

# Major Entity Identification: A Generalizable Alternative to Coreference Resolution

Anonymous ACL submission

## Abstract

The limited generalization of coreference resolution (CR) models has been a major bottleneck in the task’s broad application. Prior work has identified annotation differences, especially for mention detection, as one of the main reasons for the generalization gap and proposed using additional annotated target domain data. Rather than relying on this additional annotation, we propose an alternative formulation of the CR task, **Major Entity Identification (MEI)**, where we: (a) assume the target entities to be specified in the input, and (b) limit the task to only the frequent entities. Through extensive experiments, we demonstrate that MEI models generalize well across domains on multiple datasets with supervised models and LLM-based few-shot prompting. Additionally, the MEI task fits the classification framework, which enables the use of classification-based metrics that are more robust than the current CR metrics. Finally, MEI is also of practical use as it allows a user to search for all mentions of a particular entity or a group of entities of interest.

## 1 Introduction

Coreference resolution (CR) is the task of finding text spans that refer to the same entity. CR is a fundamental language understanding task relevant to various downstream NLP applications, such as question-answering (Dhingra et al., 2018), building knowledge graphs (Koncel-Kedziorski et al., 2019), and summarization (Sharma et al., 2019). Despite the importance of CR and the progress made by neural coreference models (Dobrovolskii, 2021; Bohnet et al., 2023; Zhang et al., 2023), domain generalization remains an issue even with the best-performing supervised models (Xia and Van Durme, 2021; Toshniwal et al., 2021).

The lack of domain generalization in CR models can largely be attributed to differences in annotation guidelines of popular CR benchmarks, specifically annotation guidelines about what con-

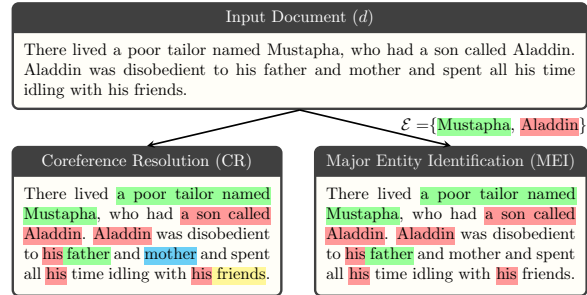


Figure 1: CR vs. MEI. The CR task aims to detect and cluster all mentions into different entities, shown in various colors. MEI takes major entities as additional input and aims to detect and classify the mentions that refer only to these entities.

stitutes a mention (Porada et al., 2023). For example, OntoNotes (Pradhan et al., 2013) does not annotate singletons, confounding mention identity with being referential. Thus, models trained on OntoNotes generalize poorly. The importance of mention detection for CR generalization is further highlighted by Gandhi et al. (2023), who show that solely annotating mentions is sufficient and more efficient for adapting pre-trained coreference models to new domains (in comparison to annotating coreference chains). Similarly, GPT-4 struggles with zero-/few-shot mention prediction, but given ground-truth mentions, its CR performance is competitive with the best-supervised models (Le and Ritter, 2023).

Given these observations, we hypothesize that current CR models, including large language models, generalize well at *mention clustering* but struggle to generalize on *mention detection* due to idiosyncrasies of different domains/benchmarks. We put forth an alternative formulation of the CR task where the entities of interest are provided as additional input. Assuming entities to be part of the input offloads the required domain adaptation from training to inference. Specifically, we propose the task of Major Entity Identification (MEI), where we assume the major entities of the narrative, de-

Statistics	LitBank		FantasyCoref	
	CR	MEI	CR	MEI
# of Mentions	29103	16985	56968	35938
# of Non singletons	23340	16985	56968	35938
Mean ant. dist.	55.31	36.95	57.58	30.24
# of Clusters	7927	490	5829	942
Avg. cluster size	3.67	34.66	9.77	38.15

Table 1: Comparing CR and MEI. MEI has fewer but larger clusters, and a smaller mean antecedent distance (Mean ant. dist.). Our formulation’s frequency-based criterion for deciding major entities means that singleton mentions are typically not a part of MEI.

069 fined as the most frequently occurring entities, to  
070 be provided as input along with the text (see Fig. 1).  
071 We focus on major entities for the following reasons:  
072 (a) Specifying major entities of a narrative is  
073 intuitively easier. (b) A handful of major entities  
074 often dominate any discourse. Table 1 shows that in  
075 LitBank roughly 6% of entities (490 of 7927) contribute  
076 to 60% of the mentions (16985 of 29103).

077 To test the generalizability of MEI, we adapt two  
078 literary CR benchmarks, namely LitBank (Bamman  
079 et al., 2020) and FantasyCoref (Han et al., 2021),  
080 and a state-of-the-art coreference model (Toshniwal  
081 et al., 2021) to MEI. While there is a big  
082 gap in CR performance between in- and out-of-  
083 domain models (Toshniwal et al., 2021), we show  
084 that this performance gap is much smaller for MEI  
085 (Section 5.1). To test this hypothesis further, we  
086 evaluate large language models (LLMs) for MEI in  
087 a few-shot learning setup. On CR, LLMs are shown  
088 to struggle with mention detection and perform  
089 worse than supervised models (Le and Ritter, 2023).  
090 Contrary to this, on MEI, top LLMs (e.g. GPT-4)  
091 are only slightly behind supervised models (Section  
092 5.2). These experiments in the supervised  
093 setting and the few-shot setting demonstrate that  
094 the MEI task is more generalizable than CR.

095 Additionally, we argue that MEI is easier to evaluate  
096 than CR. The MEI task can be viewed as a  
097 classification task in which any text span either  
098 refers to one of the input entities or the null class  
099 (*minor* entities and other non-mention spans). The  
100 classification formulation of MEI allows for the use  
101 of classification-based metrics that are more robust  
102 than the current CR metrics. Furthermore, MEI, by  
103 its definition, disregards insignificant and smaller  
104 clusters known to inflate the CR metrics (Moosavi  
105 and Strube, 2016; Lu and Ng, 2020; Kummerfeld  
106 and Klein, 2013). As an aside, formulating MEI as a  
107 classification task allows for a trivial parallelization

108 across candidate spans (Appendix A.1).

109 Finally, MEI’s explicit mapping of mentions to  
110 predefined entities improves its usability over CR in  
111 downstream applications that focus on mentions of  
112 specific entities. MEI effectively replaces tailored  
113 heuristics employed to extract CR cluster(s) refer-  
114 ring to entities of choice in such applications (entity  
115 understanding (Inoue et al., 2022), sentiment and  
116 social dynamics analysis (Zahiri and Choi, 2017;  
117 Antoniak et al., 2023)).

## 118 2 Task Formulation

119 **Notation.** For a document  $d$ , let  $\mathcal{E} = \{e_j\}_{j=1}^L$  be  
120 the set of  $L$  major entities that we wish to identify.  
121 We define  $\mathcal{M}_{\text{all}}$  as the set of all mentions that could  
122 refer to any entity and subsequently  $\mathcal{M}_j \subseteq \mathcal{M}_{\text{all}}$   
123 as the set of mentions that refer to a major entity  $e_j$ .  
124 Furthermore, we denote  $\mathcal{M} = \bigcup_j \mathcal{M}_j$  as the set  
125 of mentions that refer to one of the major entities  
126 while mentions that do not correspond to any major  
127 entity are designated as  $\mathcal{M}_{\text{other}} = \mathcal{M}_{\text{all}} \setminus \mathcal{M}$ .

128 **Task formulation.** In MEI, the input consists of  
129 the document  $d$  and designative phrases  $\mathcal{P} =$   
130  $\{p(e_j)\}_{j=1}^L$  where  $p(e_j)$  succinctly represents the  
131 entity  $e_j$ . For example, in Fig. 1, the phrases “*Al-*  
132 *addin*” and “*Mustapha*” uniquely represent Al-  
133 addin and his father who appear in “*Aladdin And*  
134 *The Wonderful Lamp*”. Note that in CR, the design-  
135 ative phrases  $\mathcal{P}$  are not part of the input.

136 In contrast to CR’s clustering foundations, MEI  
137 starts with a prior for each entity (the designative  
138 phrase) and can be formulated as an open set clas-  
139 sification, where every mention is either classified  
140 as one of the major entities or ignored. Formally,  
141 MEI aims to assign each mention  $m \in \mathcal{M}_j$  to  $e_j$   
142 and mentions  $m \in \mathcal{M}_{\text{other}}$  to  $\emptyset$ , a null entity.

