

# Arukikata Travelogue Dataset with Geographic Entity Mention, Coreference, and Link Annotation

Shohei Higashiyama<sup>1,2</sup>, Hiroki Ouchi<sup>2,3</sup>, Hiroki Teranishi<sup>3,2</sup>, Hiroyuki Otomo<sup>4</sup>,  
Yusuke Ide<sup>2</sup>, Aitaro Yamamoto<sup>2</sup>, Hiroyuki Shindo<sup>2,3</sup>, Yuki Matsuda<sup>2,3</sup>,  
Shoko Wakamiya<sup>2</sup>, Naoya Inoue<sup>5,3</sup>, Ikuya Yamada<sup>6,3</sup>, Taro Watanabe<sup>2</sup>

<sup>1</sup>NICT <sup>2</sup>NAIST <sup>3</sup>RIKEN <sup>4</sup>CyberAgent, Inc. <sup>5</sup>JAIST <sup>6</sup>Studio Ousia  
shohei.higashiyama@nict.go.jp, {hiroki.ouchi, ide.yusuke.ja6,  
yamamoto.aitaro.xv6, shindo.yukimat, wakamiya, taro}@is.naist.jp,  
hiroki.teranishi@riken.jp, otomo\_hiroyuki@cyberagent.co.jp,  
naoya-i@jaist.ac.jp, ikuya@ousia.jp

## Abstract

Geoparsing is a fundamental technique for analyzing geo-entity information in text, which is useful for geographic applications, e.g., tourist spot recommendation. We focus on *document-level* geoparsing that considers geographic relatedness among geo-entity mentions and present a Japanese travelogue dataset designed for training and evaluating document-level geoparsing systems. Our dataset comprises 200 travelogue documents with rich geo-entity information: 12,171 mentions, 6,339 coreference clusters, and 2,551 geo-entities linked to geo-database entries.

## 1 Introduction

Human activities, mobility, and events are often described with natural language expressions of locations or geographic entities (*geo-entities*), which indicate the geographic positions in the real world. This signifies the importance of technologies for extracting and grounding geo-entity expressions for various application domains, including tourism management, disaster management, and disease surveillance (Hu et al., 2022).

*Geoparsing* (Leidner, 2006; Gritta et al., 2020) is a fundamental technique involving two subtasks: *geotagging*, which identifies geo-entity mentions, and *geocoding*, which identifies corresponding database (DB) entries for (or the coordinates of) geo-entities. Notably, geoparsing, geotagging, and geocoding can be regarded as special cases of entity linking (EL), named entity recognition (NER) or mention recognition (MR), and entity disambiguation (ED), respectively.

This study focuses on geoparsing from the perspective of *document-level* analysis. Geo-entity mentions that co-occur in a document tend to be geographically close or related to each other; thus,

近鉄奈良駅<sup>FAC-NAME</sup><sub>(1)</sub> に到着。そこ<sup>DEICTIC</sup><sub>(1)</sub> から  
奈良公園<sup>FAC-NAME</sup><sub>(2)</sub> までは歩いてすぐです。  
お寺<sup>FAC-NOM</sup><sub>(GENERIC)</sub> が好きなので最初に興福寺<sup>FAC-NAME</sup><sub>(3)</sub>  
に行きました。境内<sup>FAC-NOM</sup><sub>(3)</sub> で鹿と遭遇し、  
奈良<sup>LOC-NAME</sup><sub>(4)</sub> に来たことを実感しました。

I arrived at Kintetsu Nara Station<sup>FAC-NAME</sup><sub>(1)</sub>.  
From there<sup>DEICTIC</sup><sub>(1)</sub> it's a short walk to  
Nara Park<sup>FAC-NAME</sup><sub>(2)</sub>. I like temples<sup>FAC-NOM</sup><sub>(GENERIC)</sub>  
so I first went to Kofukuji Temple<sup>FAC-NAME</sup><sub>(3)</sub>.  
I encountered a deer in the precincts<sup>FAC-NOM</sup><sub>(3)</sub> and  
felt that I had come to Nara<sup>LOC-NAME</sup><sub>(4)</sub>.

- (1) <https://www.openstreetmap.org/relation/11532920>
- (2) <https://www.openstreetmap.org/way/456314269>
- (3) <https://www.openstreetmap.org/way/1134439456>
- (4) <https://www.openstreetmap.org/relation/3227707>

Figure 1: Example illustration of an annotated document with English translation. Expressions underlined in blue indicate *geo-entity mentions*, superscript strings (e.g., FAC-NAME) indicate entity types of mentions, and subscript numbers (e.g., {1}) indicate coreference cluster IDs of mentions. URLs indicate OpenStreetMap entries that correspond to coreference clusters.

information about some geo-entity mentions can be useful for specifying information about other mentions. For example, by considering the context that describes a trip to Nara Prefecture, Japan, the mention of 興福寺 *kofukuji* ‘Kofukuji Temple’ in Figure 1 {3} can be disambiguated to refer to the temple in Nara rather than temples with the same name at different locations.

This paper presents a dataset suitable for document-level geoparsing: the Arukikata Travelogue Dataset with geographic entity Mention, Coreference, and Link annotation (ATD-MCL). Our dataset includes the three types of geo-entity

information illustrated in Figure 1: (1) spans and entity types of geo-entity mentions, (2) coreference relations among mentions, and (3) links from coreference clusters to corresponding entries in a geographic DB (geo-DB).

Our dataset has two desirable characteristics for document-level geoparsing. The first characteristic is that single travelogue documents in our dataset contain a rich amount of geo-entity mentions, in contrast to short documents, e.g., social media posts (Matsuda et al., 2017; Wallgrün et al., 2018). To leverage the inherent characteristic of the original travelogues, we have adopted an annotation policy to exhaustively markup geo-entity mentions, which refer to various locations and facilities expressed by named, nominal, and deictic expressions. The second characteristic is the *geographic continuity* among co-occurring mentions; that is, mentions that refer to nearby locations in the real world tend to appear near to one another within a document. Because travel records reflect the actual trajectories of travelers, this characteristic is more notable in travelogues than other text genres, e.g., news articles (Lieberman et al., 2010; Kamaloo and Rafiei, 2018; Gritta et al., 2018, 2020).

The potential applications of our dataset (and constructed geoparsers) include but not limited to tourism management applications. This is because geoparsing of location and facility mentions with diverse surface forms is essential for gaining a detailed understanding of where some event happened from text. For example, in disaster prevention/mitigation applications, it is crucial to specify detailed geographical positions by analyzing expressions other than named locations, utilizing geographic continuity if available, from social media posts about ongoing disasters and reports on past disasters.

As a result of manual annotation, our dataset comprises 12,273 sentences from the full text of 200 travelogue documents with 12,171 mentions, 6,339 coreference clusters (geo-entities), and 2,551 linked geo-entities.<sup>1</sup> Furthermore, we have conducted two types of evaluation using our dataset. First, we have measured inter-annotator agreement (IAA) for three types of information; the results indicate the practical quality of our dataset in terms of consistency. Second, we have evaluated current entity analysis systems on our dataset

<sup>1</sup>We conducted link annotation for 100 out of 200 documents including 3,208 geo-entities as described in §3.

for benchmarking baseline performance; the results demonstrate that reasonable performance can be achieved for MR and coreference resolution (CR), but performance has room for improvement in ED. We will release our annotated dataset at <https://github.com/naist-nlp/atd-mcl> and experimental codes at <https://github.com/naist-nlp/atd-mcl-baselines>.

## 2 Dataset Annotation

**Design Strategy** For building geoparsing datasets, it has been challenging to achieve a high coverage for facility entity mentions mainly because of the limited coverage of public geo-DBs, e.g., GeoNames.<sup>2</sup> To address this DB coverage problem, we adopt OpenStreetMap (OSM),<sup>3</sup> a free, editable, and large-scale geo-DB of the world. The usefulness of OSM has been continually increasing, as evidenced by the increase in node entries from over 1.5B in 2013 to over 80B in 2023.<sup>4</sup> Furthermore, we define entity types to cover broad types of location and facility mentions, including districts, buildings, landmarks, roads, and public transport lines and vehicles, as described in §2.2.

**Annotation Flow** Following the data preparation by the authors, annotation work was performed by native Japanese annotators at a professional data annotation company according to the three-step annotation flow: (1) mention annotation, (2) coreference annotation, and (3) link annotation.

### 2.1 Data Preparation

As raw text data, we adopted the ATD<sup>5</sup> (Arukikata Co., Ltd., 2022; Ouchi et al., 2023), which was constructed from user-posted travelogues written in Japanese. We first sampled documents about Japanese domestic travel with a reasonable document length (500–3000 characters, that is, approximately 300–1800 words) from the ATD. We then applied the GiNZA NLP Library<sup>6</sup> (Matsuda et al., 2019) to the raw text for sentence segmentation and automatic annotation of named entity (NE) mention candidates.

<sup>2</sup><https://www.geonames.org/>

<sup>3</sup><https://www.openstreetmap.org/>

<sup>4</sup><https://wiki.openstreetmap.org/wiki/Stats>

<sup>5</sup><https://www.nii.ac.jp/dsc/idr/arukikata/>

<sup>6</sup><https://github.com/megagonlabs/ginza>

Type and subtype	Example mentions
LOC-NAME LOC-NOM	奈良 ‘Nara’; 生駒山 ‘Mt. Ikoma’ 町 ‘town’; 島 ‘island’
FAC-NAME FAC-NOM	大神神社 ‘Omiwa Shrine’ 駅 ‘station’; 公園 ‘park’
LINE-NAME LINE-NOM	近鉄奈良線 ‘Kintetsu Nara Line’ 国道 ‘national route’; 川 ‘river’
TRANS-NAME TRANS-NOM	特急ひのとり ‘Ltd. Exp. Hinotori’ バス ‘bus’; フェリー ‘ferry’

Table 1: Examples of NAME and NOM entity mentions.

## 2.2 Mention Annotation

In the mention annotation step, we required the annotators to identify spans of geo-entity mentions in the documents, which may or may not refer to real-world locations, and assign entity type tags to the identified mentions by modifying the auto-annotated NE mentions. We adopted the brat annotation tool<sup>7</sup> (Stenetorp et al., 2012) for mention annotation (and succeeding coreference annotation).

