MEM | RB-MEM |
| 5 | 4.5 | 70.0 | 0.3 | 3.7 |
| 10 | 0.0 | 14.2 | 0.0 | 0.4 |
| 20 | 0.0 | 0.4 | 0.0 | 0.1 |

Table 7: Error Analysis for OntoNotes dev set. CE=Conflated Entities, DE=Divided Entity, EM=Extra Mention, EE=Extra Entity, MM=Missing Mention, ME=Missing Entity.

| Model | CE | DE | EM | EE | MM | ME |
| --- | --- | --- | --- | --- | --- | --- |
| U-MEM | 853 | 496 | 515 | 904 | 545 | 603 |
| U-MEM* | 754 | 466 | 504 | 816 | 527 | 583 |
| LB-MEM | | | | | | |
| 5 cells | 706 | 381 | 340 | 972 | 844 | 1116 |
| 10 cells | 757 | 425 | 340 | 868 | 655 | 894 |
| 20 cells | 752 | 396 | 402 | 859 | 613 | 799 |
| RB-MEM | | | | | | |
| 5 cells | 672 | 369 | 365 | 1146 | 923 | 1359 |
| 10 cells | 722 | 393 | 403 | 986 | 696 | 931 |
| 20 cells | 713 | 420 | 380 | 833 | 559 | 853 |

**Memory Utilization** Table 4 compares the memory and inference time statistics for the different memory models for the LitBank cross-validation split zero.[8] For training, the bounded memory models are significantly less memory intensive than the U-MEM model. The table also shows that the bounded memory models are faster than the U-MEM memory model during inference (inference time calculated by averaging over three runs). This is because the number of entities tracked by the U-MEM memory model grows well beyond the maximum of 20 memory slots reserved for the bounded models as shown in Table 5.

Surprisingly, for inference we see that the bounded models have a slightly larger memory footprint than the U-MEM model. This is because the document encoder, SpanBERT, dominates the memory usage during inference (as also observed by Xia et al., 2020). Thus the peak memory usage during inference is determined by the mention proposal stage rather than the mention clustering stage. And during the mention proposal stage, the additional parameters of bounded memory models, which are loaded as part of the whole model, cause the slight uptick in peak inference memory. Note that using a cheaper encoder or running on a sufficiently long document, such as a book, can change these results.

**Number of Entities in Memory** Table 5 compares the maximum number of entities kept in memory by the different memory models for the LitBank cross-validation dev sets and the OntoNotes dev set. As expected, the U-MEM model keeps more entities in memory than the bounded memory models on average for both datasets. For LitBank the difference is especially stark with the U-MEM model tracking about 4/8 times more entities in memory

---

[8]Peak memory usage estimated via `torch.cuda.max_memory_allocated()`

on average/worst case, respectively. The difference between the U-MEM and U-MEM* model is striking, with U-MEM* tracking more than 10 times the entities of U-MEM in both the average and worst case. Also, while some OntoNotes documents do not use even the full 5 memory cell capacity, all LitBank documents fully utilize even the 20 memory cell capacity. This is because LitBank documents are more than four times as long as OntoNotes documents, and LitBank has singletons marked. These results also justify our initial motivation that with long documents, the memory requirement will increase even if we only keep the entity representations.

**LB-MEM vs. RB-MEM** Table 6 compares the number of mentions ignored by LB-MEM and RB-MEM. The LB-MEM model ignores far fewer mentions than RB-MEM. This is because while the RB-MEM model can only evict the LRU entity, which might not be optimal, the LB-MEM model can choose any entity for eviction. These statistics combined with the fact that the LB-MEM model typically outperforms RB-MEM mean that the LB-MEM model is able to anticipate which entities are important and which are not.

**Error Analysis** Table 7 presents the results of automated error analysis done using the Berkeley Coreference Analyzer (Kummerfeld and Klein,

2013) for the OntoNotes dev set. As the memory capacity of models increases, the errors shift from missing mention, missing entity, and divided entity categories, to conflated entities, extra mention, and extra entity categories. The LB-MEM model outperforms RB-MEM in terms of tracking more entities.

# 7 Conclusion and Future Work

We propose a memory model which tracks a small, bounded number of entities. The proposed model guarantees a linear runtime in document length, and in practice significantly reduces peak memory usage during training. Empirical results on LitBank and OntoNotes show that the model is competitive with an unbounded memory version and outperforms a strong rule-based baseline. In particular, we report state of the art results on LitBank. In future work we plan to apply our model to longer, book length documents, and plan to add more structure to the memory.