## 143 3 Supervised MEI models

144 We propose MEIRa, Major Entity Identification via  
145 Ranking, which draws inspiration from the entity  
146 ranking formulation (Xia et al., 2021; Toshniwal  
147 et al., 2020) and maintains an explicit representa-  
148 tion for entities. The MEIRa models consist of 3  
149 steps: encoding the document, proposing candidate  
150 mentions, and an identification (id) module that  
151 tags mentions with major entities or the null entity.

152 **Document encoding** is performed using a  
153 Longformer-Large (Beltagy et al., 2020),  $\phi$ , that  
154 we finetune for the task. Mentions (or spans) are  
155 encoded as  $\mathbf{m}_i = \phi(m_i, d)$  by concatenating the

156 first, last, and an attention-weighted average of the  
 157 token representations within the mention span. In  
 158 MEI, an additional input is the set of designative  
 159 phrases  $\mathcal{P}$  for the major entities. Since each phrase  
 160 is derived from the document itself, we also obtain  
 161 its encoding using the backbone:  $\mathbf{e}_j = \phi(p(e_j), d)$ .

162 **Mention detection.** Similar to prior efforts (Toshniwal  
 163 et al., 2021), we use a mention proposal network  
 164 that predicts high-scoring candidate mentions. This  
 165 step finds all mentions  $\mathcal{M}_{\text{all}}$  and not just the ones  
 166 corresponding to the major entities  $\mathcal{M}$ . Training a  
 167 model to only detect mentions of major entities  
 168 would confuse it leading to poor performance.

169 **Identification module.** As illustrated in Fig. 2, we  
 170 initialize a working memory  $\mathcal{E}^W = [\mathbf{e}_j]_{j=1}^L$  as a  
 171 list of  $L$  major entities based on their designative  
 172 phrase representations. Given a mention  $m_i$ , the id  
 173 module computes the most likely entity as:

$$174 [s_i^*, e_i^*] = \max_{j=1 \dots L} f([\mathbf{m}_i, \mathbf{e}_j, \chi(m_i, e_j)]), \quad (1)$$

175 where  $f()$  is an MLP that predicts the score of tag-  
 176 ging mention  $m_i$  with the entity  $e_j$ , and  $\chi(m_i, e_j)$   
 177 encodes metadata. The output  $s_i^*$  corresponds to  
 178 the highest score and  $e_i^*$  is the top-scoring entity.  
 179 Based on the score,  $m_i$  is assigned to:

$$180 y(m_i) = \begin{cases} e_i^* & \text{if } s_i^* > \tau, \\ \emptyset & \text{otherwise,} \end{cases} \quad (2)$$

181 where  $\tau$  is a threshold (set to 0 in practice).

182 The metadata  $\chi(m_i, e_j)$  contains a distance (po-  
 183 sition) embedding representing the log distance be-  
 184 tween the mention  $m_i$  and the last tagged instance  
 185 of the entity  $e_j$ . If no mention is yet associated with  
 186 the entity, we use a special learnable embedding.

187 **Updates to the working memory.** We investigate  
 188 two approaches:

189 (i) **MEIRa-Static:** As the name suggests, the  
 190 working memory  $\mathcal{E}^W$  of the entity representations  
 191 remains constant ( $\mathcal{E}^{W(0)}$ ) and is not updated with  
 192 new mention associations. This makes the approach  
 193 highly parallelizable.

194 (ii) **MEIRa-Hybrid:** Similar to traditional CR,  
 195 this variation maintains a dynamic working memory  
 196  $\mathcal{E}^W$ , which is updated with every new mention-id  
 197 association. Specifically, assuming  $m_i$  is assigned  
 198 to  $e_j^*$ , the working memory would be updated using  
 199 a weighted mean operator  $g$  as  $\mathbf{e}_j \leftarrow g(\mathbf{e}_j, \mathbf{m}_i)$ ,  
 200 similar to Toshniwal et al. (2020). To prevent error  
 201 accumulation, we evaluate the mentions against

202  $\mathcal{E}^W$  and the initial entity representations ( $\mathcal{E}^{W(0)}$ ),  
 203 then compute the average score. This hybrid ap-  
 204 proach reaps benefits from both, the initial clean  
 205 designative phrases and the dynamic updates.

206 Following Toshniwal et al. (2020), the mention  
 207 detection and identification modules are trained end-  
 208 to-end using separate cross-entropy loss functions.

## 209 4 Few-shot MEI with LLMs

210 We propose a prompting strategy to leverage LLMs  
 211 for MEI, addressing their challenges in CR.

212 **Mention detection challenges.** CR or MEI can  
 213 be addressed using separate few-shot prompting  
 214 strategies for mention detection and mention clus-  
 215 tering/identification. However, Le and Ritter (2023)  
 216 found that this strategy faced significant challenges  
 217 with mention detection, performing worse than a  
 218 deterministic mention detector. Thus, they assume  
 219 access to an oracle mention detector and focus  
 220 their study to evaluating the linking capabilities of  
 221 LLMs.

222 An alternative is to use an external supervised  
 223 mention detector instead of the oracle. However,  
 224 this requires annotated training data and may not  
 225 align with a true few-shot LLM prompt paradigm.  
 226 Additionally, supervised mention detectors often  
 227 fail to generalize across CR datasets due to annota-  
 228 tion variability (Lu and Ng, 2020).

229 **MEI with LLMs.** We demonstrate that transition-  
 230 ing from CR to MEI addresses this gap in mention  
 231 detection and proposes an end-to-end, few-shot  
 232 prompting approach for MEI. Inspired by Dobro-  
 233 volskii (2021), we develop a prompting strategy  
 234 that first performs MEI at word-level (rather than  
 235 span), followed by a prompt to retrieve the span  
 236 corresponding to the word.

237 In addition to the document  $d$  and the set of  
 238 phrases  $\mathcal{P}$ , we also provide entity identifiers (e.g. #1,  
 239 #2) to the LLM. We will use the following example:  
 240 Document: *That lady in the BMW is Alice’s mom.*  
 241 Major Entities: 1. *Alice*; 2. *Alice’s mother*.

242 **Prompt 1. Word-level MEI.** Mention detection  
 243 with LLMs is challenging due to the frequent oc-  
 244 currence of nested mentions. We overcome this  
 245 by prompting the LLM to tag each word. Specifi-  
 246 cally, through few-shot examples, we ask the LLM  
 247 to detect and tag the **syntactic heads**<sup>1</sup> (e.g., *lady*,  
 248 *Alice*, *mom*) of mentions that refer to the major

<sup>1</sup>A syntactic head of a phrase is a word (*lady*) that is central to the characteristics of the phrase (*The lady in the BMW*).

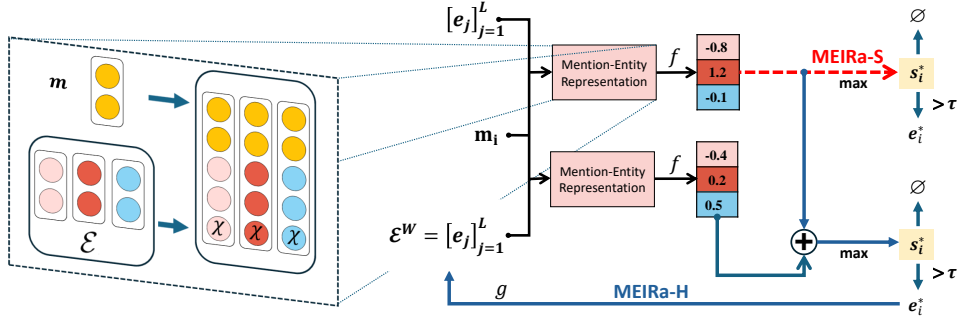


Figure 2: Identification module of MEIRa. A mention encoding  $m_i$  is concatenated with each entity’s embedding in  $\mathcal{E}^W$  and the metadata  $\chi(m_i, e_j)$ . Network  $f$  scores the likelihood of assigning  $m_i$  to each major entity. If the highest score  $s_i^*$  is above the threshold  $\tau$ ,  $m_i$  is associated with the highest scoring major entity  $e_i^*$  or discarded. In MEIRa-S, the entity memory  $\mathcal{E}^W$  remains static. For MEIRa-H (blue path), the assigned entity’s working memory is updated, and both the static (top half) and updated working memory (bottom half) are utilized to compute a final score.

entities. Other words are left untagged (implicitly assigned to  $\emptyset$ , the null entity). To create the few-shot examples, a contiguous set of words annotated with the same entity is considered as a span and its syntactic head is extracted using spaCy (Honnibal et al., 2020).