The criteria for mention annotation define the *entity types* of geo-entity mentions, along with *mention spans* explained in Appendix B. Specifically, we define the following eight main entity types, which roughly correspond to Location, Facility, and Vehicle in Sekine’s Extended Named Entity (ENE) taxonomy (version 9.0)<sup>8</sup> (Sekine et al., 2002). (1) LOC, (2) FAC, and (3) TRANS respectively represent locations, facilities, and public transport vehicles; (4) LINE represents roads, waterways/streams, or public transport lines. The above four types are further divided into NAME and NOM subtypes, corresponding to whether a mention is named or nominal, as described in Table 1. (5) LOC\_ORG and (6) FAC\_ORG indicate location and facility mentions, respectively, that metonymically refer to organizations, e.g., ホテル *hoteru* in a sentence such as “The hotel serves its lunch menu.” (7) LOC\_OR\_FAC\_NOM indicates nominal mentions that can refer to both location and facility, e.g., 観光地 *kankōchi* ‘sightseeing spot.’ Finally, (8) DEICTIC indicates deictic expressions that refer to other geo-entity mentions or real-world locations, e.g., そこ *soko* ‘there’ in Figure 1.

## 2.3 Coreference Annotation

In the coreference annotation step, we required the annotators to assign mention-level *specificity*

*tags* or mention-pair-level *relations* to mentions identified in the previous step (except for those labeled with TRANS tags) using brat.

The criteria for coreference annotation define three types of specificity tags and two types of relations. As the representative cases, we introduce here the GENERIC specificity tag and the COREF coreference relation, and explain the remaining tags and relations in Appendix B. GENERIC is assigned to a generic mention, e.g., お寺 *otera* ‘temples’ in Figure 1, to distinguish singleton mentions that refer to real-world location, but are not coreferenced with other mentions. COREF is assigned to two mentions that both refer to the same real-world location, e.g., 近鉄奈良駅 *kintetsu nara eki* ‘Kintetsu Nara Station’ and そこ *soko* ‘there’ in Figure 1 ⟨1⟩. After relation annotation, a set of mentions that is sequentially connected through binary relations is regarded as one coreference cluster. A mention without any relations or specificity tags is regarded as a singleton, e.g., Figure 1 ⟨2⟩ and ⟨4⟩.<sup>9</sup>

## 2.4 Link Annotation

In the link annotation step, we required the annotators to link each coreference cluster to the URL of the corresponding OSM entry (e.g., ⟨1⟩–⟨4⟩ in Figure 1) on the basis of OSM and web search results. For URL assignment, the annotators added URLs to the cells representing coreference clusters in TSV files, which were converted from the brat output files.

The criteria for link annotation define the annotation flow as follows. For each coreference cluster, an annotator determines one or more normalized names of the referent location, e.g., formal or common name. The annotator then searches and assigns a URL of an appropriate OSM entry to the coreference cluster using search engines.<sup>10</sup>

The specific assignment process of entries is as follows. (a) If one or more candidate entries for a coreference cluster are found, assign the most probable candidate as BEST\_REF\_URL and (up to two) other possible candidates as SECOND\_REF\_URLS. (b) If the only candidate entry geographically includes but does not exactly match with the real-world ref-

<sup>9</sup>Although singleton mentions are marked with coreference cluster IDs in Figure 1 for clarity, singletons were not annotated with any coreference information in the actual work.

<sup>10</sup>Because it was sometimes difficult to find the desired entries using the Nominatim search engine available on the official OSM site, we asked the annotators to use additional search engines: web search engines and an original search engine that we developed.

<sup>7</sup><https://github.com/nlplab/brat>

<sup>8</sup><http://ene-project.info/ene9/?lang=en>

	#Doc	#Sent	#Word	#Men	#Ent
Set-A	100	5,949	85,741	6,052	3,131
Set-B	100	6,324	87,074	6,119	3,208
Total	200	12,273	172,815	12,171	6,339

Table 2: Statistics of the ATD-MCL.

erent, assign the found entry with the PART\_OF tag. (c) If no candidate entries are found in OSM, search and assign an appropriate entry from alternative DBs: Wikidata,<sup>11</sup> Wikipedia,<sup>12</sup> and general web pages describing the real-world referent. (d) If no candidate entries are found in any DBs, assign the NOT\_FOUND tag instead of an entry URL. The annotators can skip the search steps and assign the NOT\_FOUND tag when all member mentions and surrounding context do not provide any specific information that identifies the referent.

### 3 Dataset Statistics

The annotators first annotated 200 documents with mention information, then annotated the same 200 documents with coreference information, and finally annotated 100 documents, which were randomly sampled from the 200 documents, with link information.<sup>13</sup> We call the latter 100 documents that contain link annotation Set-B and refer to the remaining 100 documents without link annotation as Set-A. The numbers of documents (#Doc), sentences (#Sent), words (#Word), mentions (#Men), and entities (coreference clusters) (#Ent) in the ATD-MCL are listed in Table 2. We used Mode B (the middle unit) of the SudachiPy tokenizer (version 0.6.7)<sup>14</sup> (Takaoka et al., 2018) for counting the number of words in the Japanese text.

The notable characteristics of our dataset are summarized below. For more details, see Appendix C.

1. As shown in Table 3, facility mentions account for 50.3% (6,090/12,114) and nominal or demonstrative expressions account for 48.4% (5,867/12,114) of geo-entity mentions.<sup>15</sup>

<sup>11</sup><https://www.wikidata.org/>

<sup>12</sup><https://ja.wikipedia.org/>

<sup>13</sup>To construct the dataset within budget, we sampled 100 articles for link annotation, which is a heavy workload. It took 60, 70, and 200 hours to annotate 100 documents with mention, coreference, and link information, respectively.

<sup>14</sup><https://github.com/WorksApplications/SudachiPy>

<sup>15</sup>57 out of 12,171 mentions were non-geo-entity mentions, i.e., FAC\_ORG and LOC\_ORG.

	LOC	FAC	LINE	TRANS	GeoOther
NAME	2,289	3,239	462	257	–
NOM	861	2,851	582	666	–
Other	–	–	–	–	907
Total	3,150	6,090	1,044	923	907

Table 3: Tag distribution of geo-entity mentions in the whole dataset. “GeoOther” mentions consist of 372 LOC\_OR\_FAC\_NOM and 535 DEICTIC mentions. Non-geo-entity mentions (23 LOC\_ORG and 34 FAC\_ORG) are excluded from this table.

2. Multi-member clusters account for 35.6% (2,256/6,339) of coreference clusters, and the average number of member mention text types (distinct strings) for the multi-member clusters is 1.85, suggesting that the same geo-entity is often repeatedly referred to by named, nominal, and deictic expressions in a document (Appendix C.2 Table 12).
3. Geo-entities assigned with some URLs account for 97.1% (1,942/2,001) of entities with NAME mentions (“HasName” entities) and 50.5% (609/1,207) of the remaining entities, suggesting that identifying the referents that are not clearly written in text is difficult even for humans (Appendix C.3 Table 14).
4. Geo-entities assigned with OSM entry URLs account for 75.7% (1,514/2,001) of all “HasName” entities and 74.0% (811/1,096) of “HasName” facility entities, indicating that OSM has reasonable coverage of various types of locations in Japan (Appendix C.3 Table 15).

## 4 Inter-Annotator Agreement

For mention, coreference, and link annotation, we requested two annotators to independently annotate the same 10, 10, and 5 documents out of the 200, 200, and 100 documents, respectively; we simply selected 10 or five documents in ascending order based on document ID.<sup>16</sup> We measured the inter-annotator agreement (IAA) for the three annotation tasks.

### 4.1 Mention Annotation

As the IAA measure for mention annotation, we calculated the F1 scores between the results of two

<sup>16</sup>For coreference annotation, 10 documents annotated by two annotators did not include any mentions with specificity tags or mention pairs with attributive coreference relations.

Tag set	F1	Token			Type	
		#W1	#W2	#M	#W1	#W2
NAME	0.835	229	243	197	162	174
NOM	0.846	214	207	178	105	109
DEICT	0.621	19	10	9	6	3
ORG	0	1	0	0	1	0
All	0.832	463	460	384	274	283

Table 4: IAA for mention annotation. NAME, NOM, DEICT, and ORG indicate the (micro-averaged) scores for all NAME mentions, all NOM mentions, DEICTIC, and both LOC\_ORG and FAC\_ORG, respectively. The token and type columns indicate the scores and numbers based on token and type frequencies of mention text, respectively.

annotators (W1 and W2), based on exact match of both spans and tags.<sup>17</sup> Table 4 shows the F1 score for each tag set and the numbers of annotated mentions by W1, W2, and both (M).

The F1 score for all mentions was 0.832. Higher F1 score for NOM mentions (0.846) than that for NAME mentions (0.835) is probably because the less variety of NOM mention text types eased the annotation work for those mentions, as suggested by the mention token/type frequencies in Table 4.

## 4.2 Coreference Annotation

To assess IAA for COREF relation annotation, we used the metrics commonly used in coreference resolution studies: MUC (Vilain et al., 1995), B<sup>3</sup> (Bagga and Baldwin, 1998), CEAF<sub>e</sub> (Luo, 2005), and the average of the three metrics (a.k.a the CoNLL score) (Pradhan et al., 2012).<sup>18</sup>

Table 5 shows the F1 scores between two annotators’ (W1 and W2) results for each IAA measure and the numbers of clusters constructed from two annotators’ results for 2×2 settings: (a) original coreference clusters with all mentions or (b) clusters where only NAME mentions are retained, and (i) clusters with size ≥ 1 or (ii) clusters with size ≥ 2. In the basic setting (a)-(i), the average F1 score was 0.802. In addition, we observed two intuitive results. One is the lower scores for (a) than for (b), indicating that it was difficult to identify which mentions coreferenced with non-NAME mentions. The other is the higher scores for (i) than for

<sup>17</sup>We did not adopt a tag-level Kappa score regarding character-level BIO tags) because it would be biased toward being higher due to the majority of tags being O tags.

<sup>18</sup>A mention-level Kappa score can be calculated by regarding the task as, for example, classifying mentions into singleton or multi-member clusters. However, we did not adopt it because the resulting scores would be biased toward preferring singletons.

	MUC	B <sup>3</sup>	CEAF <sub>e</sub>	Avg.	#W1/#W2
(a) Original clusters with all mentions					
(i)	0.797	0.827	0.782	0.802	237/297
(ii)	0.797	0.768	0.811	0.792	91/79
(b) Clusters only with NAME mentions					
(i)	0.912	0.914	0.893	0.906	142/159
(ii)	0.912	0.868	0.844	0.874	46/46

Table 5: IAA between two annotators for coreference clusters in coreference annotation. The top two rows (a) and the bottom two rows (b) show the results in the described settings. (i) and (ii) show the results in the settings where singletons are included or not, respectively.

