

LREC-COLING 2024

Constructing Indonesian-English Travelogue Dataset

Eunike Andriani Kardinata¹, Hiroki Ouchi^{1,2}, Taro Watanabe¹

¹Nara Institute of Science and Technology ²RIKEN {eunike.kardinata.ef9, hiroki.ouchi, taro}@is.naist.jp





- 1. Introduction
- 2. Dataset Construction
- 3. Dataset Evaluation
- 4. Conclusions
- 5. Future Works

CONTENT



Low-Resource Languages (LRLs)

- Languages spoken in the world with **less linguistic resources** for technologies. (Cieri et al. 2016)
- Indonesian language is lacking in labelled data collection, but having a growing presence in the digital world (cluster 3).
- There is an increasing effort to develop Indonesian **datasets** and **language models**. (Willie et al. 2020, Ariesandy et al. 2020, Aji et al. 2022, Winata et al. 2023)



Fig. 1 - Language Resource Distribution (Joshi et al. 2020)



<u>1. Introduction</u> Challenges Faced by LRLs





(Magueresse et al. 2020)

5

Importance of Geographic Data in Indonesia



- With the rise in travel and tourism in Indonesia, more information from both locals and common tourists would be needed.
- Valuable information are often contained in texts written by common people.
 - Collection of geographic data through the use of sensors and its reports in the form of statistics by certain organizations are often **limited**.







Processing Geographic Data in Texts

Geoparsing consists of two main tasks (Gritta et al. 2020):

- Toponym Extraction (Geotagging)
 Similar to named entity recognition (NER), but more focused on reference (mention) of location (toponym) in the text.
- Toponym Resolution (Geocoding) Regarded as entity linking (EL), i.e., disambiguate location mentions in the text using available knowledge bases.



Processing Geographic Data in Texts

Sample Text

Last week, we went to visit Heijo Palace. We really enjoyed our time there, it was such a big area to explore. Before that, we also bought some drinks at Family Mart near Yamato-Saidaiji station. We thought of visiting Nara Park too, but it was too late. We decided to visit the park the day after.

Processing Geographic Data in Texts

Sample Text

Last week, we went to visit Heijo Palace.

We really enjoyed our time there, it was

such a big area to explore. Before that,

we also bought some drinks at Family Mart

near Yamato-Saidaiji station. We thought of

visiting Nara Park too, but it was too late.

We decided to visit the park the day after.

Tasks:

1. Named Entity Recognition

9



Processing Geographic Data in Texts

Sample Text

Last week, we went to visit Heijo Palace. We really enjoyed our time there, it was such a big area to explore. Before that, we also bought some drinks at Family Mart near Yamato-Saidaiji station. We thought of visiting Nara Park too, but it was too late. We decided to visit the park the day after.

Tasks:

- 1. Named Entity Recognition
- 2. Nominal Expressions Recognition





Processing Geographic Data in Texts

Sample Text

Last week, we went to visit Heijo Palace. We really enjoyed our time there, it was such a big area to explore. Before that, we also bought some drinks at Family Mart near Yamato-Saidaiji station. We thought of visiting Nara Park too, but it was too late. We decided to visit the park the day after. Tasks:

- 1. Named Entity Recognition
- 2. Coreference Resolution (including Nominal Expressions Recognition)





Processing Geographic Data in Texts

Sample Text

Last week, we went to visit Heijo Palace. We really enjoyed our time there, it was such a big area to explore. Before that, we also bought some drinks at Family Mart near Yamato-Saidaiji station. We thought of visiting Nara Park too, but it was too late. We decided to visit the park the day after.

Tasks:

- 1. Named Entity Recognition
- 2. Coreference Resolution
- 3. Entity Disambiguation/Linking



Fig. 2 - Sample Capture of Google Maps



Challenges in Geoparsing



• Metonymy Resolution (Gritta et al. 2018)

A toponym word is used to substitute for something else For example, "<u>Japan</u> wins the 2023 World Baseball Classic"

Location Inference based on the surrounding context (Farzana and Hecking. 2023)
 A location is being described instead of explicitly mentioned in the text
 For example, "Famous park in Nara prefecture with many deers" [Nara Park]





State of Indonesian Language Resources

Name	Year	Domain
IndQNER	2022	religion
IndoNLU NERGrit	2020	general
NERGrit	2020	general
NERP (IndoNLU Split)	2018	news
NER UI (IndoLEM split)	2017	general
Singgalang	2017	wiki
WikiAnn (multilingual)	2017	wiki
NER UGM (IndoLEM split)	2014	news