The ideal output for the example above is:

*“That lady#2 in the BMW is Alice#1’s mom#2..”*

Note that, even though the span “BMW” might be a valid mention, it is not annotated as it does not refer to one of the major entities. The exact prompt used for this is provided in the Appendix, Table 9.

**Prompt 2. Head2Span retrieval.** The entity tagged heads are passed to the Head2Span (H2S) module, along with the document to retrieve the span. The prompt consists of the document pre-annotated with the positions of the head, where each candidate head-word is followed by a “#” and is instructed to be replaced by the complete span (including any existent determiners and adjectives). For the input:

*That lady# in the BMW is Alice#’s mom#.*

the expected ideal output is

*That lady (That lady in the BMW) in the BMW is Alice(Alice’s)’s mom (Alice’s mom).*

Table 10 in the appendix shows the H2S prompt.

**Preserving structure.** We pose MEI as a structured generation task, prompting LLMs to reproduce documents and generate MEI tags at specific locations. Proprietary models like GPT-4 generally reproduce documents faithfully but for rare failures, we use the Needleman-Wunsch algorithm (Needleman and Wunsch, 1970) to align documents and extract tags. In the case of open-source models, we employ regular expression-based constrained decoding with the outlines library (Willard and Louf, 2023)<sup>2</sup>.

<sup>2</sup><https://outlines-dev.github.io/outlines/>

## 5 Experiments

**Datasets.** We evaluate three literary datasets chosen for their longer length and identifiable major entities, particularly the key narrative elements such as characters or plot devices. Table 1 compares statistical aspects of MEI and CR, revealing that MEI features fewer clusters (entities) but larger cluster sizes (more mentions per cluster).

(i) *LitBank* (Bamman et al., 2020) annotates coreference in 100 literary texts, each averaging around 2000 words. Following prior work (Toshniwal et al., 2021), we utilize the initial cross-validation split, dividing the documents into training, validation, and test sets with an 80:10:10 ratio.

(ii) *FantasyCoref* (Han et al., 2021) provides OntoNotes (Pradhan et al., 2013)-style<sup>3</sup> coreference annotations for 211 documents from Grimm’s Fairy Tales, with an average length of approximately 1700 words. The dataset includes 171 training, 20 validation, and 20 test documents.

(iii) *Additional Fantasy Text (AFT)* (Han et al., 2021) provides annotations for long narratives: (a) Aladdin (6976 words), (b) Ali Baba and the Forty Thieves (6911 words), and (c) Alice in Wonderland (13471 words).

**Metrics.** In contrast to CR, MEI facilitates the use of simple classification metrics. We define standard precision and recall for each major entity considered as an individual class of its own.

For a dataset  $\mathcal{D} = \{d_1, \dots, d_{|\mathcal{D}|}\}$ , the evaluation metrics are defined as follows:

$$\text{Macro-F1} = \frac{\sum_{d \in \mathcal{D}} \sum_{e_j \in \mathcal{E}_d} F1(e_j)}{\sum_{d \in \mathcal{D}} |\mathcal{E}_d|}, \text{ and} \quad (3)$$

<sup>3</sup>The exact guidelines are documented [here](#)

Model	FantasyCoref		LitBank	
	Macro-F1	Micro-F1	Macro-F1	Micro-F1
Coref-ID	72.5±2.2	78.8±2.7	79.7±2.7	80.6±3.7
Coref-CM	77.7±1.8	82.4±2.2	74.1±2.5	76.0±3.0
Coref-FM	77.9±1.7	83.2±2.2	77.4±2.3	80.6±4.7
MEIRa-S	<b>80.7±0.6</b>	<b>84.9±0.5</b>	80.8±0.8	81.8±1.0
MEIRa-H	80.3±1.4	84.3±2.0	<b>82.3±1.2</b>	<b>83.2±2.5</b>

Table 2: Results for models trained jointly on FantasyCoref and LitBank.

$$\text{Micro-F1} = \frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} \frac{\sum_{e_j \in \mathcal{E}_d} F1(e_j) \cdot |\mathcal{M}_j|}{\sum_{e_j \in \mathcal{E}_d} |\mathcal{M}_j|}. \quad (4)$$

Macro-F1 is the average F1-score of entities across the dataset, while Micro-F1 is the frequency-weighted F1-score of entities within a document, averaged across the dataset.

**Major entity selection.** We select as major entities, the top- $k$  entities ranked as per the frequency of occurrences. We use  $k=5$  for LitBank and FantasyCoref after visualizing the frequency plots of their training sets. For longer documents in AFT, we select up to 9 entities to ensure coverage of all key entities from the story. We also enforce that every entity  $e_j \in \mathcal{E}$  has a mention count  $|\mathcal{M}_j| \geq 5$ . We derive the representative span for each selected  $e_j$  from the set of mentions  $\mathcal{M}_j$  by selecting the most commonly occurring name or nominal mention.

### Implementation details.

*Supervised models:* Model hyperparameters are derived from [Toshniwal et al. \(2021\)](#). To ensure consistent performance across different numbers of target entities, we randomly select a subset of major entities at each training iteration (for more details, see Appendix A.2). All supervised models were trained five times with different random seeds, and we present aggregated results as the mean and standard deviation.

*LLMs:* We follow a few-shot prompting mechanism across the setups and experiments. Prompts that perform referential tasks consist of 3 examples of 6 sentences each. These 3 examples contain a mixture of narrative styles (narratives, dialogues), types of entities (major, non-major entities), categories of mentions (names, nominals, pronouns), and plurality. Additionally, before producing the MEI output, we ask the LLM to describe each major entity briefly. We find that this additional step improves performance. For the H2S prompt, we provide 9 sentences as examples, balancing the number of pre- and post-modifiers to the head. All

Model	FantasyCoref		LitBank	
	Macro-F1	Micro-F1	Macro-F1	Micro-F1
Coref-ID	63.4±1.8	69.5±3.6	58.0±2.4	57.7±1.0
Coref-CM	72.8±0.3	76.5±0.5	61.0±5.9	61.2±5.2
Coref-FM	71.2±1.5	75.2±1.3	66.1±2.1	67.1±3.9
MEIRa-S	<b>75.7±1.5</b>	78.5±1.2	74.6±1.1	74.7±1.6
MEIRa-H	74.7±1.0	<b>78.5±0.8</b>	<b>77.2±1.9</b>	<b>78.6±2.7</b>

Table 3: Results for models trained on OntoNotes.

examples were selected from LitBank’s train set and kept constant throughout the experiments. We set the temperature to 0 for all the models to ensure consistent and reproducible outputs.

## 5.1 Experiments: Supervised Models

**Baselines.** We train the fast-coref model ([Toshniwal et al., 2021](#)) for CR and perform the following three inference-time adaptations for MEI:

*Coref-ID:* fast-coref uses active lists of entity representations, resolving coreference by associating mentions with existing clusters or generating new ones. During inference, we disable the cluster creation step and pre-fill the entity list with the encoded vector representations of the major entities. Hence, all the detected mentions either get mapped to one of the major entities or are discarded.

*Coref-Cosine Map (Coref-CM):* Since coreference clusters obtained from fast-coref lack explicit entity association, we employ the Kuhn-Munkres (KM) algorithm ([Munkres, 1957](#)) to find the optimal matching cluster for each major entity. The cost matrix uses the cosine similarity between the encoded representation of the major entities and that of the predicted cluster embeddings, both derived from fast-coref.

*Coref-Fuzzy Map (Coref-FM):* This method uses the KM algorithm to derive optimal mappings by constructing a cost matrix from accumulated fuzzy-string matching scores between designative phrases and the predicted cluster’s mention strings.

**Supervised results.** In this experiment, we train MEIRa and the baseline models on the joint training set of LitBank and FantasyCoref. Subsequently, we assess their performance on the individual test sets, with results summarized in Table 2. Overall, MEIRa models consistently outperform the baselines on both metrics while also exhibiting better stability with a lower variance. The considerable variance observed in the performance of baseline methods across all experiments underscores the non-trivial nature of identifying clusters corresponding

Model	AFT	
	Macro-F1	Micro-F1
Coref-ID	68.1±5.9	78.7±6.1
Coref-CM	71.1±2.8	82.4±4.2
Coref-FM	71.1±4.7	83.2±4.7
MEIRa-S	81.6±1.4	88.8±1.3
MEIRa-H	<b>82.8±1.1</b>	<b>89.5±1.0</b>

Table 4: Results on the AFT dataset.

to major entities within the output clusters provided by the CR algorithms. MEIRa-H and MEIRa-S exhibit competitive parity on FantasyCoref (children stories), while MEIRa-H edges out on LitBank dataset, showcasing its adaptability in elaborate sentence constructions.