	F1	$\kappa$	#W1	#W2	#M
(a) Original URL					
URL	0.718	–	81	75	56
NOT_FOUND	0.737	–	16	22	14
All	0.722	0.707	97	97	70
(b) Grouped URL					
URL	0.821	–	81	75	64
NOT_FOUND	0.737	–	16	22	14
All	0.804	0.793	97	97	78

Table 6: IAA between two annotators for link annotation in (a) the original URL and (b) the grouped URL settings. The “URL” and NOT\_FOUND columns show the results for the assigned URLs and tag, respectively.

(ii); this is because leaving mentions as singletons is more likely to agree, since each mention is a singleton by default.

## 4.3 Link Annotation

As the IAA measure for link annotation, we calculated the F1 score and the Kappa score  $\kappa$  of OSM (or other DB) entry URL assignment for the same entities between two annotators (W1 and W2), which is similar to cluster-level hard F1 score (Zaporojets et al., 2022).<sup>19</sup>

Table 6 shows the agreement scores along with the numbers of entities to which URLs or the NOT\_FOUND tags were assigned by W1, W2, and both (M).<sup>20</sup> We used two settings about the equivalence for assigned URLs. (a) The original URL

<sup>19</sup>The same coreference information were provided to the annotators, but W1 and W2 merged or split three and one clusters, respectively, as a result of adopting the editable policy of clusters. We then evaluated link agreement only for clusters in which all members matched between the two annotators’ results.

<sup>20</sup>We regarded an entity as a matched URL instance when both annotators assigned the same URL and as a matched NOT\_FOUND instance when both annotators assigned NOT\_FOUND.

setting compares raw URL strings assigned by the annotators. (b) The grouped URL setting treats OSM entries or web pages representing practically the same locations as the same and compares the grouped URL sets instead of original URLs.<sup>21</sup>

The F1 scores for URLs and NOT\_FOUND were over 0.7 in both settings, indicating that the annotator could assign the same URL (or the NOT\_FOUND tag) to the majority of geo-entities in spite of the huge number of candidate URLs. The lower agreement scores in (a) the original setting than those in (b) the grouped setting is because the annotators assigned different but practically equivalent entry URLs to eight entities.

## 5 Experiments

We conducted experiments on the ATD-MCL for three tasks: MR, CR, and ED. The purpose of the experiments is to clarify the performance level of current entity analysis systems, including off-the-shelf and finetuned models, on our dataset.

### 5.1 Data Split

We regarded all Set-A documents as train-a and split the Set-B documents into train-b, development, and test sets at a ratio of 10:10:80. The union of train-a and train-b was used as the training set for both MR and CR, whereas train-b was used as the training set for ED. Thus, the data split of 110:10:80 was used for MR and CR, and that of 10:10:80 was used for ED. We determined to assign the large part of datasets to the test set to obtain less biased and more reliable evaluation results.<sup>22</sup>

### 5.2 Database Preprocessing

To the OSM data file consisting of Japanese domestic location entries,<sup>23</sup> we applied preprocessing to group together entries that refer to almost the same real-world locations by assigning the same group ID string, which resulted in 1.8M entry groups. Thus, we adopted a setting where entry groups are considered as linking units rather than individual entries. Detailed processing is described in Appendix D.3.

<sup>21</sup>The first author manually judged the practical equivalence of different OSM entries and web pages for entities unmatched between two annotators.

<sup>22</sup>The unsupervised ED systems in our experiments did not actually use any training examples. Different data split that includes more training examples can also be useful for future experiments involving supervised ED systems.

<sup>23</sup>We used `japan-230601.osm.bz2`, which was available at <http://download.geofabrik.de/asia/>.

Examples of entry group IDs are as follows.

- “name=スターバックス|branch=None|prefecture=奈良県|city=奈良市|quarter=樽井町|road=猿沢遊歩道|amenity=cafe” (Starbucks Coffee at Sarusawa pathway, Tarui-cho, Nara City, Nara Prefecture)
- “name=ローソン|branch=京王多摩川駅|prefecture=東京都|city=調布市|shop=convenience” (Lawson Keio Tamagawa Station store at Chofu City, Tokyo Prefecture)
- “name=首都高速湾岸線|prefecture=千葉県,東京都,神奈川県|city=None|route=road” (The Metropolitan Expressway Bayshore Route passing through Chiba Prefecture, Tokyo Prefecture, and Kanagawa Prefecture)
- “name=JR予讃線|prefecture=愛媛県,香川県|city=None|route=railway” (The JR Yoson line passing through Ehime Prefecture and Kagawa Prefecture)

Whereas the first two groups contain only one entry, the third and fourth groups contain 140 and 718 entries, respectively.

### 5.3 Mention Recognition

**Task Setting** We treat MR as the task of identifying spans and entity types of mentions in given documents. As the evaluation measure, we use the F1 score between the gold and predicted mentions based on exact match of both spans and entity types.

**Systems** We evaluated two systems that we finetuned models on our training set (spaCy-MR and mLUKE-MR) and two off-the-shelf systems without model finetuning (KWJA and GiNZA). spaCy-MR indicates a transition-based parsing model on the spaCy NLP library<sup>24</sup> that we built using a pretrained Japanese ELECTRA (Clark et al., 2020) model.<sup>25</sup> This corresponds to the finetuned version of the GiNZA model. mLUKE-MR is our implementation of a span-based MR system using a pretrained multilingual LUKE (mLUKE) (Ri et al., 2022) model.<sup>26</sup> As the off-the-shelf systems, we used KWJA “base” (version 2.1.1)<sup>27,28</sup> (Ueda et al., 2023) and GiNZA “ja\_ginza\_electra” (version 5.1.2). GiNZA and KWJA follow the ENE and IREX (Sekine and Isahara, 2000) tag sets, which are different from ours. Thus, we applied

<sup>24</sup><https://spacy.io/api/architectures#parser>

<sup>25</sup><https://huggingface.co/megagonlabs/transformers-ud-japanese-electra-base-discriminator>

<sup>26</sup><https://huggingface.co/studio-ousia/mluke-large-lite>

<sup>27</sup><https://github.com/ku-nlp/kwja>

<sup>28</sup>There was no KWJA documentation describing how to train a custom model, and we attempted but failed to perform training/finetuning.

System	Tag	P	R	F1
KWJA	Overall	0.279	0.352	0.311
	NAME	0.279	0.695	0.398
GiNZA	Overall	0.574	0.277	0.374
	NAME	0.574	0.548	0.560
spaCy-MR	Overall	0.752	0.732	0.742
	NAME	0.733	0.719	0.726
	NOM	0.790	0.753	0.771
	DEICTIC	0.645	0.721	0.681
	ORG	0.353	0.250	0.293
mLUKE-MR	Overall	<b>0.813</b>	<b>0.817</b>	<b>0.815</b>
	NAME	0.828	0.813	0.821
	NOM	0.826	0.818	0.822
	DEICTIC	0.616	0.896	0.730
	ORG	0.833	0.417	0.556

Table 7: System performance for mention recognition: precision (P), recall (R), and F1.

tag conversion rules to their outputs. Because the LOCATION tag in IREX semantically includes LOC\_NAME, FAC\_NAME, and LINE\_NAME tags, we converted each KWJA output mention with the LOCATION tag into three mention instances with the same span and with one of the three tags, which prioritizes recall over precision. More detailed settings are described in Appendix D.

**Results** Table 7 shows the performance of the MR systems for the test set. The off-the-shelf systems, GiNZA and KWJA, achieved the recall of 0.55–0.70 for NAME mentions, indicating moderate coverage for named geo-entity mentions. However, the two systems failed to extract non-NAME mentions (the F1 scores were 0), which is natural because these systems had been trained on only NE annotations (not nominal phrases). Owing to our finetuning, spaCy-MR and mLUKE-MR improved the performance: the overall F1 scores of 0.74–0.82. More specifically, both finetuned models achieved F1 scores of 0.73–0.82 for NAME and NOM, but they exhibited lower F1 scores for DEICTIC and ORG. These results are likely because it is difficult for the models to learn from a limited number of training examples whether DEICTIC mentions refer to real-world locations or not, and whether ORG mentions metonymically refer to organizations or not. For the fine-grained results for each tag, see Appendix E.

## 5.4 Coreference Resolution

**Task Setting** We define CR as the task of clustering the given gold mentions that corefer the same real-world locations. We use the same evaluation

System	Size	MUC	B <sup>3</sup>	CEAF <sub>e</sub>	Avg.
Rule-CR-1	≥ 1	0	0.755	0.639	0.465
	≥ 2	0	0	0	0
Rule-CR-2	≥ 1	0.622	0.840	0.790	0.750
	≥ 2	0.622	0.613	0.629	0.621
KWJA	≥ 1	0.694	0.839	0.793	0.775
	≥ 2	0.694	0.661	0.658	0.671
mLUKE-CR	≥ 1	<b>0.753</b>	<b>0.875</b>	<b>0.839</b>	<b>0.822</b>
	≥ 2	<b>0.753</b>	<b>0.733</b>	<b>0.737</b>	<b>0.741</b>

Table 8: System performance for coreference resolution.

metrics as the IAA measures.

**Systems** We evaluated one finetuned system (mLUKE-CR), one off-the-shelf system (KWJA), and two rule-based systems (Rule-CR-1 and 2). mLUKE-CR is our implementation of an end-to-end CR model based on a pretrained mLUKE model,<sup>29</sup> which identifies the antecedent (preceding coreference mention) for a given mention following Lee et al. (2017). We used the KWJA ‘base’ model and applied a modification rule to the KWJA’s output clusters so that the union of all output clusters matched the set of all gold mentions.<sup>30</sup> Simple rule-based systems are as follows. Rule-CR-1 treats all given mentions as singletons. Rule-CR-2 groups together sets of mentions with the same surface form in a document into clusters and treats the remaining mentions as singletons.

**Results** Table 8 shows the performance of the CR systems for the test set. The simplest rule-based system, Rule-CR-1, appears to have achieved the moderate B<sup>3</sup> and CEAF<sub>e</sub> scores for clusters with size ≥ 1 (although resulted in the zero score for the link-based MUC metric), due to the dataset distribution biased toward a high population of singletons. Thus, it is necessary to pay attention to the improvement from these baseline scores as meaningful performance evaluation measures. Another rule-based system, Rule-CR-2, achieved the scores of 0.61–0.84 for the three metrics, indicating that the simple heuristic regarding surface forms was a strong clue for finding coreferent mentions. The superior performance of KWJA and mLUKE-CR over Rule-CR-2 indicates that these two systems

<sup>29</sup><https://huggingface.co/studio-ousia/mluke-large>

<sup>30</sup>The modification rule removes predicted mentions that do not match any gold mentions from the output clusters and adds gold mentions that do not match any predicted mentions as singletons on the basis of mention span overlapping.

