



#MEDINF023

Towards Structuring Clinical Texts: Joint Entity and Relation Extraction from Japanese Case Report Corpus

### Yoshimasa KAWAZOE<sup>1</sup>

Associate Professor <sup>I</sup> Graduate School of Medicine, The University of Tokyo

Co-authors: Daisaku SHIBATA<sup>1</sup>, Emiko SHINOHARA<sup>1</sup>, Kiminori SHIMAMOTO<sup>1</sup>







### Background: Challenges with clinical text usage

- Clinical text contains valuable information for research purposes
- it is still challenging to extract specific or relevant information by NLP

Spelling inconsistency / va	riants ———					
血圧110-120 mmHg	<b>収縮期</b> 120mmHg	転倒した ベッ	ドから落ちた			
Blood pressure	Systolic BP	Fall	Fall from bed			
「Factuality(positive or Negative) 関節炎を認めた 筋力低下や錐体路徴候を伴わない 胸水あるも腹水なく Arthritis present Muscles weakness and pyramidal signs absent Pleural effusion present, ascites absent						
Factuality (General knowledge It did not occur in the patient )						
術後に出血が予想され	る 網膜症がを	今任しやすい				

Bleeding is expected after surgery.

Retinopathy tends to occur as a complication

@TheInstituteDH





### Background: Information retrieval requires a text corpus

- Datasets (text + annotation) for developing NLP technology to extract information from text.
- Entity tags represent the type of entity, while relationship tags represent the type of relationship between entities.





#### 8 – 12 JULY 2023 | SYDNEY, AUSTRALIA

# Background: iCorpus<sup>1</sup>

- A corpus consists cases reports of rare and intractable disease.
- Based on annotation criteria of versatility, comprehensiveness, and consistency.
- Formulated as NLP tasks (Named Entity Recognition and Relation Extraction)
  - 183 Documents, 56 Entity types, 35 Relationships
- Publicly available for research use

https://ai-health.m.u-tokyo.ac.jp/home/research/corpus

@TheInstituteDH #MEDINF023



Shinohara E et al. J Biomed Inform. 2022 Oct;134:104200.







## Objective

- Showing the baseline performance of Named Entity Recognition (NER) and Relation Extraction (RE) using iCorpus.
- 2. Also, showing the difference in performance between two types of BERT<sup>1</sup>
  - Clinical-BERT: Pretrained on clinical text
  - Wikipedia-BERT: Pretrained on Wikipedia (JP) text

1. Jacob Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv:1810.04805

@TheInstituteDH #MEDINF023



## Method: Compare the two types of BERTs

BERT is pre-trained on a large corpus of text, then fine-tuned with just one additional layer for specific tasks, such as NER and RE.

### Clinical-BERT<sup>1</sup>

- BERT-base: 12 (L), 12 (A), 768 (E)
- Clinical text (120 million lines)
- 25,000 vocabularies (BPE)

#MEDINF023

• Domain: Clinical

@TheInstituteDH

Kawazoe Y, et al. PLoS One. 2021 Nov 9;16(11):e0259763.

### Wikipedia-BERT<sup>2</sup>

- BERT-base: 12 (L), 12 (A), 768 (E)
- Wikipedia (ja) (17 million lines)
- 32,000 vocabularies (BPE)
- Domain: General
- 2. https://alaginrc.nict.go.jp/nict-bert/index.html



#### 8 – 12 JULY 2023 | SYDNEY, AUSTRALIA



## Method: Joint NER-RE Model<sup>1</sup>





- NER assigns a suitably named entity tag to every token
- RE classifies the relation labels between tokens
- Joint NER-RE model simultaneously conducts NER and RE and optimizes the model weight.

@TheInstituteDH

#MEDINF023



8 – 12 JULY 2023 | SYDNEY, AUSTRALIA



# Method: Experiment settings

### Table 1. Statistics of iCorpus

Documents		183	
Average	Characters (S.D)	1,692 (593)	
	Entities (S.D)	394 (129)	
	Relations (S.D)	387 (127)	
Entity tags		113 (B-, I- + D)	
Relation tags		36 (35 + None)	

### Evaluate the two Joint NER-RE models

**Metrics**: The average of Macro-F1 and Micro-F1 through 5-fold cross-validation

• Train 64%, Val 16%, Test 20%

**Macro-F1:** Calculated by the total true positives, false positives, and false negatives.

Micro-F1: The average of F1 scores across multiple classes







### 8 - 12 JULY 2023 | SYDNEY, AUSTRALIA

### Results

	NER		RE	
	<b>Micro-F1</b>	<b>Macro-F1</b>	<b>Micro-F1</b>	<b>Macro-F1</b>
	(95%C1)	(95%C1)	(95%CI)	(95%CI)
Clinical-BERT	<mark>0.912</mark> *	<b>D.601</b>	<mark>0.759 *</mark>	<b>D.611</b>
	(0.907-0.916)	(0.580-0.621)	(0.757-0.761)	(0.585-0.637)
Wiki-BERT	0.892	0.586	0.741	0.591
	(0.884-0.899)	(0.567-0.606)	(0.734-0.747)	(0.571-0.611)



### Discussions: Comparison with related research

<ol> <li>i2b2/VA dataset <sup>1</sup></li> <li>Discharge summary (en)</li> <li>Entity: 3, Relationships: 8</li> <li>NER 0.859, RE: 0.737 (Micro F1)</li> </ol>	<ul> <li>JaMIE <sup>2</sup></li> <li>Medical history Reports (jp)</li> <li>Entity types: 9, Relationships: 10</li> <li>NER 0.855, RE: 0.710 (Micro F1)</li> </ul>
<ul> <li>JaMIE <sup>2</sup></li> <li>Radiology reports (jp)</li> <li>Entity: 9, Relationships: 10</li> <li>NER 0.956, RE: 0.865 (Micro F1)</li> </ul>	<ul> <li>4. iCorpus (this study)</li> <li>Case Reports (jp)</li> <li>Entity types: 56, Relationships: 35</li> <li>NER 0.912, RE: 0.759 (Micro F1)</li> </ul>

1. Uzuner Ö et al. 2010 i2b2/VA challenge on concepts, assertions, and relations in clinical text. J Am Med Inform Assoc. 2011 Sep-Oct;18(5):552-6.

2. Fei Cheng et al. JaMIE: A Pipeline Japanese Medical Information Extraction System. arXiv:2111.04261.

@TheInstituteDH #MEDINF023



## Conclusions

- 1. We showed the baseline performances of NER and RE on iCorpus.
  - iCorpus: Manually annotated 56 types of entities and 35 types of relationships.
  - NER: Micro-F1 of 0.913, Macro-F1 of 0.601 (Clinical-BERT)
  - RE: Micro-F1 of 0.759 and Macro-F1 of 0.611 (Clinical-BERT)
- 2. We also showed the domain-specific clinical BERT tended to show better performance than the Wikipedia BERT.