**Generalization across datasets.** To evaluate the generalization capabilities of MEIRa and baseline models, we train them on the OntoNotes dataset and then test their performance on LitBank and FantasyCoref. The results are presented in Table 3. When compared with Table 2, we observe a significant performance drop across the baseline models (*e.g.* for Coref-ID, the average Micro-F1 scores drop from 80.6 to 57.7 on LitBank). The performance gap for the baseline models is more pronounced on LitBank than on FantasyCoref because LitBank’s annotation strategies differ more significantly from those of OntoNotes. The observations aligns with previous work (Toshniwal et al., 2021), that showcase poor generalization of models trained for CR. In contrast, MEIRa models recover most of the underlying performance on both the datasets (MEIRa-H drops a little from 83.2 to 78.6 on LitBank Micro-F1), demonstrating MEI as a more adaptable task, bringing robustness over varying annotation strategies.

**Long documents.** Table 4 presents results on the AFT dataset of the models trained using a combined training set of LitBank and FantasyCoref. MEIRa models significantly outperform the baseline models, with MEIRa-H gaining 11.7% in Macro-F1 over the best baseline. The results demonstrate the efficacy of MEIRa models on resolving key entities in longer narratives.

**Computational performance.** MEIRa-S supports parallel batched processing since it does not update the working memory after associating mentions, *i.e.* the mentions need not be processed sequentially from left to right. Hence, post-mention detection (common to all models), MEIRa-S is about 25× faster than fast-coref when assessed across LitBank, FantasyCoref and AFT datasets on an

Model	FantasyCoref		LitBank	
	Macro-F1	Micro-F1	Macro-F1	Micro-F1
MEIRa-H	88.5	91.0	86.1	85.4
GPT-4	<b>90.7</b>	<b>92.0</b>	<b>88.8</b>	<b>91.6</b>
GPT-3.5	65.6	70.4	74.3	75.8
Code Llama-34B	63.4	70.8	68.3	72.7
Llama3-8B	50.5	57.8	46.3	52.1
Mistral-7B	62.1	71.1	61.2	70.9

Table 5: Few-shot LLM prompting results assuming the availability of ground-truth mentions.

NVIDIA RTX 4090 (see Fig. 3 in the appendix). Additionally, with the model’s small memory footprint during inference, the entire process can also be parallelized across chunks of documents making it extremely efficient. Hence, we pose MEIRa-S as a faster while competitive alternative to MEIRa-H (that requires dynamic updates and has similar computational performance as fast-coref).

## 5.2 Experiments: Few-shot prompting

**Models.** We experiment with GPT-4<sup>4</sup> (OpenAI, 2024), GPT-3.5<sup>5</sup>, Code Llama-34B (Rozière et al., 2024), Mistral-7B (Jiang et al., 2023), and Llama3-8B.<sup>6</sup> Following Le and Ritter (2023), we use the instruction-tuned versions for open-source models. These models were chosen for their ability to handle the extended context required for our benchmarks.

### 5.2.1 Linking Performance w/ Gold Mentions

We first evaluate all the models assuming the availability of an oracle mention detector. The experimental configuration is aligned with that of Le and Ritter (2023), albeit with the distinction that we assess them for the MEI task rather than for CR. The prompt used in our setup is provided in Table 11 of Appendix. For comparison, we also perform inference on golden mentions with MEIRa-H.

The results in Table 5 show that GPT-4 surpasses the supervised MEIRa-H model in this setup. Among LLMs, GPT-4 is easily the best-performing model. Code Llama-34B performs the best among open-source models, closely followed by Mistral-7B. While Code Llama-34B is tailored for the code domain, surprisingly, it outperforms strong LLMs suited for natural language. This result corroborates a similar finding by Le and Ritter (2023) for CR and related evidence regarding code pretraining aiding entity tracking (Kim et al., 2024). We find

<sup>4</sup>Specifically, gpt-4-1106-preview

<sup>5</sup>Specifically, gpt-3.5-turbo-1106

<sup>6</sup><https://ai.meta.com/blog/meta-llama-3/>

Model	FantasyCoref		LitBank	
	Macro-F1	Micro-F1	Macro-F1	Micro-F1
MEIRa-H	<b>80.3</b>	<b>84.3</b>	82.3	83.2
GPT-4 w/ Ext det	80.1	82.2	78.6	83.9
<b>GPT-4 with varying prompting strategies</b>				
Single prompt	63.0	66.2	64.4	72.8
Two-stage prompt	70.5	74.9	76.5	81.3
<b>Word-level MEI + spaCy H2S</b>				
GPT-4	77.4	79.4	<b>82.5</b>	<b>85.5</b>
GPT-3.5	50.1	54.4	60.1	63.1
Code Llama-34B	19.4	23.4	9.4	16.2
Llama3-8B	29.2	32.8	24.5	27.1
Mistral-7B	28.0	30.9	14.9	15.3

Table 6: Results on LLMs with different mention detection and linking strategies.

that Code Llama-34B performs close to GPT-3.5 for FantasyCoref, though a sizable gap remains for LitBank, potentially due to its linguistic complexity.

### 5.2.2 MEI Task Performance with LLMs

In this section, we present the results for the end-to-end MEI task using LLMs. We compare all the approaches from Section 4 and relevant baselines with the results summarized in Table 6. To limit the combinations of LLMs and approaches for our experiments, we first compare all the approaches in tandem with GPT-4 and then present results for the best-performing approach with other LLMs.

The first straightforward approach of using a *Single Prompt* to retrieve all the mentions of major entities in a single pass results in a significant performance drop compared to MEIRa-H (prompt in Table 12 of Appendix). The reason is that while GPT-4 outperforms MEIRa-H at mention linking, its mention detection performance, especially with nested mentions, is much worse compared to MEIRa-H.<sup>7</sup>

To further underscore the importance of mention detection, we also compare against the baseline *GPT-4 w/ Ext det*, which utilizes an external pre-trained mention detector followed by prompt-based linking (prompt in Table 11 of Appendix). We train the mention detector on the PreCo dataset (Chen et al., 2018), which achieves a 93.8% recall and 53.1% precision on the combined FantasyCoref and LitBank validation sets. We observe that *GPT-4 w/ Ext det* is almost at par with the fully supervised MEIRa-H, again highlighting the strong mention linking capabilities of GPT-4.

Next, we present the results of our proposed

<sup>7</sup>The failure to detect nested mentions is despite best efforts to provide illustrative examples in the few-shot prompt. Le and Ritter (2023) report similar findings with earlier GPT versions.

Error Type	MEIRa-H	GPT-4
Missing Major	162	793
Major-Major	210	154
Major-Other	243	0
Other-Major	200	516
Extra-Major	461	896
Total	1276	2359

Table 7: Breakdown of errors by MEIRa-H and GPT-4 on the combined LitBank and FantasyCoref test set.

*Two-stage prompt*, motivated by the *Single prompt* method’s failure with nested mentions. The first prompt asks GPT-4 to perform word-level MEI, by limiting the task to syntactic heads only. The second prompt then performs the task of mapping the identified syntactic heads to full mention spans. The results strongly validate our proposed approach with a relative improvement of more than 7% over the *Single prompt* method across all metrics and datasets. We also explore replacing the second step, i.e., head-to-span (H2S) retrieval, with an external tool. Specifically, we invert spaCy’s span-to-head mapping to obtain a head-to-span retriever.<sup>8</sup>

GPT-4 significantly improves in this setup, outperforming even the supervised model on LitBank. Given the strong performance of *GPT-4 + spaCy H2S*, we evaluate the open-source LLMs in only this setting. We observe a wide gap between GPT-4 and the open-source models. Llama3-8B surpasses other open-source models on both datasets, whereas the larger Code Llama-34B underperforms on the end-to-end task. This contrasts with the findings of the idealized golden mention setting, which assesses purely the model’s linking capabilities. The discrepancy between these results highlights the importance of evaluating in the realistic end-to-end setup.

### 5.3 Error Analysis

We classify MEI errors into five categories: (1) *Missing Major*: Not detecting a mention  $m \in \mathcal{M}$ . (2) *Major-Major*: Assigning a mention  $m \in \mathcal{M}_j$  to any other major entity  $\mathcal{E} \setminus e_j$ . (3) *Major-Other*: Assigning a mention  $m \in \mathcal{M}$  to  $\emptyset$ . (4) *Other-Major*: Assigning a mention  $m \in \mathcal{M}_{\text{other}}$  to any major entity in  $\mathcal{E}$ . (5) *Extra-Major*: Detecting extra mentions  $m \notin \mathcal{M}_{\text{all}}$  and assigning to any major entity in  $\mathcal{E}$ .