System	R@1	R@5	R@10	R@100
Rule-ED	0.221	0.323	0.345	0.362
BERT-ED	0.245	0.401	0.443	0.555

Table 9: System performance for entity disambiguation.

identified (part of) coreferent mentions with different surface forms, although mLUKE-CR expectedly performed better owing to finetuning.

## 5.5 Entity Disambiguation

**Task Setting** We define ED as the task of selecting appropriate entry group IDs from all entry groups for each given geo-entity. As the evaluation measure, we use recall@ $k$  ( $R@k$ ) for the given entities; the prediction is regarded as correct if one of the predicted  $k$  entity groups contains the gold OSM entry URL for each geo-entity.

**Systems** We evaluated an unsupervised system (BERT-ED) and a rule-based system (Rule-ED). For an input entity, both systems regard the longest mention surface among its member mentions with NAME entity subtype tags as the entity name and predict DB entry groups based on the entity name. The systems return no entry groups if the entity contains no NAME mentions. BERT-ED is our implementation of an ED system without hyperparameters based on a pretrained Japanese BERT (Devlin et al., 2019) model.<sup>31</sup> BERT-ED calculates the similarity between each entity’s name and “name” attribute value of each candidate entry group, and then ranks the candidates. For the similarity score, we used the cosine similarity score between vector representations, that is, the average of hidden states at the last layer for input words within the name string.<sup>32</sup> Rule-ED extracts entry groups whose “name” attribute values exactly match the entity’s name for each given entity, and then ranks them in lexicographic order of full group ID strings.

**Results** Table 9 presents the performance of the ED systems for the test set. Overall, BERT-ED achieved better scores than Rule-ED owing to soft matching and ranking using vector representations. In particular, BERT-ED outperformed Rule-ED by a larger margin on  $R@k$  with larger  $k$ . Although this result suggests the effectiveness of vector rep-

<sup>31</sup><https://huggingface.co/cl-tohoku/bert-base-japanese-whole-word-masking>

<sup>32</sup>We also tried an entity representation calculated from the full sentence where its representative mention occurred, but confirmed its poor performance.

resentations, the performance for  $R@1$  can be improved by introducing more sophisticated disambiguation strategies that consider the geography-related content in a document, including location and facility types identified from the surrounding context, and geographic areas mentioned within the document.

## 5.6 Discussion

For MR and CR, the finetuned systems achieved the reasonable performance in our experiments. For ED, in contrast, the simple unsupervised systems did not achieve practical performance. A possible solution is training supervised ED systems on in-domain training data. However, we suppose that predicting appropriate DB entries for unknown instances would remain a main challenge due to limits to improving coverage by increasing training instances.

Another challenge in geographic ED is that natural language descriptions of geo-DB entries are unavailable, different from general DBs represented by Wikipedia. This also makes it difficult to directly apply state-of-the-art general ED systems using entry description text (Wu et al., 2020; Yamada et al., 2022) to geographic ED, i.e., geocoding. Instead, OSM entries have rich information of semantic attributes and geographic relations, such as distance and hierarchy. A prospective direction is learning mention/entry representations that leverage or encode such geographic information, as well as entity type and population information (Zhang and Bethard, 2023). For example, if some geographic relations between two mentions are indicated by calculation based on their representations, geo-entities referred to by them may also have similar relations, which would be useful for CR and ED.

## 6 Related Work

**Entity Analysis Datasets** For over two decades, efforts have been devoted to developing annotated corpora for English entity analysis tasks, including NER (Tjong Kim Sang, 2002; Ling and Weld, 2012; Baldwin et al., 2015), anaphora/coreference resolution (Grishman and Sundheim, 1996; Doddington et al., 2004; Pradhan et al., 2011; Ghaddar and Langlais, 2016), and ED and EL (McNamee et al., 2010; Hoffart et al., 2011; Ratinov et al., 2011; Rizzo et al., 2016). For Japanese text, annotated corpora have been developed for general



Dataset Name		Lang	Text Genre	Geo-DB	#Men	Facility	Nominal
LGL Corpus	(Lieberman et al., 2010)	en	News	GeoNames	4.8K	✗	✗
TR-News	(Kamalloo and Rafiei, 2018)	en	News	GeoNames	1.3K	✗	✗
GeoVirus	(Gritta et al., 2018)	en	News	Wikipedia	2.2K	✗	✗
GeoWebNews	(Gritta et al., 2020)	en	News	GeoNames	2.7K	△	✓
SemEval-2019 T12	(Weissenbacher et al., 2019)	en	Science	GeoNames	8.4K	✗	✗
CLDW	(Rayson et al., 2017)	en	Historical	Unlock	3.7K	△	✗
GeoCorpora	(Wallgrün et al., 2018)	en	Microblog	GeoNames	3.0K	△	✗
LRE Corpus	(Matsuda et al., 2017)	ja	Microblog	ISJ & Orig.	1.0K	△	✓
ATD-MCL	(Ours)	ja	Travelogue	OSM	12.3K	✓	✓

Table 10: Characteristics of representative geoparsing datasets and ours. “#Men” indicates the number of annotated mentions in each dataset. The facility and nominal columns show the availability of geoparsed facility mentions and nominal mentions, respectively: ✓ (available), ✗ (not available), and △ (available to a limited extent).

NER (Sekine et al., 2002; Hashimoto and Nakamura, 2010; Iwakura et al., 2016), coreference resolution (Kawahara et al., 2002; Hashimoto et al., 2011; Hangyo et al., 2014), and EL (Jargalsaikhan et al., 2016; Murawaki and Mori, 2016).

**Geoparsing Datasets** Table 10 summarizes the characteristics of representative geoparsing datasets and the ATD-MCL. For English geoparsing, annotated corpora have been developed and used as benchmarks for system evaluation. The Local Global Lexicon (LGL) Corpus (Lieberman et al., 2010), TR-News (Kamalloo and Rafiei, 2018), and GeoWebNews (Gritta et al., 2020) contain approximately 100–600 news articles from global and local news sources. Although GeoWebNews contains facility mentions, which account for 8% of the total, Gritta et al. (2020) estimated their coordinates using the Google Maps API due to the absence of GeoNames entries, and excluded them from their experiments. GeoVirus (Gritta et al., 2018) comprises 229 WikiNews articles focusing on viral infections. The SemEval-2019 Task 12 dataset (Weissenbacher et al., 2019) comprises 150 biomedical journal articles on the epidemiology of viruses. The GeoCorpora project (Wallgrün et al., 2018) constructed a geo-microblog corpus that comprises 6,711 tweets with the very limited amount of facility mentions.<sup>33</sup> The Corpus of Lake District Writing (CLDW) (Rayson et al., 2017) consists of 80 historical texts, including travelogues and tourist guidebooks. The location and facility mentions in their gold standard subset of 28 texts were manually checked, but the coordinates were not. For Japanese geoparsing, Matsuda et al. (2017) constructed the Location Reference Expres-

<sup>33</sup>According to their supplemental material, the proportion of mentions referring to facilities, such as buildings and airports, is less than 3%.

sions (LRE) corpus, comprising 10,000 Japanese tweets, 951 of which have geo-entity-related tags. They used Ichi Sansho Joho (ISJ) ‘City-block-level location reference information’ and their original gazetteer of facilities, but the latter gazetteer has not been available due to licensing reasons.

## 7 Conclusion

This paper has described the ATD-MCL dataset, which is designed for document-level geoparsing, along with the annotation criteria, IAA assessment, and performance evaluation of the baseline systems. Our dataset enables other researchers to conduct reproducible experiments through the public release of our annotated data. We expect that our dataset contributes to fostering future research and advancing geoparsing techniques.

In future work, we plan to (1) develop a document-level geoparser that leverages both characteristics of geo-entity mentions in text and geo-DB entries, (2) enhance our dataset with additional semantic information, such as the movement trajectories of travelogue writers, for more advanced analytics, and (3) construct annotated travelogue datasets in other languages by extending our annotation guidelines.

## Limitations

**Optimization of Database Preprocessing** As the preprocessed DB for ED, we used 2.8M OSM entries of Japanese domestic locations with “name” attributes. While checking a portion of the generated entry groups, we performed rule engineering to make the original DB more desirable for our ED task, which means entries that can be regarded as practically equivalent to each other belong to the same groups. Over- and under-aggregated groups in the final DB could produce the evaluation results

with underestimated or overestimated system performance. This would have a greater influence on the recall@ $k$  scores with smaller  $k$  for evaluating disambiguation accuracy, but a lesser influence on the scores with larger  $k$  for evaluating extraction coverage.

**Optimization of System Performance** We performed not systematic but minimum hyperparameter search for mLUKE-based models due to time and resource limitations. Similarly, we used the fixed hyperparameters for spaCy-MR, which correspond to those used for GiNZA. Thus, performing optimized experiments has potential for further performance improvement in these systems.

**Independent Experiments on Geoparsing Subtasks** As a first step toward comprehensive evaluation of geoparsing techniques, we independently evaluated the baseline systems on each subtask in the gold input setting; that is, gold mention spans were given in the CR experiments and gold entities were given in the ED experiments. However, it is also necessary to explore developing and evaluating more practical systems in the full geoparsing setting, which requires systems to predict mentions, coreference clusters, and links from raw documents.

## Ethics Statement

As a potential risk associated with our dataset, a model trained on the dataset has the ability, to some extent, to identify locations mentioned in input texts and could be applied to link the content of individual posts containing private information with the mentioned locations. In addition, regardless of the purpose of use, the predicted locations may be inaccurate due to the limitations of the model’s performance or the discrepancy of domains, writing styles, and mentioned regions between our dataset and input texts.

Consistently with their intended use, we used existing language resources and tools to develop or evaluate NLP datasets or models under the specified license or terms of use. As for the dataset that we constructed, its intended use is for academic research purposes related to information science, similarly to that of the ATD. The text in our dataset is a subset of the original ATD data, and the original data does not contain any information about the travelogue authors.

The annotation work was performed by anno-

tators at a professional data annotation company. The payment amount to the company was based on the estimate submitted by the company. The actual annotators and the payment amount to each annotator was determined by the company. For mention, coreference, and link annotation, the annotation work were performed by five (four men and one woman), five (four men and one woman), and seven (five men and two women) annotators, respectively. The age range of the annotators is from their 20s to 50s. All of them are native Japanese speakers. Before commencing the annotation work to construct our dataset, we explained to the annotators that we or other researchers would use the annotated data for future research related to NLP.