## Acknowledgments

## References

Alan Baddeley. 1986. *Working Memory*. Oxford University Press.

David Bamman, Olivia Lewke, and Anya Mansoor. 2020. An Annotated Dataset of Coreference in English Literature. In *LREC*.

Kevin Clark and Christopher D. Manning. 2015. Entity-Centric Coreference Resolution with Model Stacking. In *ACL*.

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. SpanBERT: Improving Pre-training by Representing and Predicting Spans. *TACL*, 8.

Mandar Joshi, Omer Levy, Luke Zettlemoyer, and Daniel Weld. 2019. BERT for Coreference Resolution: Baselines and Analysis. In *EMNLP*.

Frank Keller. 2010. Cognitively Plausible Models of Human Language Processing. In *ACL*.

Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In *ICLR*.

Jonathan K. Kummerfeld and Dan Klein. 2013. Error-Driven Analysis of Challenges in Coreference Resolution. In *EMNLP*.

Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. End-to-end Neural Coreference Resolution. In *EMNLP*.

Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018. Higher-Order Coreference Resolution with Coarse-to-Fine Inference. In *NAACL-HLT*.

Fei Liu, Luke Zettlemoyer, and Jacob Eisenstein. 2019. The Referential Reader: A Recurrent Entity Network for Anaphora Resolution. In *ACL*.

Xiaoqiang Luo, Abe Ittycheriah, Hongyan Jing, Nanda Kambhatla, and Salim Roukos. 2004. A Mention-Synchronous Coreference Resolution Algorithm Based On the Bell Tree. In *ACL*.

Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. CoNLL-2012 Shared Task: Modeling Multilingual Unrestricted Coreference in OntoNotes. In *CoNLL*.

Jack W. Rae, Jonathan J. Hunt, Ivo Danihelka, Timothy Harley, Andrew W. Senior, Gregory Wayne, Alex Graves, and Tim Lillicrap. 2016. Scaling Memory-Augmented Neural Networks with Sparse Reads and Writes. In *NeurIPS*.

Altaf Rahman and Vincent Ng. 2011. Narrowing the modeling gap: a cluster-ranking approach to coreference resolution. *JAIR*, 40:469–521.

Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy P. Lillicrap. 2016. One-shot Learning with Memory-Augmented Neural Networks. In *ICML*.

Veselin Stoyanov and Jason Eisner. 2012. Easy-first Coreference Resolution. In *COLING*.

MK Tanenhaus, MJ Spivey-Knowlton, KM Eberhard, and JC Sedivy. 1995. Integration of visual and linguistic information in spoken language comprehension. *Science*, 268(5217).

Shubham Toshniwal, Allyson Ettinger, Kevin Gimpel, and Karen Livescu. 2020. PeTra: A Sparsely Supervised Memory Model for People Tracking. In *ACL*.

Kellie Webster and James R. Curran. 2014. Limited memory incremental coreference resolution. In *COLING*.

Sam Wiseman, Alexander M. Rush, and Stuart M. Shieber. 2016. Learning Global Features for Coreference Resolution. In *NAACL*.

Wei Wu, Fei Wang, Arianna Yuan, Fei Wu, and Jiwei Li. 2020. Coreference Resolution as Query-based Span Prediction. In *ACL*.

Patrick Xia, João Sedoc, and Benjamin Van Durme. 2020. Revisiting Memory-Efficient Incremental Coreference Resolution. In *EMNLP*.
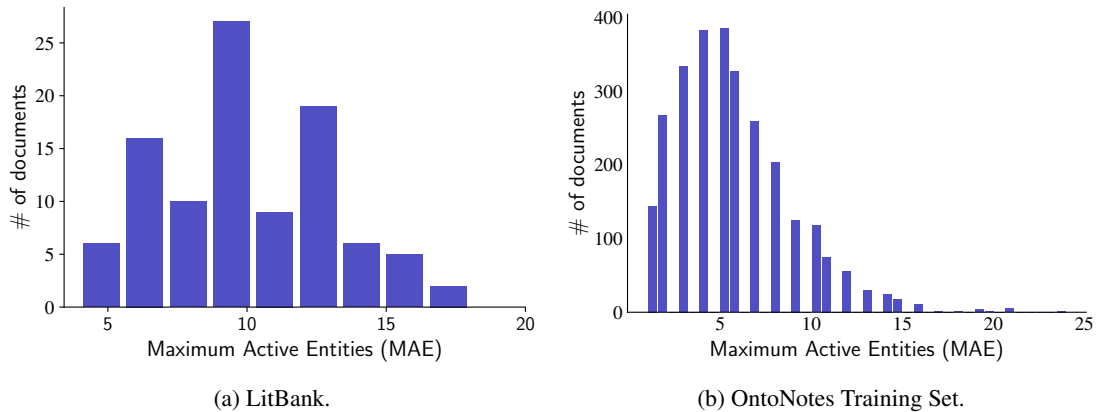
(a) LitBank.