<sup>8</sup>For the test set gold mentions of the two datasets, there were only two cases where spans had the same head. We handled these two cases manually.

<b>Golden Mentions</b>	Presently [a small boy] <sub>0</sub> came walking along the path – [an urchin of nine or ten] <sub>0</sub> . . . . . [Winterbourne] <sub>1</sub> had immediately perceived that [he] <sub>1</sub> might have the honor of claiming [him] <sub>2</sub> as a fellow countryman. "Take care [you] <sub>2</sub> don't hurt [your] <sub>2</sub> teeth," [he] <sub>1</sub> said, paternally . . . . . [My] <sub>2</sub> mother counted them last night, and one came out right afterwards. She said she'd slap [me] <sub>2</sub> if any more came out. [I] <sub>2</sub> can't help it. It's this old Europe . . . . . If [you] <sub>2</sub> eat three lumps of sugar, [your] <sub>2</sub> mother will certainly slap [you] <sub>2</sub> ," [he] <sub>1</sub> said. "She's got to give [me] <sub>2</sub> some candy, then," rejoined [[his] <sub>1</sub> young interlocutor] <sub>2</sub> .
<b>GPT-4 Output</b>	Presently [a small boy] <sub>0</sub> came walking along the path – [an urchin of nine or ten] <sub>0</sub> . . . . . [Winterbourne] <sub>1</sub> had immediately perceived that [he] <sub>1</sub> might have the honor of claiming [him] <sub>2</sub> as a fellow countryman. "Take care you don't hurt your teeth," [he] <sub>1</sub> said, paternally . . . . . [My] <sub>2</sub> mother counted them last night, and one came out right afterwards. [She] <sub>2</sub> said [she] <sub>2</sub> d slap [me] <sub>2</sub> if any more came out. [I] <sub>2</sub> can't help it. [It] <sub>2</sub> 's this old Europe . . . . . If you eat three lumps of sugar, [your] <sub>2</sub> mother will certainly slap [you] <sub>2</sub> ," [he] <sub>1</sub> said. "[She] <sub>2</sub> 's got to give [me] <sub>2</sub> some candy, then," rejoined [his] <sub>2</sub> young interlocutor.
<b>MEIRa-H Output</b>	Presently a small boy came walking along the path – [an urchin of nine or ten] . . . . . [Winterbourne] <sub>1</sub> had immediately perceived that [he] <sub>1</sub> might have the honor of claiming [him] <sub>2</sub> as a fellow countryman. "Take care [you] <sub>2</sub> don't hurt [your] <sub>2</sub> teeth," [he] <sub>1</sub> said, paternally . . . . . [My] <sub>2</sub> mother counted them last night, and one came out right afterwards. She said she'd slap [me] <sub>2</sub> if any more came out. [I] <sub>2</sub> can't help it. It's this old Europe . . . . . If [you] <sub>2</sub> eat three lumps of sugar, [your] <sub>2</sub> mother will certainly slap [you] <sub>2</sub> ," [he] <sub>1</sub> said. "She's got to give [me] <sub>2</sub> some candy, then," rejoined [[his] <sub>1</sub> young interlocutor] <sub>2</sub> .

Table 8: Qualitative Analysis showcasing different errors made by GPT-4 and MEIRa-H. Errors are color-coded as follows: **Missing Major**, **Others-Major**, **Extra-Major**, **Major-Major**, and **Major-Other**.

Results combined over the LitBank and FantasyCoref test sets are presented in Table 7. Missing Major and Extra-Major contribute most of the errors for GPT-4, highlighting the scope for improvement in mention detection and span retrieval. Mention detection also remains a challenge in MEIRa-H, the model making most of the mistakes in the Extra-Major category. GPT-4 distinguishes major entities more clearly than MEIRa-H but tends to over-associate other mentions with major entities, resulting in higher Other-Major and Extra-Major errors. Note that GPT-4 has zero errors in the Major-Other category due to the prompt design, which only allows annotating major entities. Examples of these errors are visualized in Table 8.

## 6 Related Work

**Neural models for CR** have become the *de facto* choice in supervised settings (Lee et al., 2017; Kantor and Globerson, 2019; Joshi et al., 2020; Otmazgin et al., 2023). Efforts to enhance model efficiency include reducing candidate mentions to word-level spans (Dobrovolskii, 2021) and using single dense representations for entity clusters (Xia et al., 2021; Toshniwal et al., 2020).

**Generalization in CR** remains a lingering prob-

lem (Moosavi and Strube, 2017; Zhu et al., 2021; Porada et al., 2023). Current solutions include feature addition (Aralikatte et al., 2019; Otmazgin et al., 2023), joint training (Xia and Van Durme, 2021; Toshniwal et al., 2021), and active learning (Zhao and Ng, 2014; Yuan et al., 2022; Gandhi et al., 2023). Rather than relying on additional training data, we argue for an alternative formulation where the burden of domain adaptation is offloaded from training to inference.

**Evaluation of LLMs for CR** has largely been conducted in limited settings, such as the sentence-level Winograd Schema Challenges (WSC) (Brown et al., 2020), clinical pronoun resolution (Agrawal et al., 2022) and instance-level Q&A (Yang et al., 2022). Le and Ritter (2023) conducted the first document-level evaluation of LLMs for CR but assumed an oracle-mention detector. In contrast, we conduct end-to-end evaluations.

**Character Identification** deals with specific characters from transcripts of TV shows and trains a model tailored to these constrained inputs (Chen and Choi, 2016; Zahiri and Choi, 2017; Jiang et al., 2019). Baruah and Narayanan (2023) introduced a dataset annotated with referent mentions of specific characters of interest. We differ from these works by adopting a generalized task formulation independent of annotation strategies and entity selection.

## 7 Conclusion

CR models are limited in their generalization capabilities owing to annotation differences and general challenges of domain adaptation. We propose MEI as an alternative to CR, where the entities relevant to the input text are provided as input along with the text. Our experiments demonstrate that MEI is more suited for generalization than CR. Additionally, MEI can be viewed as a classification task that (a) enables the use of more robust classification-based metrics and (b) a trivially parallelizable model across document spans, which gives a 25x speedup over a comparable coreference model, making MEI more suitable for longer narratives. Unlike CR, the formulation of MEI allows few-shot prompted LLMs to effectively compete with trained models. Our novel two-stage prompting and robust baseline methods empower top-performing LLMs like GPT-4 to achieve this. Our analysis indicates that this task holds promise for effectively evaluating the long-context referential capabilities of LLMs in an end-to-end manner.



## 8 Limitations

Major Entity Identification (MEI) is proposed as a generalizable alternative to the coreference resolution (CR) task, and is not a replacement of CR. MEI limits itself to major entities and only caters to applications that are interested in a particular pre-defined set of entities. Our experiments follow certain thresholds that might not be universally applicable, and results and performance might vary slightly along this decision (refer Appendix A.2). Our current few-shot prompting evaluations are limited only to a few models that accommodate a large context window. Optimizing prompts and architecture to allow for a piece-wise aggregation of outputs across chunks of documents is left for future work.

## References

Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. Large Language Models are Few-Shot Clinical Information Extractors. In *EMNLP*.

Maria Antoniak, Anjalie Field, Jimin Mun, Melanie Walsh, Lauren Klein, and Maarten Sap. 2023. Riveter: Measuring Power and Social Dynamics Between Entities. In *ACL (Volume 3: System Demonstrations)*.

Rahul Aralikatte, Heather Lent, Ana Valeria Gonzalez, Daniel Herscovich, Chen Qiu, Anders Sandholm, Michael Ringgaard, and Anders Søgaard. 2019. Rewarding Coreference Resolvers for Being Consistent with World Knowledge. In *EMNLP-IJCNLP*.

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

Sabyasachee Baruah and Shrikanth Narayanan. 2023. Character Coreference Resolution in Movie Screenplays. In *Findings of ACL*.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The Long-Document Transformer. *arXiv preprint arXiv:2004.05150*.

Bernd Bohnet, Chris Alberti, and Michael Collins. 2023. Coreference Resolution through a seq2seq Transition-Based System. *TACL*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish,

Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *NeurIPS*.

Hong Chen, Zhenhua Fan, Hao Lu, Alan Yuille, and Shu Rong. 2018. PreCo: A Large-scale Dataset in Preschool Vocabulary for Coreference Resolution. In *EMNLP*.

Yu-Hsin Chen and Jinho D. Choi. 2016. Character Identification on Multiparty Conversation: Identifying Mentions of Characters in TV Shows. In *Annual Meeting of the Special Interest Group on Discourse and Dialogue*.