## Acknowledgments

We would like to thank the anonymous reviewers and meta reviewers for their constructive comments. We used the Arukikata Travelogue Dataset to construct our dataset. This study was supported by JSPS KAKENHI Grant Number JP22H03648.

## References

- Arukikata. Co., Ltd. 2022. Arukikata travelogue dataset. Informatics Research Data Repository, National Institute of Informatics. <https://doi.org/10.32130/idr.18.1>.
- Amit Bagga and Breck Baldwin. 1998. Algorithms for scoring coreference chains. In *The First International Conference on Language Resources and Evaluation Workshop on Linguistics Coreference*, volume 1, pages 563–566.
- Timothy Baldwin, Marie Catherine de Marneffe, Bo Han, Young-Bum Kim, Alan Ritter, and Wei Xu. 2015. *Shared tasks of the 2015 workshop on noisy user-generated text: Twitter lexical normalization and named entity recognition*. In *Proceedings of the Workshop on Noisy User-generated Text*, pages 126–135, Beijing, China. Association for Computational Linguistics.
- Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. *ELECTRA: Pre-training text encoders as discriminators rather than generators*. In *8th International Conference on Learning Representations*, Addis Ababa, Ethiopia.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. *BERT: Pre-training of deep bidirectional transformers for language understanding*. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*,

- pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- George Doddington, Alexis Mitchell, Mark Przybocki, Lance Ramshaw, Stephanie Strassel, and Ralph Weischedel. 2004. [The automatic content extraction \(ACE\) program – tasks, data, and evaluation](#). In [Proceedings of the Fourth International Conference on Language Resources and Evaluation \(LREC'04\)](#), Lisbon, Portugal. European Language Resources Association (ELRA).
- Abbas Ghaddar and Phillippe Langlais. 2016. [WikiCoref: An English coreference-annotated corpus of Wikipedia articles](#). In [Proceedings of the Tenth International Conference on Language Resources and Evaluation \(LREC'16\)](#), pages 136–142, Portorož, Slovenia. European Language Resources Association (ELRA).
- Ralph Grishman and Beth Sundheim. 1996. [Message Understanding Conference- 6: A brief history](#). In [COLING 1996 Volume 1: The 16th International Conference on Computational Linguistics](#).
- Milan Gritta, Mohammad Taher Pilehvar, and Nigel Collier. 2018. [Which Melbourne? augmenting geocoding with maps](#). In [Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 1285–1296, Melbourne, Australia. Association for Computational Linguistics.
- Milan Gritta, Mohammad Taher Pilehvar, and Nigel Collier. 2020. [A pragmatic guide to geoparsing evaluation: Toponyms, named entity recognition and pragmatics](#). [Language Resources and Evaluation](#), 54:683–712.
- Masatsugu Hangyo, Daisuke Kawahara, and Sadao Kurohashi. 2014. [Building and analyzing a diverse document leads corpus annotated with semantic relations](#). [Journal of Natural Language Processing](#), 21(2):213–247.
- Chikara Hashimoto, Sadao Kurohashi, Daisuke Kawahara, Keiji Shinzato, and Masaaki Nagata. 2011. [Construction of a blog corpus with syntactic, anaphoric, and sentiment annotations](#). [Journal of Natural Language Processing](#), 18(2):175–201.
- Taiichi Hashimoto and Shun'ichi Nakamura. 2010. [Kakuchō koyū hyōgen tag tsuki corpus-no kōchiku—hakusho, shoseki, Yahoo! chiebukuro core data—\(Construction of an extended named entity-annotated corpus—white papers, books, Yahoo! chiebukuro core data\)](#). In [Proceedings of the 16th Annual Meeting of the Association for Natural Language Processing](#).
- Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenu, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. [Robust disambiguation of named entities in text](#). In [Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing](#), pages 782–792, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Xuke Hu, Zhiyong Zhou, Hao Li, Yingjie Hu, Fuqiang Gu, Jens Kersten, Hongchao Fan, and Friederike Klan. 2022. [Location reference recognition from texts: A survey and comparison](#). [Computing Research Repository](#), arXiv:2207.01683.
- Tomoya Iwakura, Kanako Komiya, and Ryuichi Tachibana. 2016. [Constructing a Japanese basic named entity corpus of various genres](#). In [Proceedings of the Sixth Named Entity Workshop](#), pages 41–46, Berlin, Germany. Association for Computational Linguistics.
- Davaajav Jargalsaikhan, Naoaki Okazaki, Koji Matsuda, and Kentaro Inui. 2016. [Building a corpus for Japanese wikification with fine-grained entity classes](#). In [Proceedings of the ACL 2016 Student Research Workshop](#), pages 138–144, Berlin, Germany. Association for Computational Linguistics.
- Mandar Joshi, Omer Levy, Luke Zettlemoyer, and Daniel Weld. 2019. [BERT for coreference resolution: Baselines and analysis](#). In [Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing \(EMNLP-IJCNLP\)](#), pages 5803–5808, Hong Kong, China. Association for Computational Linguistics.
- Ehsan Kamaloo and Davood Rafiei. 2018. [A coherent unsupervised model for toponym resolution](#). In [Proceedings of the 2018 World Wide Web Conference, WWW '18](#), page 1287–1296, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
- Daisuke Kawahara, Sadao Kurohashi, and Kōiti Hasida. 2002. [Construction of a Japanese relevance-tagged corpus](#). In [Proceedings of the Third International Conference on Language Resources and Evaluation \(LREC'02\)](#), Las Palmas, Canary Islands - Spain. European Language Resources Association (ELRA).
- Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. [End-to-end neural coreference resolution](#). In [Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing](#), pages 188–197, Copenhagen, Denmark. Association for Computational Linguistics.
- Jochen L Leidner. 2006. [An evaluation dataset for the toponym resolution task](#). [Computers, Environment and Urban Systems](#), 30(4):400–417.
- Michael D. Lieberman, Hanan Samet, and Jagan Sankaranarayanan. 2010. [Geotagging with local lexicons to build indexes for textually-specified spatial data](#). In [2010 IEEE 26th International Conference on Data Engineering](#), pages 201–212. IEEE.
- Xiao Ling and Daniel S Weld. 2012. [Fine-grained entity recognition](#). In [Proceedings of the 26th AAAI Conference on Artificial Intelligence](#).

- Xiaoqiang Luo. 2005. [On coreference resolution performance metrics](#). In [Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing](#), pages 25–32, Vancouver, British Columbia, Canada. Association for Computational Linguistics.
- Hiroshi Matsuda, Mai Omura, and Masayuki Asahara. 2019. [Tantan’i hinshi-no yōhō aimaisē kaiketsu-to ison kankē labeling-no dōji gakushū \(Simultaneous learning of usage disambiguation of parts-of-speech for short unit words and dependency relation labeling\)](#). [Proceedings of the 25th Annual Meeting of the Association for Natural Language Processing](#).
- Koji Matsuda, Akira Sasaki, Naoaki Okazaki, and Kentaro Inui. 2017. [Geographical entity annotated corpus of japanese microblogs](#). [Journal of Information Processing](#), 25:121–130.
- Paul McNamee, Hoa Trang Dang, Heather Simpson, Patrick Schone, and Stephanie M. Strassel. 2010. [An evaluation of technologies for knowledge base population](#). In [Proceedings of the Seventh International Conference on Language Resources and Evaluation \(LREC’10\)](#), Valletta, Malta. European Language Resources Association (ELRA).
- Yugo Murawaki and Shinsuke Mori. 2016. [Wikification for scriptio continua](#). In [Proceedings of the Tenth International Conference on Language Resources and Evaluation \(LREC’16\)](#), pages 1346–1351, Portorož, Slovenia. European Language Resources Association (ELRA).
- Hiroki Ouchi, Hiroyuki Shindo, Shoko Wakamiya, Yuki Matsuda, Naoya Inoue, Shohei Higashiyama, Satoshi Nakamura, and Taro Watanabe. 2023. [Arukikata travelogue dataset](#). arXiv:2305.11444.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. [CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in OntoNotes](#). In [Joint Conference on EMNLP and CoNLL - Shared Task](#), pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.
- Sameer Pradhan, Lance Ramshaw, Mitchell Marcus, Martha Palmer, Ralph Weischedel, and Nianwen Xue. 2011. [CoNLL-2011 shared task: Modeling unrestricted coreference in OntoNotes](#). In [Proceedings of the Fifteenth Conference on Computational Natural Language Learning: Shared Task](#), pages 1–27, Portland, Oregon, USA. Association for Computational Linguistics.
- Lev Ratinov, Dan Roth, Doug Downey, and Mike Anderson. 2011. [Local and global algorithms for disambiguation to Wikipedia](#). In [Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies](#), pages 1375–1384, Portland, Oregon, USA. Association for Computational Linguistics.
- Paul Rayson, Alex Reinhold, James Butler, Chris Donaldson, Ian Gregory, and Joanna Taylor. 2017. [A deeply annotated testbed for geographical text analysis: The corpus of lake district writing](#). In [Proceedings of the 1st ACM SIGSPATIAL Workshop on Geospatial Humanities, GeoHumanities’17](#), page 9–15, New York, NY, USA. Association for Computing Machinery.
- Ryokan Ri, Ikuya Yamada, and Yoshimasa Tsuruoka. 2022. [mLUKE: The power of entity representations in multilingual pretrained language models](#). In [Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 7316–7330, Dublin, Ireland. Association for Computational Linguistics.
- Giuseppe Rizzo, Marieke van Erp, Julien Plu, and Raphaël Troncy. 2016. [Making sense of microposts \(#Microposts2015\) named entity recognition and linking \(NEEL\) challenge](#). In [Proceedings of the 6th Workshop on ‘Making Sense of Microposts’](#), pages 50–59.
- Satoshi Sekine and Hitoshi Isahara. 2000. [IREX: IR & IE evaluation project in Japanese](#). In [Proceedings of the Second International Conference on Language Resources and Evaluation \(LREC’00\)](#), Athens, Greece. European Language Resources Association (ELRA).
- Satoshi Sekine, Kiyoshi Sudo, and Chikashi Nobata. 2002. [Extended named entity hierarchy](#). In [Proceedings of the Third International Conference on Language Resources and Evaluation \(LREC’02\)](#), Las Palmas, Canary Islands - Spain. European Language Resources Association (ELRA).
- Pontus Stenetorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou, and Jun’ichi Tsujii. 2012. [brat: a web-based tool for NLP-assisted text annotation](#). In [Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics](#), pages 102–107, Avignon, France. Association for Computational Linguistics.
- Kazuma Takaoka, Sorami Hisamoto, Noriko Kawahara, Miho Sakamoto, Yoshitaka Uchida, and Yuji Matsumoto. 2018. [Sudachi: a Japanese tokenizer for business](#). In [Proceedings of the Eleventh International Conference on Language Resources and Evaluation \(LREC 2018\)](#), Miyazaki, Japan. European Language Resources Association (ELRA).
- Erik F. Tjong Kim Sang. 2002. [Introduction to the CoNLL-2002 shared task: Language-independent named entity recognition](#). In [COLING-02: The 6th Conference on Natural Language Learning 2002 \(CoNLL-2002\)](#).
- Nobuhiro Ueda, Kazumasa Omura, Takashi Kodama, Hirokazu Kiyomaru, Yugo Murawaki, Daisuke Kawahara, and Sadao Kurohashi. 2023. [KWJA: A unified japanese analyzer based on foundation models](#).