Table 1: NER Datasets in NusaCrowd
(Cahyawijaya et al. 2023)

Dataset Name	Year	Language	Size	Domain
ATD-MCL	2023	ia	12K	Travelogue
Event Geoparsing	2020	id	1.1K	News
GeoWebNews	2020	en	2.4K	News
SemEval-2019 T12	2019	en	8.4K	Science
GeoCorpora	2018	en	3.1K	Microblog
TR-News	2018	en	1.3K	News
GeoVirus	2018	en	2.2K	News
CLDW	2017	en	3.7K	Historical
LRE Corpus	2017	ja	1.0K	Microblog

 Table 2: Geoparsing Datasets (Size in Mentions)

Focus on Geographic

Indonesian Travelogue







2. Dataset Construction

2. Dataset Construction Scope of the Dataset



Referring to the recent Japanese travelogue dataset (Higashiyama et al. 2023, Ouchi et al. 2023)



Note:

Mention \rightarrow both named entities and **nominal expressions**

2. Dataset Construction Scope of the Dataset





Note:

Mention \rightarrow both named entities and **nominal expressions**





→ Obtained 88 articles (44 ID + 44 EN). (see detailed statistics here)

 \rightarrow Only articles with very **similar** content are included.

2. Dataset Construction

How similar are the ID and EN articles?



Fig. 3 - Sample of Annotated Sentences in Indonesian and English







MAIST KON

 \rightarrow Using brat rapid annotation tool with some tags as defined in JP dataset.



How the Annotation was Conducted

Tags included in the annotation:

- 1. LOC_NAME : naturally existing locations, e.g., country, mountain, lake
- 2. **FAC_NAME** : man-made structures or area, e.g., park, building, station
- 3. **TRANS_NAME :** transportation modes or vehicles, e.g., bus, train, ship
- 4. **LINE_NAME** : roads and waterways, e.g., street, river, route
- 5. And another four for the nominal expression of each of the above. LOC_NOM, FAC_NOM, TRANS_NOM, LINE_NOM

2. Dataset Construction



How the Annotation was Conducted

Overview:

- 1. Based on the guideline, did **trial** annotation for 15 articles.
- 2. **Revised** the guideline based on what was discovered in the trial round.
- 3. Continued the annotation and asked for other **independent** annotators.
- 4. Calculated inter-annotator agreement and observe the results.
- 5. Conducted some **experiments** to evaluate the dataset.





Inter-Annotator Agreement (Exact Match)

		F1	Ann1	Ann2	Both	
	Named	0.839	328	309	294	
id	Nominal	0.757	225	191	165	
	All	0.792	553	500	459	
2	Named	0.828	268	256	224	
en	Nominal	0.719	187	195	127	
	All	0.766	455	451	351	

Table 3: Inter-Annotator Agreement (Exact Match)

Labol	Indonesian			English				
Laber	F1	Ann1	Ann2	Both	F1	Ann1	Ann2	Both
LOC_NAME	0.949	207	215	203	0.864	160	174	149
FAC_NAME	0.817	106	85	82	0.788	98	73	68
TRANS_NAME	-	-	-	-	-	-	-	-
LINE_NAME	0.750	15	9	9	0.833	10	9	7
LOC_NOM	0.844	88	85	79	0.551	72	82	48
FAC_NOM	0.767	86	74	62	0.633	81	76	49
TRANS_NOM	0.805	25	13	13	0.900	13	12	12
LINE_NOM	0.613	26	19	11	0.792	21	25	18

Table 4: Inter-Annotator Agreement by Labels (Exact Match)

3. Dataset Evaluation Different Spans in Annotation







Fig. 4 - Sample of Annotated Sentences with Different Spans in English

3. Dataset Evaluation Different Spans in Annotation



ANNOTATOR 1 LOC_NOM FAC NAME FAC_NAME FAC_NAME FAC_NAME FAC_NAME Water spring can be found in post 1, 3, 5, 8. 10. Amazing view LOC_NAME LOC_NOM LOC_NAME LOC_NOM from the top of Mount Bawakaraeng, with range of hills in Sulawesi view. and also city view.