Bhuwan Dhingra, Qiao Jin, Zhilin Yang, William Cohen, and Ruslan Salakhutdinov. 2018. Neural Models for Reasoning over Multiple Mentions Using Coreference. In *NAACL-HLT*.

Vladimir Dobrovolskii. 2021. Word-Level Coreference Resolution. In *EMNLP*.

Nupoor Gandhi, Anjalie Field, and Emma Strubell. 2023. Annotating Mentions Alone Enables Efficient Domain Adaptation for Coreference Resolution. In *ACL*.

Sooyoun Han, Sumin Seo, Minji Kang, Jongin Kim, Nayoung Choi, Min Song, and Jinho D. Choi. 2021. FantasyCoref: Coreference Resolution on Fantasy Literature Through Omniscient Writer’s Point of View. In *Fourth Workshop on Computational Models of Reference, Anaphora and Coreference*.

Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. *spaCy: Industrial-strength Natural Language Processing in Python*.

Naoya Inoue, Charuta Pethe, Allen Kim, and Steven Skiena. 2022. Learning and Evaluating Character Representations in Novels. In *Findings of ACL*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7B. *arXiv preprint arXiv:2310.06825*.

Hang Jiang, Xianzhe Zhang, and Jinho D. Choi. 2019. Automatic Text-based Personality Recognition on Monologues and Multiparty Dialogues Using Attentive Networks and Contextual Embeddings.

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.

Ben Kantor and Amir Globerson. 2019. Coreference Resolution with Entity Equalization. In *ACL*.

Najoung Kim, Sebastian Schuster, and Shubham Toshniwal. 2024. Code Pretraining Improves Entity Tracking Abilities of Language Models.

723	Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text Generation from Knowledge Graphs with Graph Transformers. In <i>NAACL-HLT</i> .	Eva Sharma, Luyang Huang, Zhe Hu, and Lu Wang. 2019. An Entity-Driven Framework for Abstractive Summarization. In <i>EMNLP-IJCNLP</i> .	775
724			776
725			777
726			
727	Jonathan K. Kummerfeld and Dan Klein. 2013. Error-Driven Analysis of Challenges in Coreference Resolution. In <i>EMNLP</i> .	Shubham Toshniwal, Sam Wiseman, Allyson Ettinger, Karen Livescu, and Kevin Gimpel. 2020. Learning to Ignore: Long Document Coreference with Bounded Memory Neural Networks. In <i>EMNLP</i> .	778
728			779
729			780
730	Nghia T. Le and Alan Ritter. 2023. <a href="#">Are Large Language Models Robust Coreference Resolvers?</a>	Shubham Toshniwal, Patrick Xia, Sam Wiseman, Karen Livescu, and Kevin Gimpel. 2021. On Generalization in Coreference Resolution. In <i>Workshop on Computational Models of Reference, Anaphora and Coreference</i> .	782
731			783
732	Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. End-to-end Neural Coreference Resolution. In <i>EMNLP</i> .		784
733			785
734			786
735	Jing Lu and Vincent Ng. 2020. Conundrums in Entity Coreference Resolution: Making Sense of the State of the Art. In <i>EMNLP</i> .	Brandon T Willard and Rémi Louf. 2023. Efficient Guided Generation for LLMs. <i>arXiv preprint arXiv:2307.09702</i> .	787
736			788
737			789
738	Nafise Sadat Moosavi and Michael Strube. 2016. Which Coreference Evaluation Metric Do You Trust? A Proposal for a Link-based Entity Aware Metric. In <i>ACL</i> .	Congying Xia, Wenpeng Yin, Yihao Feng, and Philip Yu. 2021. Incremental Few-shot Text Classification with Multi-round New Classes: Formulation, Dataset and System. In <i>NAACL-HLT</i> .	790
739			791
740			792
741			793
742	Nafise Sadat Moosavi and Michael Strube. 2017. Lexical features in coreference resolution: To be used with caution. In <i>ACL</i> .	Patrick Xia and Benjamin Van Durme. 2021. Moving on from OntoNotes: Coreference Resolution Model Transfer. In <i>EMNLP</i> .	794
743			795
744			796
745	James Munkres. 1957. Algorithms for the assignment and transportation problems. <i>Journal of the society for industrial and applied mathematics</i> .	Xiaohan Yang, Eduardo Peynetti, Vasco Meerman, and Chris Tanner. 2022. What GPT Knows About Who is Who. In <i>Workshop on Insights from Negative Results in NLP</i> .	797
746			798
747			799
748	Saul B Needleman and Christian D Wunsch. 1970. A general method applicable to the search for similarities in the amino acid sequence of two proteins. <i>Journal of molecular biology</i> .	Michelle Yuan, Patrick Xia, Chandler May, Benjamin Van Durme, and Jordan Boyd-Graber. 2022. Adapting Coreference Resolution Models through Active Learning. In <i>ACL</i> .	801
749			802
750			803
751			804
752	OpenAI. 2024. <a href="#">GPT-4 Technical Report</a> .	Sayyed M. Zahiri and Jinho D. Choi. 2017. <a href="#">Emotion Detection on TV Show Transcripts with Sequence-based Convolutional Neural Networks</a> .	805
753	Shon Otmazgin, Arie Cattan, and Yoav Goldberg. 2023. LingMess: Linguistically Informed Multi Expert Scorers for Coreference Resolution. In <i>EACL</i> .		806
754			807
755		Wenzheng Zhang, Sam Wiseman, and Karl Stratos. 2023. Seq2seq is All You Need for Coreference Resolution. In <i>EMNLP</i> .	808
756	Ian Porada, Alexandra Olteanu, Kaheer Suleman, Adam Trischler, and Jackie Chi Kit Cheung. 2023. <a href="#">Investigating Failures to Generalize for Coreference Resolution Models</a> .		809
757			810
758		Shanheng Zhao and Hwee Tou Ng. 2014. Domain Adaptation with Active Learning for Coreference Resolution. In <i>Workshop on Health Text Mining and Information Analysis (Louhi)</i> .	811
759			812
760	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. 2013. Towards Robust Linguistic Analysis using OntoNotes. In <i>CONLL</i> .		813
761			814
762		Yilun Zhu, Sameer Pradhan, and Amir Zeldes. 2021. OntoGUM: Evaluating Contextualized SOTA Coreference Resolution on 12 More Genres. In <i>ACL-IJCNLP</i> .	815
763			816
764			817
765	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. <a href="#">Code Llama: Open Foundation Models for Code</a> .		818
766			
767			
768			
769			
770			
771			
772			
773			
774			

819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829

## A Appendix

### A.1 Linking Speed Comparison

This section compares the computational performance of fast-coref with the proposed MEIRa-S architecture. The classification formulation and the lack of an update step in MEIRa-S makes it a more efficient alternative to MEIRa-H and CR models. Fig. 3 displays the speed-up obtained in the identification module when assessed across documents with varying numbers of mentions. MEIRa-S consistently clocks a 20x efficiency across all ranges.

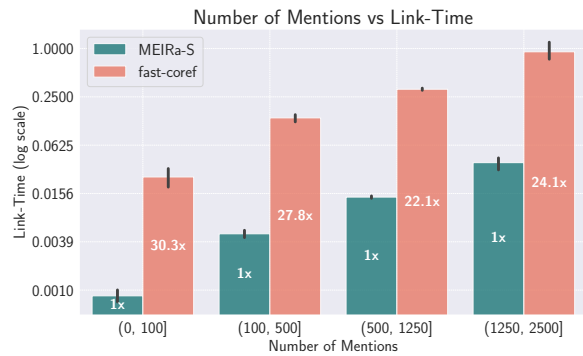


Figure 3: Linking speed comparison between MEIRa-S and fast-coref for the combined LitBank and FantasyCoref test set. There exists 6 documents with (0, 100] mentions, 19 with (100, 500] mentions, 5 with (500, 1250] mentions and 3 with (1250, 2500] mentions.

### A.2 Performance across number of entities

For consistency, the experiments of the main paper are evaluated across all the selected major entities (chosen using the thresholds defined in Section 5). A natural extension is to assess the model’s performance with varying numbers of entities of choice. For instance, if one is interested in only two key characters, can these models maintain consistency when provided with their designative phrases?

In this section, we address this concern and evaluate the MEI models with varying numbers of input entities. We present the per-entity F1-score of all entities across the AFT dataset. The results for MEIRa-H are showcased in Fig. 4, Fig. 5 and Fig. 6. The first column of the heatmap shows the per-entity F1-score when it is the sole target entity in the document. For e.g., the value in the first column in Fig. 4 corresponding to the entity *Baba Mustapha* (0.93) indicates the performance of the model when *Baba Mustapha* is the only target entity.