- In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics: System Demonstrations, Toronto, Canada. Association for Computational Linguistics.
- Marc Vilain, John Burger, John Aberdeen, Dennis Connolly, and Lynette Hirschman. 1995. [A model-theoretic coreference scoring scheme](#). In Sixth Message Understanding Conference (MUC-6): Proceedings of a Conference Held in Columbia, Maryland, November 6-8, 1995.
- Jan Oliver Wallgrün, Morteza Karimzadeh, Alan M MacEachren, and Scott Pezanowski. 2018. [GeoCorpora: building a corpus to test and train microblog geoparsers](#). International Journal of Geographical Information Science, 32(1):1–29.
- Davy Weissenbacher, Arjun Magge, Karen O’Connor, Matthew Scotch, and Graciela Gonzalez-Hernandez. 2019. [SemEval-2019 task 12: Toponym resolution in scientific papers](#). In Proceedings of the 13th International Workshop on Semantic Evaluation, pages 907–916, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. [Scalable zero-shot entity linking with dense entity retrieval](#). In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6397–6407, Online. Association for Computational Linguistics.
- Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. [LUKE: Deep contextualized entity representations with entity-aware self-attention](#). In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6442–6454, Online. Association for Computational Linguistics.
- Ikuya Yamada, Koki Washio, Hiroyuki Shindo, and Yuji Matsumoto. 2022. [Global entity disambiguation with BERT](#). In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 3264–3271, Seattle, United States. Association for Computational Linguistics.
- Klim Zaporozets, Johannes Deleu, Yiwei Jiang, Thomas Demeester, and Chris Develder. 2022. [Towards consistent document-level entity linking: Joint models for entity linking and coreference resolution](#). In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 778–784, Dublin, Ireland. Association for Computational Linguistics.
- Zeyu Zhang and Steven Bethard. 2023. [Improving toponym resolution with better candidate generation, transformer-based reranking, and two-stage resolution](#). In Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (\*SEM 2023), pages 48–60, Toronto, Canada. Association for Computational Linguistics.

## A Licenses of Used Resources

We used some existing NLP software and language resources as described in the main sections. The licenses of the used resources are as follows. The Arukikata Travelogue Dataset is available via the Informatics Research Data Repository, National Institute of Informatics under specific terms of use.<sup>34</sup> brat, spaCy, GiNZA, KWJA, the pretrained Japanese ELECTRA model are available under the MIT License. SudachiPy and the pretrained mLUKE models are available under the Apache License 2.0. The pretrained Japanese BERT model is available under CC BY-SA 4.0. OpenStreetMap data files are available via Geofabrik<sup>35</sup> under the Open Database License 1.0.

## B Detailed Annotation Criteria

### B.1 Mention Span Annotation

The spans of geo-entity mentions are determined as follows. Generally, a noun phrase (NP) in which a head  $h$  is modified by a nominal modifier  $m$  is treated as a single mention (Table 11-a). An appositive compound of two nouns  $n_1$  and  $n_2$  is treated as a single mention (Table 11-b) unless there is some expression (e.g., *no*-particle “の”) or separator symbol (e.g., *tōten* “、”) inserted between them. A common name is treated as a single mention even if it is not a simple NP (Table 11-c). For an NP with an affix or affix-like noun  $a$  representing directions or relative positions, a cardinal direction prefix preceding a location name is included in the span (Table 11-d-1), but other affixes are excluded from the span (Table 11-d-2). There may be instances in which a modifier  $m$  represents a geo-entity, but its NP head  $h$  does not. In such cases, the modifier is treated as a single mention if the head is a verbal noun that means move, stay, or habitation (Table 11-e-1), but the NP is not treated as a mention if not (Table 11-e-2). In the case that a geo-entity name  $g$  is embedded in a non-geo-entity mention  $n$ , the inner geo-entity name is treated as a geo-entity mention if the external entity corresponds to an event held in the real world (Table 11-f). If the external entity corresponds to other types of entities, such as an organization or the title of a work, the inner geo-entity name is not treated as a geo-entity mention.

<sup>34</sup><https://www.nii.ac.jp/dsc/idr/arukikata/documents/arukikata-policy.html> (in Japanese)

<sup>35</sup><http://www.geofabrik.de/data/download.html>

(a)	<u>[山頂]<sub>m</sub> [駐車場]<sub>h</sub></u> <u>[parking area]<sub>h</sub> [on top of the mountain]<sub>m</sub></u>
(b)	<u>[駅ビル]<sub>n1</sub> [「ビエラ奈良」]<sub>n2</sub></u> <u>[station building]<sub>n1</sub> [Vierra Nara]<sub>n2</sub></u>
(c)	<u>天国への階段</u> <u>Stairway to Heaven</u>
(d-1)	<u>[東]<sub>a</sub> [東京]</u> <u>[East]<sub>a</sub> [Tokyo]</u>
(d-2)	<u>[北海道] [全域]<sub>a</sub></u> <u>[the whole area of]<sub>a</sub> [Hokkaido]</u>
(e-1)	<u>[京都]<sub>m</sub> [旅行]<sub>h</sub></u> <u>[Kyoto]<sub>m</sub> [Travel]<sub>h</sub></u>
(e-2)	<u>[三輪]<sub>m</sub> [そうめん]<sub>h</sub></u> <u>[Miwa]<sub>m</sub> [somen noodles]<sub>h</sub></u>
(f)	<u>[保津川]<sub>g</sub> 下り]<sub>n</sub></u> <u>[Hozugawa river]<sub>g</sub> boat tour]<sub>n</sub></u>

Table 11: Examples of mention spans.

### B.2 Coreference Annotation

We consider coreference and link annotation for TRANS mentions to be outside the scope of this study. This is because how to treat the identity of those mentions is not obvious, and OSM does not contain such type of entries. However, TRANS (-NAME) mentions would be helpful to identify the referents of other types of mentions that are not clearly written.

Following (or concurrently with) specificity tag annotation, relations are assigned to pairs of mentions that have not been labeled with either specificity tag.

**Specificity Tags** Specificity tags can be either GENERIC, SPEC\_AMB, or HIE\_AMB. GENERIC is assigned to a generic mention, as explained in §2.3. SPEC\_AMB (which means “specific but ambiguous”) is assigned to a mention that refers to a specific real-world location, but there is some ambiguity about the detailed area to which it refers, e.g., 海 *umi* in a sentence like “You can see a beautiful sea from this spot.” HIE\_AMB (which means “hierarchically ambiguous”) is assigned to an ambiguously described mention with multiple potential referents at both higher and lower-level locations, e.g., 奈良 in a sentence like “We are heading to Nara.” Annotators were instructed to annotate with coreference and link information, operating under the hypothesis that such mentions refer to the lowest-level location among candidate referents, e.g., not Nara

<sup>1</sup>世界遺産・<sup>2</sup>白川郷は素敵<sup>3</sup>ところでした。  
A <sup>1</sup>world heritage site, <sup>2</sup>Shirakawago was a nice <sup>3</sup>place.

Figure 2: Examples of attributive mentions.

Prefecture but Nara City.

**Coreference Relations** Coreference relations can be either the identical coreference relation COREF or the attributive coreference relation COREF\_ATTR. The coreference relation COREF is assigned to two mentions that both refer to the same real-world location, as explained in §2.3. The directed relation COREF\_ATTR is assigned to mention pairs in which one expresses the attribute of the other, either in appositive phrases or copular sentences. For example, a sentence in Figure 2 is annotated with COREF\_ATTR relations from mention 2 to mention 1 and from mention 2 to mention 3. This schema is similar to that in WikiCoref (Ghaddar and Langlais, 2016).

Notably, no coreference relations are assigned to mentions whose referents geographically overlap but are not identical; e.g., 首都高速道路 *shuto kōsoku dōro* ‘Metropolitan Expressway’ and 湾岸線 *wangansen* ‘Bayshore Route,’ which have a whole-part relation.

## C Detailed Dataset Statistics

### C.1 Mention Annotation

In the mention annotation step, 12,171 mentions were identified; they consist of 12,114 geo-entity and 57 non-geo-entity mentions (23 LOC\_ORG and 34 FAC\_ORG mentions). Table 3 shows the distribution of geo-entity mentions for entity type tags. The tag distribution represents some characteristics of travelogue documents of our dataset. First, the documents contain the largest number of facility mentions, which is even more than the number of location mentions. Second, the documents also contain the similar number of non-NAME (5,867)<sup>36</sup> to NAME mentions (6,247).

### C.2 Coreference Annotation

As a result of the coreference annotation step, 289 GENERIC mentions and 322 SPEC\_AMB mentions along with 923 TRANS mentions were excluded from the coreference relation annotation. Out of the remaining 10,580 mentions, 6,497 mentions

<sup>36</sup>Non-NAME mentions include \*-NOM, and DEICTIC mentions, in addition to all NOM mentions.

Size	1	2	3	4	5	6	≥7
#Cls	4,083	1,278	507	240	103	58	70
#Typ	1.0	1.5	2.0	2.3	2.6	2.8	3.3

Table 12: Number of geo-entity coreference clusters (#Cls) and the average number of member mention text types (#Typ) for each size.

	LOC	FAC	LINE	MIX	UNK
Set-A	819	1,823	327	29	133
Set-B	852	1,819	370	22	145
Total	1,671	3,642	697	51	278

Table 13: Tag distribution of geo-entities.

were annotated with one or more COREF and/or COREF\_ATTR relations among other mentions, of which 350 mention pairs were annotated with COREF\_ATTR relations. These mentions comprise coreference clusters with size  $\geq 2$ , and the remaining 4,083 mentions correspond to singletons. Table 12 shows the number of clusters and the average number of mention text types (distinct strings) among members<sup>37</sup> for each cluster size. This indicates that 35.6% (2,256/6,339) of coreference clusters have more than one member; that is, multiple mentions in a document often refer to the same referent.