(b) OntoNotes Training Set.

Figure 1: Histograms of Maximum Active Entities for documents in LitBank and OntoNotes.



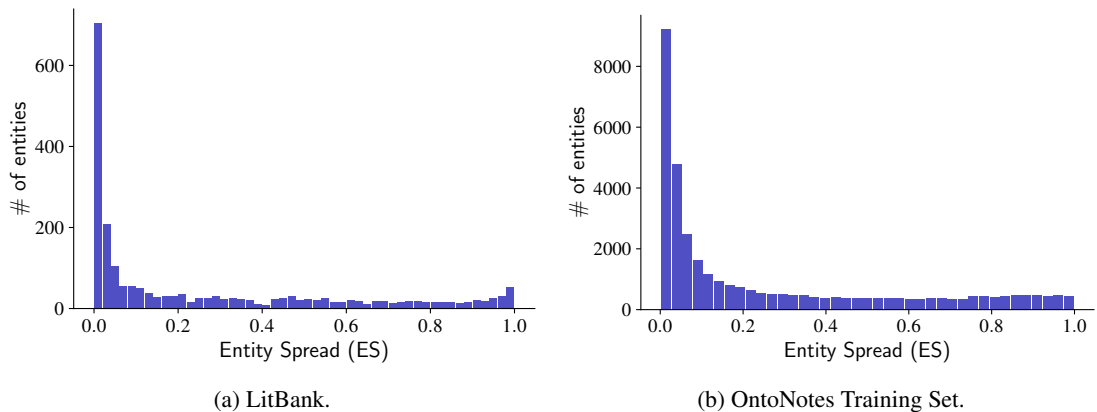(a) LitBank.



(b) OntoNotes Training Set.

Figure 2: Histograms of Entity Spread as fraction of document length for LitBank and OntoNotes.

## A  Appendix

### A.1  Maximum Active Entities

Figure 1 visualizes the histograms of length of Entity Spread (ES), defined in Section 2, as a fraction of document length for documents in LitBank and OntoNotes. For LitBank we only visualize the entity spread of non-singleton clusters because otherwise the histogram is too skewed towards one. Figure 2 visualizes the histograms of Maximum Active Entity Count (MAE), defined in Section 2, for documents in LitBank and OntoNotes.

Table 8: Hyperparameter options for OntoNotes with preferred choices highlighted in bold.

| Parameter | Range |
|---|---|
| Dropout | {0.4} |
| FFNN hidden layer | {3000} |
| FFNN # of hidden layers | 1 |
| Document Encoding | {Independent, **Overlap**} |
| Label Smoothing | {0.0, 0.01, 0.1} |
| Sampling Prob. Invalid Mentions | {0.25, 0.5, 0.75, 1.0} |
| Max. # of BERT Segments | {**3**, 5} |
| Non-coreferent entity weight | {1.0} |

### A.2  Model Details

**Hyperparameter Choices**  We stick with the hyperparameters for feedforward neural network (FFNN) size and depth from Joshi et al. (2020). We didn't do much exploration with dropout but with the limited experiments our finding was that there was little separating dropout probability of 0.3 and 0.4. Among choices for how to segment document into BERT windows, we found overlapping windows to work better than independent BERT windows. Two very important hyperparameters that affect peak memory usage during training are: (a) maximum number of BERT segments, and (b) sampling probability of invalid mentions. Truncating the document by selecting a chosen maximum number of contiguous BERT segments essentially caps the length of documents during training. And the second hyperparameter of sampling invalid mentions controls the number of invalid mentions, which happens to be the overwhelming category of proposed mentions, the model sees during training. We also explore two hyperparameters for the

Table 9: Hyperparameter options for LitBank with preferred choices highlighted in bold.

| Parameter | Range |
|---|---|
| Dropout | {0.3} |
| FFNN hidden layer | {3000} |
| FFNN # of hidden layers | 1 |
| Document Encoding | {Independent, **Overlap**} |
| Label Smoothing | {0.0} |
| Sampling Prob. Invalid Mentions | {0.25, 0.5, 0.75, 1.0} |
| Max. # of BERT Segments | {3, **5**} |
| Non-coreferent entity weight | {1.0, **2.0**} |

cross-entropy loss of the first step of mention clustering: (a) label smoothing for regularization, and (b) weight of the non-coreferent term in the cross-entropy loss.