As we move across the columns of a particular row (ignoring the first column), the column number indicates the number of target entities used at

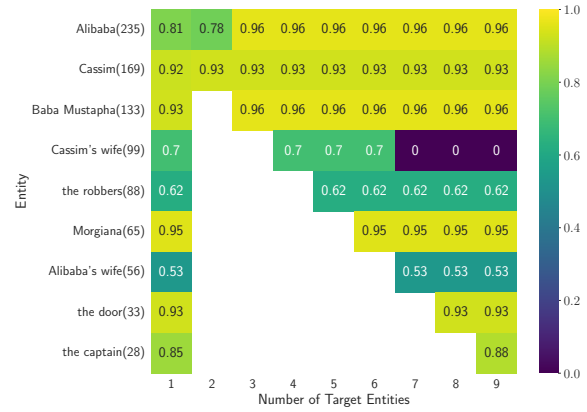


Figure 4: Performance of MEIRa-H across number of target entities for the document Ali Baba and the Forty Thieves.



Figure 5: Performance of MEIRa-H across number of target entities for Aladdin.

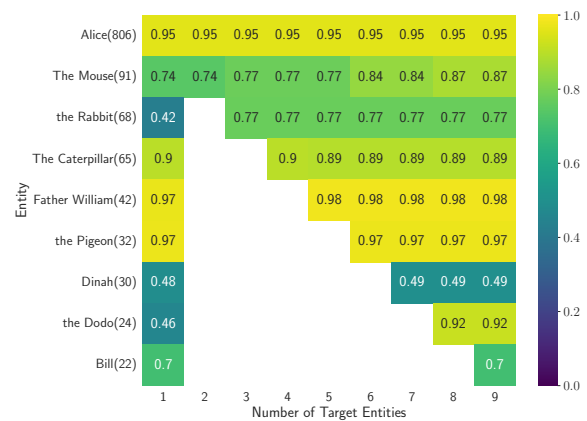


Figure 6: Performance of MEIRa-H across number of target entities for Alice in Wonderland.

inference. For instance, if the column number is  $k$ , the target entities are the top- $k$  frequent entities. Again, the 4<sup>th</sup> column in the row corresponding to *Baba Mustapha* indicates its individual F1-score in the experiment where the four input entities are *Alibaba*, *Cassim*, *Baba Mustapha* and *Cassim's*

853  
854  
855  
856  
857  
858

859 *wife*.

860 There are a few individual cases where the per-  
861 formance significantly varies with modifying the  
862 number of input entities. For example, *Cassim’s*  
863 *wife* is confused with *Alibaba’s wife* after the latter’s  
864 introduction. However, overall, the per-entity F1-  
865 score remains consistent across varying numbers  
866 of input entities across all three documents. These  
867 results demonstrate the effectiveness of MEIRa-H  
868 for applications requiring variable numbers of tar-  
869 get entities. This consistency is mainly due to the  
870 variable entity training, where a randomly chosen  
871 subset of major entities is selected in each iteration.  
872 Excluding this procedure leads to significant fluctu-  
873 ation in performance while modifying the number  
874 of target entities.

### 875 **A.3 Prompts**

876 We provide exact prompts for all the few-shot  
877 prompting experiments. Please note that not all  
878 the major entities listed in the few shot examples  
879 are necessary to be present in the text.

### 880 **A.4 Budget and Hardware details**

881 The supervised models were trained on a 24GB  
882 NVIDIA RTX 4090Ti GPU. For experiments with  
883 the open source language models, we used two  
884 48GB NVIDIA RTX A6000 GPU’s. For GPT-4  
885 and GPT-3.5 experiments, we spent approximately  
886 175\$ in total, covering both initial explorations and  
887 the computation of final results.

---

**Instruction**

---

You will receive a Text along with a list of Key Entities and their corresponding Cluster IDs as input. Your task is to perform Coreference Resolution on the provided text to categorize "each word belonging to a cluster" with its respective cluster id. Also briefly describe the key entities in 1-2 sentences before starting the coreference task.

Follow the format below to label a word with its cluster ID:

word#cluster\_id

Please keep in mind:

- Ensure the output adheres to the specified format for easy parsing.
  - Classify the words in the given text without altering any of the other content.
- 

**Example Input:**

---

Key Entities:

1. Katharine Hilbery (#katharine\_hilbery)
2. Mr. Denham (#mr.\_denham)
3. Mrs. Hilbery (#mrs.\_hilbery)
4. Mr. Hilbery (#mr.\_hilbery)
5. Mr. Fortescue (#mr.\_fortescue)

Text:

CHAPTER I It was a Sunday evening in October , and in common with many other young ladies of her class , Katharine Hilbery was pouring out tea . Perhaps a fifth part of her mind was thus occupied , and the remaining parts leapt over the little barrier of day which interposed between Monday morning and this rather subdued moment , and played with the things one does voluntarily and normally in the daylight . But although she was silent , she was evidently mistress of a situation which was familiar enough to her , and inclined to let it take its way for the six hundredth time , perhaps , without bringing into play any of her unoccupied faculties . A single glance was enough to show that Mrs. Hilbery was so rich in the gifts which make tea-parties of elderly distinguished people successful , that she scarcely needed any help from her daughter , provided that the tiresome business of teacups and bread and butter was discharged for her . Considering that the little party had been seated round the tea-table for less than twenty minutes , the animation observable on their faces , and the amount of sound they were producing collectively , were very creditable to the hostess . It suddenly came into Katharine 's mind that if someone opened the door at this moment he would think that they were enjoying themselves ; he would think , " What an extremely nice house to come into ! "

---

**Example Output:**

---

Description of Key Entities present in the text:

#1 - Katharine Hilbery: A young and apparently rich lady and the daughter of Mrs. Hilbery. She and Mrs. Hilbery were organising a party for some distinguished elders.

#3 - Mrs. Hilbery: She is the mother of Katharine Hilbery and is a well-to-do member of the society and a very efficient and able hostess

Coreference:

CHAPTER I It was a Sunday evening in October , and in common with many other young ladies of **her#1** class , **Katharine#1 Hilbery#1** was pouring out tea . Perhaps a fifth part of **her#1** mind was thus occupied , and the remaining parts leapt over the little barrier of day which interposed between Monday morning and this rather subdued moment , and played with the things one does voluntarily and normally in the daylight . But although **she#1** was silent , **she#1** was evidently mistress of a situation which was familiar enough to **her#1** , and inclined to let it take its way for the six hundredth time , perhaps , without bringing into play any of **her#1** unoccupied faculties . A single glance was enough to show that **Mrs.#3 Hilbery#3** was so rich in the gifts which make tea-parties of elderly distinguished people successful , that **she#3** scarcely needed any help from **her#3 daughter#1** , provided that the tiresome business of teacups and bread and butter was discharged for **her#1** . Considering that the little party had been seated round the tea-table for less than twenty minutes , the animation observable on their faces , and the amount of sound they were producing collectively , were very creditable to the **hostess#3** . It suddenly came into **Katharine#1 's#1** mind that if some one opened the door at this moment he would think that they were enjoying themselves ; he would think , " What an extremely nice house to come into ! "

---

Table 9: Prompt for WL Coreference

---

**Instruction**

---

Any word marked with # is supposed to be the head of a noun phrase. Expand this head to contain determiner and adjective phrases. Do not remove or add new words while expanding. Stick to the format.

---

**Example Input:**

---

Montraville# was a Lieutenant# in the army# : Belcour# was his brother officer# : they had been to take leave of their friends# previous to their departure for America# , and were now returning to Portsmouth# , where the troops# waited orders for embarkation

---

**Example Output:**

---

Montraville (**Montraville**) was a Lieutenant (**a Lieutenant in the army**) in the army (**the army**) : Belcour (**Belcour**) was his brother officer (**his brother officer**) : they had been to take leave of their friends (**their friends**) previous to their departure for America (**America**) , and were now returning to Portsmouth (**Portsmouth**) , where the troops (**the troops**) waited orders for embarkation

---

**Example Input:**

---

Arriving at the verge of the town# , he dismounted , and sending the servant# forward with the horses , proceeded toward the place# , where , in the midst of an extensive pleasure ground# , stood the mansion# which contained the lovely Charlotte Temple# .

---

**Example Output:**

---

Arriving at the verge of the town (**the town**) , he dismounted , and sending the servant (**the servant**) forward with the horses , proceeded toward the place (**the place**) , where , in the midst of an extensive pleasure ground (**an extensive pleasure ground**) , stood the mansion (**the mansion which contained the lovely Charlotte Temple**) which contained the lovely Charlotte Temple (**the lovely Charlotte Temple**) .