In addition, we automatically assign an entity type tag to each coreference cluster, i.e., entity, from the tags of its member mentions.<sup>38</sup> Table 13 shows the tag distribution of entities, which is similar to the tag distribution of mentions shown in Table 3.

### C.3 Link Annotation

As shown in Table 14, in the link annotation step for Set-B, 79.5% (2,551) and 64.2% (2,059) of 3,208 entities have been annotated with any URLs and OSM entry URLs, respectively, including entities annotated with PART\_OF tags. For “HasName” entities in which at least one member mention is labeled as NAME, any URLs and OSM entry URLs

<sup>37</sup>For example, for clusters  $C_1 = \{\text{“Nara Station”, “Nara Sta.”, “Nara”}\}$  and  $C_2 = \{\text{“Kyoto Pref.”, “Kyoto”, “Kyoto”}\}$ , the numbers of mention text types are three and two, respectively, and their average is 2.5.

<sup>38</sup>(a) LOC, FAC, or LINE is assigned to an entity that the members’ tags include only one of the three types and optionally include DEICTIC or LOC\_OR\_FAC\_NOM (for LOC and FAC). (b) UNK is assigned to an entity that all members’ tags are DEICTIC or LOC\_OR\_FAC\_NOM. (c) MIX is assigned to an entity that the members’ tags include two or three of LOC, FAC, and LINE.

	All	HasRef	HasOSMRef
HasName	2,001	1,942	1,574
HasNoName	1,207	609	485
Total	3,208	2,551	2,059

Table 14: Numbers of Set-B entities that have names and/or references in the PART\_OF-*inclusive* setting where entities assigned with PART\_OF (along with URLs) are counted as instances of “Has(OSM)Ref.”

	All	HasRef	HasOSMRef
HasName	2,001	1,861	1,514
HasNoName	1,207	298	221
Total	3,208	2,159	1,735

Table 15: Numbers of Set-B entities that have names and/or referents in the PART\_OF-*exclusive* setting where entities assigned with PART\_OF (along with URLs) are NOT counted as instances of “Has(OSM)Ref.”

are assigned to 97.1% (1,942/2,001) and 78.7% (1,574/2,001) of them, respectively. This indicates that the real-world referents can be easily identified for most of the entities explicitly written with their names. For the remaining “HasNoName” entities, any URLs and OSM entry URLs are assigned to 50.5% (609/1,207) and 40.2% (485/1,207) of them, respectively. This suggests that identifying the referents from unclearly written mentions and context is difficult even for humans.

As shown in Table 15, the percentages of referent-identified entities decrease in the setting where entities assigned with PART\_OF are excluded. The result indicates the reasonable coverage of OSM for various types of locations in Japan. Overall, entities assigned with OSM entries account for 75.7% (1,514/2,001) of “HasName” entities. For details on each entity type tag of LOC, FAC, LINE, and the others, entities assigned with OSM entries account for 79.3% (811/1,096), 74.0% (544/686), 72.7% (144/198), and 71.4% (15/21) of “HasName” entities with the specified tag, respectively.

#### C.4 Geographical Distribution of Linked Entities

As we expected, most of the mentions in our (Set-B) dataset refer to locations in Japan, except for 34 mentions that refer to overseas locations. Figure 3 shows the geographical distribution of linked entities in our dataset, namely, the number of entities located in each prefecture among entities annotated with OSM entry URLs. For example, there are 45

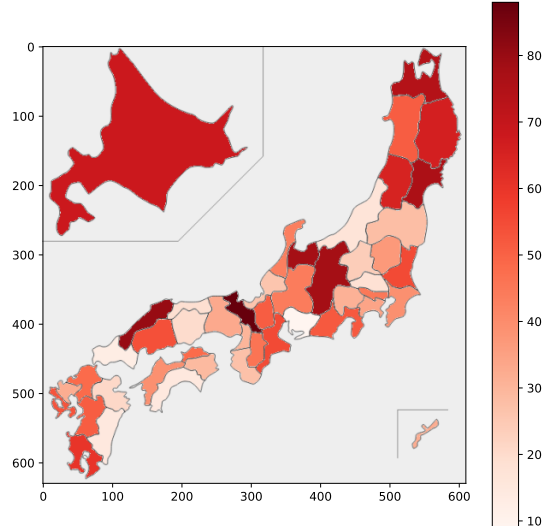


Figure 3: Numbers of linked entities located in each prefecture. Deeper red indicates the larger number. The units of the numerical values on the vertical and horizontal axes of the map are kilo-miles.

linked entities to which the coordinates of OSM entries are linked within the area of Tokyo Prefecture in all annotated travelogue documents, and thus the count of Tokyo Prefecture is 45. The minimum, maximum, and average numbers of entity counts in all 47 prefectures are 9 (Aichi), 88 (Kyoto), and 42.8, respectively.

Figure 4 shows actual examples of mentions with *geographic continuity*; that is, mentions that refer to nearby locations in the real world tend to appear near to one another within a document (§1). The example text in document ID 00019 describes five geo-entities located nearby in the real world. Table 16 further shows actual sentences, being the first five sentences that include at least one annotated mention, extracted from three documents with the smallest ID values in the development set. Including the examples depicted in Figure 4, we can observe mentions with geographic continuity.

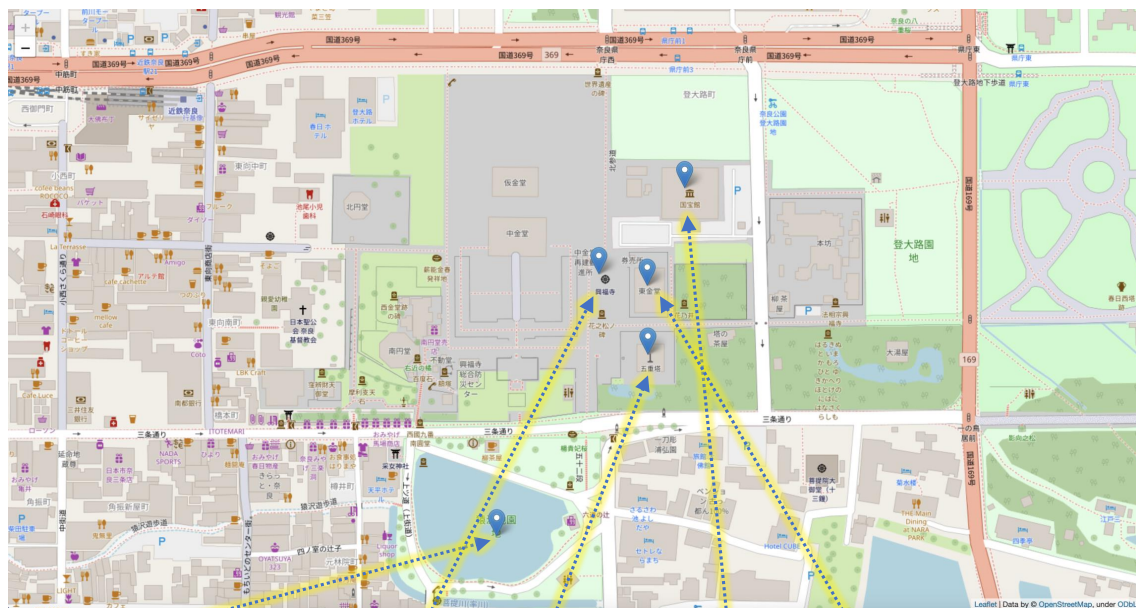
## D Details on Experimental Settings

### D.1 Evaluation Scripts

We used our code that calculates general precision, recall, and F1 score in the mention recognition and entity disambiguation experiments. We used our code that calculates the MUC, B<sup>3</sup>, and CEAF<sub>e</sub> scores in the manner equivalent to an existing evaluation tool<sup>39</sup> in the coreference resolution experiments.

<sup>39</sup><https://github.com/ns-moosavi/coval/blob/master/coval/eval/evaluator.py>





写真は猿沢池からも見える興福寺の五重塔です。国宝館と東金堂に行く場合は、...

Sarusawaikae Pond Kohfukuji Temple Five-storied Pagoda National Treasure Hall Eastern Golden Hall

Figure 4: Example of actual text, including mentions with *geographic continuity* in document ID 00019 (sentence IDs 009–010, the English translation is given in Table 16). The map depicts part of the Nara Park area, a popular sightseeing area in Nara City, Japan.

SentID	Text	English Translation
001	奈良公園 <sup>FAC-NAME way/456314269</sup> のアイドル「しか」で〜す。	There are deers, the idols in the <u>Nara Park</u> .
004	奈良 <sup>LOC-NAME (HIE_AMB), relation/3227787</sup> の有名スポット <sup>LOC_OR_FAC-NOM way/456314269</sup> ですよ!	It's a <u>famous spot</u> in <u>Nara</u> , right?
005	大仏 <sup>FAC-NOM way/43558119</sup> 様はとっても大きかったなあ〜	The <u>Great Buddha</u> was really huge.
009	写真は猿沢池 <sup>LOC-NAME way/59465653</sup> から見える興福寺 <sup>FAC-NAME way/1134439456</sup> の五重塔 <sup>FAC-NOM way/98093571</sup> です。	It's a photo of the <u>five-storied pagoda</u> at <u>Kofukuji Temple</u> visible from <u>Sarusawaikae Pond</u> .
010	国宝館 <sup>FAC-NAME way/98093576</sup> と東金堂 <sup>FAC-NAME way/98093572</sup> に行く場合は...	If you go to the <u>National Treasure Museum</u> and <u>Eastern Golden Hall</u> . . .
001	...瀬戸大橋 <sup>LINE-NAME relation/10375178</sup> をようやく通ります。	I'm finally crossing <u>Seto Ohashi Bridge</u> . . .
002	四国 <sup>LOC-NAME relation/2906044</sup> にも初上陸。	I just landed in <u>Shikoku</u> for the first time, too.
009-01	こんぴら猫 <sup>FAC-NAME general_page</sup> 。	<u>Kompira Dog</u> .
010	みやげ屋 <sup>FAC-NOM (GENERIC)</sup> が連なる参道 <sup>LINE-NOM (SPEC_AMB)</sup> もまた、...	The approach lined with <u>souvenir shops</u> is. . .
012	3~4年前に浪速餃子スタジアム <sup>FAC-NAME general_page</sup> で...	About 3–4 years ago at the <u>Naniwa Gyoza Stadium</u> . . .
001-01	二社一寺は日光山内 <sup>LOC-NAME Wikidata:Q1063133</sup> ともいいますが...	The “two shrines and one temple” are also called <u>Nikko San'nai</u> . . .
002	まずは、輪王寺 <sup>FAC-NAME way/699236460</sup> の金堂 <sup>FAC-NOM way/388017115</sup> ・三仏堂 <sup>FAC-NAME way/388017115</sup> 。	First, the <u>main holl</u> , <u>Sambutsudo</u> at <u>Rin'noji Temple</u> .
003-02	三仏堂 <sup>FAC-NAME way/388017115</sup> では干支のお守りも購入できます。	At <u>Sambutsudo</u> , you can purchase zodiac charms.
004	三仏堂 <sup>FAC-NAME way/388017115</sup> の裏手にある護摩堂 <sup>FAC-NAME way/388017145</sup> で...	At <u>Gomado</u> located behind <u>Sambutsudo</u> . . .
005-01	次は徳川家康公を祭る日光東照宮 <sup>FAC-NAME way/388017091</sup> です。	Next is <u>Nikko Toshogu Shrine</u> , where Tokugawa Ieyasu is enshrined.