**Specific Hyperparameter Choices for OntoNotes** We didn't see any gain by increasing the maximum number of BERT segments from 3 to 5 in our initial experiments. The U-MEM and bounded models preferred lower sampling probabilities for invalid mentions but no clear winner in label smoothing weight. The U-MEM* model preferred low label smoothing weight and higher sampling probabilities for invalid mentions.

**Specific Hyperparameter Choices for LitBank** Initial experiments with cross validation splits {0, 1, 2} showed that models preferred maximum number of BERT segments to be 5 in comparison to 3. This might be because most of the LitBank documents are really long, and training on a maximum of 3 BERT segments might lead to a bigger mismatch between training and inference. Another hyperparameter that proved important for LitBank was the non-coreferent entity weight of 2.0. Due to explosion of combinations driven by the fact that there are 10 cross validation splits, we didn't explore label smoothing for LitBank.

### A.3 Ground Truth Generation

In this section we explain how the ground truth action sequence is generated corresponding to the predicted mention sequence. The ground truth for U-MEM model is fairly straight forward. For the bounded memory models, we keep growing the number of entities till we hit the memory ceiling. For all the entities in memory, we maintain the number of mentions remaining in the ground truth cluster. For example, a cluster with a total of five mentions, two of which have already been processed by the model, has three remaining mentions.

Table 10: Number of model parameters (in millions).

| | LitBank | OntoNotes |
|---|---|---|
| U-MEM | 37.36 | 37.42 |
| LB-MEM | 46.83 | 46.95 |
| RB-MEM | 46.83 | 46.95 |

Suppose now a mention corresponding to a currently untracked entity comes in and the memory is already at full capacity. Then for the LB-MEM model, we compare the number of mentions of this new entity (along with the current mention) against the number of mentions remaining for all the entities currently being tracked. If there are entities in memory with number of remaining mentions less than or equal to the number of mentions of this currently untracked entity, then the untracked entity replaces the entity with the least number of remaining mentions. Ties among the entities with least number of remaining mentions are broken by the least recently seen entity. If there's no such entity in the memory, then the mention is ignored. For the RB-MEM model, the comparison is done in a similar way but is limited to just the LRU entity.

### A.4 Miscellany

**Computing Infrastructure & Runtime** All the models for a single cross validation split of LitBank can be trained within 4 hours. Training on OntoNotes finishes within 12-20 hours. The U-MEM* model where all invalid mentions are seen during training is the only configuration that requires 24GB memory GPUs, all other configurations can be trained on 12GB memory GPUs.

**Number of model parameters.** Table 10 shows the number of trainable parameters for all the model and dataset combinations. LB-MEM and RB-MEM have additional parameters in comparison to U-MEM for predicting a score corresponding to the number of remaining mentions for an entity. Comparing across datasets, the OntoNotes models have a few additional parameters than their LitBank counterparts for modeling the document genre.

**Evaluation Metric Code.** We use the coreference scorer Perl script available at https://github.com/conll/reference-coreference-scorers. We also use the Python implementation by Kenton Lee available at https://github.com/kentonl/e2e-coref/blob/master/metrics.py. The two scripts can have some rounding differences.

Table 11: Detailed results of the proposed models on the aggregated LitBank cross-validation test set.