Fig. 4 - Sample of Annotated Sentences with Different Spans in English

3. Dataset Evaluation Dataset Statistics



Number of	Indor	nesian	English				
Number of	Total	Average	Total	Average			
Sentences	1,391	31	1,914	43			
Words	47,415	1,077	47.902	1,088			
Named	3,937	89	2,756	62			
Nominal	2,062	46	2,243	50			
Named (Unique)	1,156	26	1,053	23			
Nominal (Unique)	430	9	760	17			

Table 5: Dataset Statistics for Indonesian and English

~11K mentions

Dataset Name	Year	Language	Size	Domain
ATD-MCL	2023	ja	12K	Travelogue
Event Geoparsing	2020	id	1.1K	News
GeoWebNews	2020	en	2.4K	News
SemEval-2019 T12	2019	en	8.4K	Science
GeoCorpora	2018	en	3.1K	Microblog
TR-News	2018	en	1.3K	News
GeoVirus	2018	en	2.2K	News
CLDW	2017	en	3.7K	Historical
LRE Corpus	2017	ja	1.0K	Microblog

Table 2: Geoparsing Datasets (Size in Mentions)

Experiment and Results

- Trained classifiers to recognise named entities and nominal expressions.
- Used spaCy NER with transformers for Indonesian and English.
- For each language, split 44 articles into the train, validation, and test set in the ratio of 8:1:1, giving 35, 4, and 5 articles.

		Precision	Recall	F1
	Named	0.881	0.841	0.853
id	Nominal	0.910	0.914	0.912
	Overall	0.923	0.938	0.931
2	Named	0.877	0.859	0.866
en	Nominal	0.902	0.910	0.906
~	Overall	0.922	0.922	0.922

Table 6: Results of Experiments (Macro Ave.)



Comparison with SpaCy

- Simple comparison using different texts:
 - 1. Travelogue by the same author
 - 2. Travelogue by a different author
 - 3. Wikipedia article
 - 4. News article
- Labels related to geographic entities in SpaCy:
 - 1. FAC : Buildings, airports, highways, bridges, etc.
 - 2. **GPE** : Countries, cities, states.
 - 3. **LOC** : Non-GPE locations, mountain ranges, bodies of water.



Comparison with SpaCy on Unseen Sample



Fig. 5 - Our Classifier on Travelogue Article by the Same Author



2 misses

story.

Fig. 6 - SpaCy Classifier on Travelogue Article by the Same Author



3. Dataset Evaluation Comparison with SpaCy on Another Author



Fig. 7 - Our Classifier on Travelogue Article by a Different Author

Fortresses and defensive walls pepper the island, as do churches and traditional villages. Being an island there are plenty of coves and little beaches to discover. The diving is superb here, with underwater caves, and plenty of wrecks to discover. Go souvenir hunting around Valletta GPE , hitch a ride with a donkey on Gozo GPE , eat delicious sea-food in a small fishing village, or admire the sunset from one of the many cliffs.

Fig. 8 - SpaCy Classifier on Travelogue Article by a Different Author



1 miss

Comparison with SpaCy on Wikipedia Article



Fig. 10 - SpaCy Classifier on Wikipedia Article



3. Dataset Evaluation Comparison with SpaCy on News Article



1 miss

2 misses

Fig. 12 - Our Classifier on News Article



<u>3. Dataset Evaluation</u> Further Analysis

- Types of **mistakes** occurring in our classifier:
 - 1. Not recognising new mentions
 - 2. Not recognising existing mentions
 - 3. Misclassification of mentions
 - 4. Incomplete span of mentions
- Further investigation is necessary to figure out the characteristics and tendencies of our classifier, i.e., whether the problem mainly lies in the dataset or the classifier itself.







4. Conclusions



- We have constructed an Indonesian-English travelogue dataset to with a new annotation scheme that includes nominal expressions.
- We have conducted **experiments** and observe the **potential use** of our dataset for research.
- Further expansion of scope is necessary.

SUMMARY



5. Future Works

5. Future Works

Expansion of the Dataset



Within the same scope of use:

- Add more travelogue articles from diverse authors (writing styles).
- + Cover more languages \rightarrow foreign or local languages in Indonesia.
- Conduct more experiments, e.g., comparing with NER datasets, etc.
- Improve the guideline or the tools to handle multi-spans of mention.

Beyond the current scope of use:

- Include coreference resolution and entity linking (geocoding).
- Extend the use of the dataset, e.g., including human behaviours, etc.





End of Presentation

We would like to thank the anonymous reviewers and meta reviewers for their constructive comments. This study was supported by JSPS KAKENHI Grant Number JP22H03648.

Thank you for your attention.