Table 16: Examples of actual sentences and annotated mention (blue underline and superscript) and coreference/link information (subscript). The displayed sentences are the first five sentences that include at least one annotated mention in each document: ID 00019 (top), 01158 (middle), and 03088 (bottom).

## D.2 Entity Type Conversion Rules

**IREX** We used the following rules to convert the IREX tags to our entity type tags. (1) Each output mention with the LOCATION tag was converted into three mention instances with the same span and with one of LOC\_NAME, FAC\_NAME, and LINE\_NAME tags. (2) ARTIFACT was converted into TRANS\_NAME.

**ENE** We used the following rules to convert the ENE tags (version 7.1.0),<sup>40</sup> which GiNZA adopted, to our entity type tags. (1) The Location subtype tags except for the Astral\_Body subtype tags, the Address subtype tags and River were converted to LOC\_NAME. (2) The Facility subtype tags except for the Line subtype tags were converted to FAC\_NAME. (3) River and the Line subtype tags were converted to LINE\_NAME. (4) Service and the Vehicle subtype tags were converted to TRANS\_NAME.

## D.3 Details of Database Preprocessing

The original OSM data contains a huge number of entries, and multiple entries can refer to almost the same real-world locations; for example, we found 72 entries named 東京 ‘Tokyo,’ including four railway stations, two railway station platforms, one ferry terminal, 30 train stop positions, and 27 footway sections, 8 flights of steps on footways, some of which can be equated with each other. For practical evaluation of ED systems, different entries that can be treated as equivalent should be grouped together, and such groups should be considered as linking units rather than individual entries.

Therefore, we reorganized the raw OSM data as follows. (1) We downloaded an OSM data file consisting of Japanese domestic location entries. (2) We extracted 2.8M entries with “name” attributes from the total of 2.6B entries. (3) We added 14 out of 16 entries without name attributes that were assigned to domestic geo-entities in the Set-B data, but were not contained in the extracted entries (the remaining two entries had been deleted from OSM). This resulted in DB coverage of 99.86% for the Set-B entities annotated with OSM entry URLs. (4) We then generated a *group ID string* from the original name attribute for each entry by concatenating part of the address and notable OSM tags, such as the branch name and amenity type. (5) Finally, we grouped entries with the same group ID into the

<sup>40</sup>[https://nlp.cs.nyu.edu/ene/version7\\_1\\_0Beng.html](https://nlp.cs.nyu.edu/ene/version7_1_0Beng.html)

same entry group. This series of processes resulted in 1.8M entry groups.<sup>41</sup>

## D.4 Settings of spaCy-MR

For building our custom MR model with spaCy, namely, spaCy-MR, we used almost the same settings as GiNZA,<sup>42</sup> including model architecture and hyperparameters, tokenizer, and training settings except that we disabled unnecessary pipelines other than “transformer” and “ner.” We reported the result of a single run of spaCy-MR in §5.3 and Appendix E.

## D.5 Implementation and Settings of mLUKE-MR/CR

We reported the results of single runs of mLUKE-MR and mLUKE-CR in §5.3 and Appendix E.

**Mention Recognition** Following Yamada et al. (2020), we tackle the task by enumerating and classifying all possible spans in each sentence. The representation of each candidate span is a concatenation of the word representations of the first and last tokens of the span, and the entity representation corresponding to the span, all of which are computed by the LUKE Transformer model. We employ a linear classifier to classify spans into the target entity types or *non-entity* type. We restrict candidate spans to the positions where their first and last tokens correspond to word boundaries (obtained using Sudachi Mode B), and exclude spans longer than 16 tokens.<sup>43</sup> Following Devlin et al. (2019) and Yamada et al. (2020), we prepend/append the surrounding tokens to a target sentence (up to 512 tokens in total) to give sufficient contextual information to the model.

**Coreference Resolution** Following Lee et al. (2017), we solve the task as antecedent identification for each mention. We follow the architecture proposed by Joshi et al. (2019) except that we do not use a unary score for each mention or coarse-to-fine inference because gold mentions are given in our setting.<sup>44</sup> The representation of each mention

<sup>41</sup>We will publish the preprocessed database at <https://github.com/naist-nlp/atd-mcl-baselines>.

<sup>42</sup>[https://github.com/megagonlabs/ginza/blob/develop/config/ja\\_ginza\\_electra.cfg](https://github.com/megagonlabs/ginza/blob/develop/config/ja_ginza_electra.cfg)

<sup>43</sup>We also enforce word boundaries on the mLUKE tokenizer because (word-level) mention annotation in the ATD-MCL does not align with unigram segmentation used in the tokenizer.

<sup>44</sup>We also omit discrete features based on the metadata available only in some datasets.

Task	Name	Value
MR	Learning rate	1e-5
	Batch size	8
	Training epochs	10
CR	Learning rate	5e-5
	Batch size	4
	Training epochs	20
Common	Learning rate decay	linear
	Warmup ratio	0.06
	Dropout	0.1
	Weight decay	0.01
	Gradient clipping	none
	Adam $\beta_1$	0.9
	Adam $\beta_2$	0.98
Adam $\epsilon$	1e-6	

Table 17: Hyperparameter values used in the mLUKE-MR/CR experiments.

is computed in the same way as the MR model. The model is trained by optimizing the marginal log-likelihood of the possibly correct antecedents including a dummy antecedent, which indicates no antecedents associated with a target mention. Because CR in the ATD-MCL is a document-level task and documents in the dataset are too long to be processed by a Transformer-based model for computational reasons, we independently feed each sentence in a document to the LUKE model, but optimization/prediction is made in each document.

**Hyperparameters** The hyperparameter values used in the experiments using mLUKE-MR/CR are listed in Table 17. Because our computational resources were limited, we did not conduct hyperparameter tuning except learning rate. We chose the best setting of learning rate and the number of training epochs from the search space of  $\{1e-5, 2e-5, 3e-5, 4e-5, 5e-5\}$  and  $\{5, 10, 20\}$ , respectively. We specifically selected batch size for each task, but we followed Yamada et al. (2020) for the other hyperparameters.

## D.6 Size of Used Models

Table 18 shows the numbers of model parameters in the systems that we used in the experiments. For KWJA, we report the number of parameters (112M) in the pretrained model<sup>45</sup> used in the KWJA base model (while the actual number of parameters in the whole model would be larger).

<sup>45</sup><https://huggingface.co/ku-nlp/deberta-v2-base-japanese>

Tasks	System	#Params
MR	mLUKE-MR	561M
MR	spaCy-MR	109M
MR	GiNZA (ja_ginza_electra)	110M
MR, CR	KWJA (base)	112M+
CR	mLUKE-CR	877M
ED	BERT-ED	111M

Table 18: Numbers of model parameters in evaluated systems.

## D.7 Computational Budget for Finetuning

In our experiments, mLUKE-MR was finetuned for 130 minutes (10 epochs) using four NVIDIA Tesla V100 GPUs with 16GB memory. mLUKE-CR was finetuned for 15 minutes (20 epochs) using four NVIDIA A100 Tensor Core GPUs with 40GB memory. spaCy-MR was finetuned for 17.4 hours (20000 steps) using a four-core Intel Xeon Gold 6150 CPU (32 cores total).

## E Detailed Experimental Results on Mention Recognition

Table 19 shows detailed performance of mention recognition systems. The finetuned systems spaCy-MR and mLUKE-MR achieved F1 scores higher than 0.6 and 0.7, respectively, for all tags except for TRANS\_NAME and FAC\_ORG.

Tag	#	KWJA			GiNZA			spaCy-MR <sup>o</sup>			mLUKE-MR <sup>o</sup>		
		P	R	F1	P	R	F1	P	R	F1	P	R	F1
Overall	4,958	.279	.352	.311	.574	.277	.374	.752	.732	.742	<b>.813</b>	<b>.817</b>	<b>.815</b>
NAME	2,509	.279	.695	.398	.574	.548	.560	.733	.719	.726	.828	.813	.821
NOM	2,203	0	0	0	0	0	0	.790	.753	.771	.826	.818	.822
ORG	24	0	0	0	0	0	0	.353	.250	.293	.833	.417	.556
LOC_NAME	881	.378	.857	.525	.617	.717	.664	.727	.822	.771	.830	.863	.846
FAC_NAME	1,285	.409	.635	.497	.589	.504	.543	.770	.689	.727	.843	.807	.825
LINE_NAME	195	.061	.621	.110	.425	.405	.415	.673	.677	.675	.804	.800	.802
TRANS_NAME	148	.193	.358	.251	.176	.101	.129	.525	.432	.474	.707	.588	.642
LOC_NOM	349	0	0	0	0	0	0	.739	.691	.714	.748	.808	.777
FAC_NOM	1,135	0	0	0	0	0	0	.816	.757	.785	.855	.819	.837
LINE_NOM	236	0	0	0	0	0	0	.749	.822	.784	.865	.818	.841
TRANS_NOM	334	0	0	0	0	0	0	.840	.817	.829	.830	.877	.853
LOC_OR_FAC_NOM	149	0	0	0	0	0	0	.676	.617	.646	.731	.711	.721
DEICTIC	222	0	0	0	0	0	0	.645	.721	.681	.616	.896	.730
LOC_ORG	11	0	0	0	0	0	0	.750	.545	.632	.900	.818	.857
FAC_ORG	13	0	0	0	0	0	0	0	0	0	.500	.077	.133

Table 19: System performance for mention recognition. “<sup>o</sup>” indicates the models finetuned on the ATD-MCL training set. “#” indicates the number of mentions for each tag in the test set.