| Model | MUC | | | B³ | | | CEAF$_{\phi_4}$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Prec. | Rec. | F1 | Prec. | Rec. | F1 | Prec. | Rec. | F1 | Avg. F1 |
| U-MEM | 90.8 | 85.7 | 88.2 | 80.0 | 72.1 | 75.9 | 65.1 | 66.0 | 65.5 | 76.5 |
| LB-MEM | | | | | | | | | | |
|   5 cells | 90.9 | 80.0 | 85.1 | 77.4 | 64.0 | 70.1 | 57.8 | 53.8 | 55.7 | 70.3 |
|   10 cells | 90.0 | 84.6 | 87.2 | 78.1 | 70.8 | 74.2 | 64.2 | 61.1 | 62.6 | 74.7 |
|   20 cells | 90.3 | 85.0 | 87.6 | 79.2 | 70.9 | 74.8 | 64.1 | 62.0 | 63.0 | 75.1 |
| RB-MEM | | | | | | | | | | |
|   5 cells | 91.4 | 74.6 | 82.2 | 75.7 | 51.1 | 61.0 | 52.2 | 21.3 | 30.3 | 57.8 |
|   10 cells | 91.1 | 81.3 | 85.9 | 78.5 | 62.1 | 69.3 | 56.3 | 47.8 | 51.7 | 69.0 |
|   20 cells | 90.5 | 85.1 | 87.7 | 79.7 | 69.8 | 74.4 | 61.1 | 61.0 | 61.1 | 74.4 |

Table 12: Detailed results of the proposed models on the OntoNotes 5.0 test set.

| Model | MUC | | | B³ | | | CEAF$_{\phi_4}$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Prec. | Rec. | F1 | Prec. | Rec. | F1 | Prec. | Rec. | F1 | Avg. F1 |
| U-MEM | 84.6 | 84.1 | 84.3 | 77.2 | 76.2 | 76.7 | 72.5 | 74.3 | 73.4 | 78.1 |
| U-MEM* | 85.5 | 85.1 | 85.3 | 78.7 | 77.3 | 78.0 | 74.2 | 76.5 | 75.3 | 79.6 |
| LB-MEM | | | | | | | | | | |
|   5 cells | 76.4 | 86.2 | 81.0 | 66.4 | 78.4 | 71.9 | 62.0 | 72.7 | 66.9 | 73.3 |
|   10 cells | 81.7 | 85.9 | 83.8 | 72.8 | 77.9 | 75.3 | 67.0 | 76.4 | 71.4 | 76.8 |
|   20 cells | 83.2 | 86.2 | 84.7 | 74.8 | 78.9 | 76.8 | 70.0 | 76.7 | 73.2 | 78.2 |
|   30 cells | 83.8 | 85.6 | 84.7 | 76.1 | 78.2 | 77.1 | 70.4 | 77.1 | 73.6 | 78.5 |
| RB-MEM | | | | | | | | | | |
|   5 cells | 72.0 | 85.7 | 78.3 | 60.1 | 78.9 | 68.2 | 57.0 | 68.9 | 62.4 | 69.6 |
|   10 cells | 80.1 | 85.7 | 82.8 | 70.5 | 78.3 | 74.2 | 66.0 | 73.4 | 69.5 | 75.5 |
|   20 cells | 82.8 | 85.9 | 84.3 | 74.8 | 78.3 | 76.5 | 68.0 | 77.4 | 72.4 | 77.8 |
|   30 cells | 84.0 | 85.2 | 84.6 | 76.2 | 78.2 | 77.2 | 72.1 | 75.6 | 73.8 | 78.5 |

Table 13: Spearman correlation of F1 score with document length and # of entities in OntoNotes dev set.

| Model | Document Length | # of Entities |
|---|---|---|
| U-MEM | -0.31 | -0.28 |
| U-MEM* | -0.28 | -0.25 |
| LB-MEM | | |
|   5 cells | -0.36 | -0.37 |
|   10 cells | -0.34 | -0.33 |
|   20 cells | -0.34 | -0.31 |
| RB-MEM | | |
|   5 cells | -0.37 | -0.41 |
|   10 cells | -0.29 | -0.30 |
|   20 cells | -0.31 | -0.29 |

Table 14: Spearman correlation of F1 score with document length and # of entities in aggregated LitBank cross-validation dev set.

| Model | Document Length | # of Entities |
|---|---|---|
| U-MEM | -0.07 | -0.36 |
| LB-MEM | | |
|   5 cells | 0.01 | -0.23 |
|   10 cells | -0.06 | -0.35 |
|   20 cells | -0.03 | -0.32 |
| RB-MEM | | |
|   5 cells | -0.00 | -0.41 |
|   10 cells | 0.01 | -0.36 |
|   20 cells | -0.02 | -0.33 |

**Effect of Document Length and Number of Entities.** Table 13 and 14 presents the Spearman correlation of document F1 score with document length and number of entities in the document. The correlations are mostly negative because the task becomes more challenging with an increase in document length and entities, though for LitBank the length of the document doesn't seem to be a great indicator of the hardness of the task. The increase in memory capacity for bounded models results in less negative correlations, suggesting improved performance for challenging documents.