---

**Example Input:**

---

"You are a benevolent fellow# ," said a young officer# to him one day and I have a great mind to give you a fine subject to exercise the goodness of your heart upon.

---

**Example Output:**

---

"You are a benevolent fellow (**a benevolent fellow**) ," said a young officer (**a young officer**) to him one day and I have a great mind to give you a fine subject to exercise the goodness of your heart upon.

---

Table 10: Prompt for H2S Retrieval

---

**Instruction**

---

Annotate all the entity mentions in the following text with coreference clusters. Use Markdown tags to indicate clusters in the output, with the following format [mention] (#cluster\_name). Do not modify any text outside (), only add text inside parenthesis. The cluster names of the key entities are already provided, mark the mentions of the entity with the corresponding cluster name. Mark the mentions of the other entities with (#others). Also briefly describe the key entities in 1-2 sentences before starting the coreference task.

---

**Example Input:**

---

Key Entities:

1. Katharine Hilbery (#katharine\_hilbery)
2. Mr. Denham (#mr.\_denham)
3. Mrs. Hilbery (#mrs.\_hilbery)
4. Mr. Hilbery (#mr.\_hilbery)
5. Mr. Fortescue (#mr.\_fortescue)

Text:

CHAPTER I It was a Sunday evening in October, and in common with [many other young ladies of [her] (#) class] (#) , [Katharine Hilbery] (#) was pouring out tea . Perhaps a fifth part of [her] (#) mind was thus occupied , and the remaining parts leapt over the little barrier of day which interposed between Monday morning and this rather subdued moment , and played with the things one does voluntarily and normally in the daylight . But although [she] (#) was silent , [she] (#) was evidently [mistress] (#) of a situation which was familiar enough to [her] (#) , and inclined to let it take its way for the six hundredth time , perhaps , without bringing into play any of [her] (#) unoccupied faculties . A single glance was enough to show that [Mrs. Hilbery] (#) was so rich in the gifts which make tea-parties of [elderly distinguished people] (#) successful , that [she] (#) scarcely needed any help from [[her] (#) daughter] (#) , provided that the tiresome business of teacups and bread and butter was discharged for [her] (#) . Considering that [the little party] (#) had been seated round the tea-table for less than twenty minutes , the animation observable on [their] (#) faces , and the amount of sound [they] (#) were producing collectively , were very creditable to [the hostess] (#) . It suddenly came into [Katharine 's] (#) mind that if [some one] (#) opened the door at this moment [he] (#) would think that [they] (#) were enjoying [themselves] (#) ; [he] (#) would think , " What [an extremely nice house] (#) to come into ! "

---

**Example Output:**

---

Description of Key Entities present in the text:

#1 - Katharine Hilbery: A young and apparently rich lady and the daughter of Mrs. Hilbery. She and Mrs. Hilbery were organising a party for some distinguished elders.

#3 - Mrs. Hilbery: She is the mother of Katharine Hilbery and is a well-to-do member of the society and a very efficient and able hostess

Coreference:

CHAPTER I It was a Sunday evening in October , and in common with [many other young ladies of [her] (#katharine\_hilbery) class] (#others) , [Katharine Hilbery] (#katharine\_hilbery) was pouring out tea . Perhaps a fifth part of [her] (#katharine\_hilbery) mind was thus occupied , and the remaining parts leapt over the little barrier of day which interposed between Monday morning and this rather subdued moment , and played with the things one does voluntarily and normally in the daylight . But although [she] (#katharine\_hilbery) was silent , [she] (#katharine\_hilbery) was evidently [mistress] (#others) of a situation which was familiar enough to [her] (#katharine\_hilbery) , and inclined to let it take its way for the six hundredth time , perhaps , without bringing into play any of [her] (#katharine\_hilbery) unoccupied faculties . A single glance was enough to show that [Mrs. Hilbery] (#mrs.\_hilbery) was so rich in the gifts which make tea-parties of [elderly distinguished people] (#others) successful , that [she] (#mrs.\_hilbery) scarcely needed any help from [[her] (#mrs.\_hilbery) daughter] (#katharine\_hilbery) , provided that the tiresome business of teacups and bread and butter was discharged for [her] (#katharine\_hilbery) . Considering that [the little party] (#others) had been seated round the tea-table for less than twenty minutes , the animation observable on [their] (#others) faces , and the amount of sound [they] (#others) were producing collectively , were very creditable to [the hostess] (#mrs.\_hilbery) . It suddenly came into [Katharine 's] (#katharine\_hilbery) mind that if [some one] (#others) opened the door at this moment [he] (#others) would think that [they] (#others) were enjoying [themselves] (#others) ; [he] (#others) would think , " What [an extremely nice house] (#others) to come into ! "

---

Table 11: Prompt for evaluating linking performance

---

**Instruction**

---

Annotate all the entity mentions that refer to the key entities provided. The mention needs to include determiners and adjectives, if present. Use Markdown tags to indicate clusters in the output, with the following format [mention] (#cluster\_name). The cluster names of the key entities are already provided. Mark the mentions of the entity with the corresponding cluster name. Also briefly describe the key entities in 1-2 sentences before starting the coreference task.

---

**Example Input:**

---

Key Entities:

1. Katharine Hilbery (#katharine\_hilbery)
2. Mr. Denham (#mr.\_denham)
3. Mrs. Hilbery (#mrs.\_hilbery)
4. Mr. Hilbery (#mr.\_hilbery)
5. Mr. Fortescue (#mr.\_fortescue)

Text:

CHAPTER I It was a Sunday evening in October , and in common with many other young ladies of her class , Katharine Hilbery was pouring out tea . Perhaps a fifth part of her mind was thus occupied , and the remaining parts leapt over the little barrier of day which interposed between Monday morning and this rather subdued moment , and played with the things one does voluntarily and normally in the daylight . But although she was silent , she was evidently mistress of a situation which was familiar enough to her , and inclined to let it take its way for the six hundredth time , perhaps , without bringing into play any of her unoccupied faculties . A single glance was enough to show that Mrs. Hilbery was so rich in the gifts which make tea-parties of elderly distinguished people successful , that she scarcely needed any help from her daughter , provided that the tiresome business of teacups and bread and butter was discharged for her . Considering that the little party had been seated round the tea-table for less than twenty minutes , the animation observable on their faces , and the amount of sound they were producing collectively , were very creditable to the hostess . It suddenly came into Katharine 's mind that if someone opened the door at this moment he would think that they were enjoying themselves ; he would think , " What an extremely nice house to come into ! "

---

**Example Output:**

---

Description of Key Entities present in the text:

#1 - Katharine Hilbery: A young and apparently rich lady and the daughter of Mrs. Hilbery. She and Mrs. Hilbery were organising a party for some distinguished elders.

#3 - Mrs. Hilbery: She is the mother of Katharine Hilbery and is a well-to-do member of the society and a very efficient and able hostess

Coreference:

CHAPTER I It was a Sunday evening in October , and in common with many other young ladies of [her] (#katharine\_hilbery) class , [Katharine Hilbery] (#katharine\_hilbery) was pouring out tea . Perhaps a fifth part of [her] (#katharine\_hilbery) mind was thus occupied , and the remaining parts leapt over the little barrier of day which interposed between Monday morning and this rather subdued moment , and played with the things one does voluntarily and normally in the daylight . But although [she] (#katharine\_hilbery) was silent , [she] (#katharine\_hilbery) was evidently mistress of a situation which was familiar enough to [her] (#katharine\_hilbery) , and inclined to let it take its way for the six hundredth time , perhaps , without bringing into play any of [her] (#katharine\_hilbery) unoccupied faculties . A single glance was enough to show that [Mrs. Hilbery] (#mrs.\_hilbery) was so rich in the gifts which make tea-parties of elderly distinguished people successful , that [she] (#mrs.\_hilbery) scarcely needed any help from [[her] (#mrs.\_hilbery) daughter] (#katharine\_hilbery) , provided that the tiresome business of teacups and bread and butter was discharged for [her] (#katharine\_hilbery) . Considering that the little party had been seated round the tea-table for less than twenty minutes , the animation observable on their faces , and the amount of sound they were producing collectively , were very creditable to [the hostess] (#mrs.\_hilbery) . It suddenly came into [Katharine 's] (#katharine\_hilbery) mind that if some one opened the door at this moment he would think that they were enjoying themselves ; he would think , " What an extremely nice house to come into ! "

---

Table 12: Prompt for Direct version of E2E MEI