





# Building a Document-Level Relation Extraction Dataset Assisted by Cross-Lingual Transfer

Youmi Ma, An Wang, Naoaki Okazaki Tokyo Institute of Technology {youmi.ma@nlp., an.wang@nlp., okazaki@}c.titech.ac.jp



#### Outline

# **Building a DocRE Dataset w. xLingual Transfer**

### Purpose

 explore how to build Document-level Relation Extraction (DocRE) datasets with minimal human efforts

### Method

- automatic annotation using cross-lingual transfer -> 12
- human annotation assisted by cross-lingual transfer ->

### Contributions & Findings

- collected the first Japanese DocRE dataset
- showed that although the automatic annotation is not ready for use on its own, it serves as a good start point for human annotation

#### Outline

# **Building a DocRE Dataset w. xLingual Transfer**

### Background & Motivation

- task definition: DocRE
- existing DocRE datasets

### Approaches

- automatic annotation: Re-DocRED<sup>ja</sup>
- semi-automatic annotation: JacRED
- Dataset Analysis & Experiments
  - superiority of our annotation approach
  - usefulness of our collected dataset

# Background Relation Extraction (RE)

Task to extract relations from natural language texts

- Identify relations between entity pairs within a sentence



#### Background

## **Document-level Relation Extraction (DocRE)**

- Sentence-level RE is over-simplified<sub>(Yao+, 2019)</sub>
  - Relations exist beyond sentence boundaries



# **Document-level Relation Extraction (DocRE)**

Task to decide relations between all entity pairs in a document

### challenge

decide relations based on information from the whole document

#### The Archbishop

[1] "The Archbishop" is the third episode of the first series of the <u>BBC</u> sitcom *Blackadder* (*The Black Adder*). [2] It is set in <u>England</u> in the late <u>15th century</u>, and follows the exploits of the fictitious *Prince Edmund* as he is invested as <u>Archbishop of Canterbury</u> amid a <u>Machiavellian plot</u> by the King to acquire lands from the <u>Catholic Church</u>. [3] ... [5] *Edmund*, faced with the threat of assassination, attempts to escape to <u>France</u> into self-imposed exile; and in a later scene, two drunk knights overhear <u>King</u> <u>Richard IV</u> exclaiming "Who will rid me of this turbulent priest?" [6] The words attributed to <u>King Henry II</u> which led to <u>Becket</u>'s death in <u>1170</u>, and embark on a mission to murder *Edmund*. [7] ...

Subject: *Prince Edmund* Object: *Blackadder*  Relation: present in work

Example from DocRED (Yao+, 2019)

# **Evidence in DocRE**<sub>(Yao+, 2019)</sub>

**Evidence:** Minimal set of sentences enough for relation decision

### challenge

decide relations based on information from the whole document

### information filtering

Focus more on **evidence sentences** relevant to current entity pair

#### The Archbishop

[1] "The Archbishop" is the third episode of the first series of the <u>BBC</u> sitcom *Blackadder* (*The Black Adder*). [2] It is set in <u>England</u> in the late <u>15th century</u>, and follows the exploits of the fictitious *Prince Edmund* as he is invested as <u>Archbishop of Canterbury</u> amid a <u>Machiavellian</u> plot by the King to acquire lands from the <u>Catholic Church</u>. [3] ... [5] *Edmund*, faced with the threat of assassination, attempts to escape to <u>France</u> into self-imposed exile; and in a later scene, two drunk knights overhear <u>King</u> <u>Richard IV</u> exclaiming "Who will rid me of this turbulent priest?" [6] The words attributed to <u>King Henry II</u> which led to <u>Becket</u>'s death in <u>1170</u>, and embark on a mission to murder *Edmund*. [7] ...

Subject: *Prince Edmund* Object: *Blackadder*  Relation: *present in work* Evidence: 1,2

#### **Motivation**

# **Annotating DocRE from Scratch is Difficult**

Heavy burden for annotators affect the quality of collected dataset

"The Archbishop" is the third episode of the first series of the BBC sitcom Blackadder. Identify all relations

in this sentence.



#### The Archbishop

[1] "The Archbishop" is the third episode of the first series of the <u>BBC</u> sitcom *Blackadder* (*The Black Adder*). [2] It is set in <u>England in the late 15th century</u>, and follows the exploits of the fictitious *Prince Edmund* as he is invested as <u>Archbishop of Canterbury</u> amid a <u>Machiavellian</u> plot by the King to acquire lands from the <u>Catholic Church</u>. [3] ... [5] *Edmund*, faced with the threat of assassination, attempts to escape to <u>France</u> into self-imposed exile; and in a later scene, two drunk knights overhear <u>King Richard IV</u> exclaiming "Who will rid me of this turbulent priest?" [6] The words attributed to <u>King Henry II</u> which led to <u>Becket</u>'s death in <u>1170</u>, and embark on a mission to murder *Edmund*. [7] ...

Identify all relations

in this document.



#### **Motivation**

### **Existing DocRE Datasets**

Datasets in English, Chinese and Korean are available

- collected individually despite of high human annotation costs
- Hinders DocRE research from scaling up

Dataset	Language	# Instances	# Docs	Evidence
DocRED <sub>(Yao+, 2019)</sub>	en.	50,503	4,051	Y
Re-DocRED <sub>(Tan+, 2022)</sub>	en.	120,664	4,053	N
HacRED <sub>(Cheng+, 2021)</sub>	zh.	56,798	7,731	N
HistRED <sub>(Yang+, 2023)</sub>	kr.	9,965	5,816	Y The c Same Docf son

#### **Motivation**

# **Cross-Lingual Transfer (Projection)**

- Dataset for sentence-level RE has been successfully created with translation-based cross-lingual transfer
  - Human evaluation ensured the quality of obtained dataset



Figure 1: Example translations from English to German, Polish, Turkish and Chinese with XML markup for the head and tail entities to project relation argument annotations.

(Hennig+, 2023)

# **Purpose of This Work**

To explore how DocRE dataset in one language could help collecting DocRE dataset in another using cross-lingual transfer

- To publish a <u>Japanese</u> DocRE dataset ready for use
  - one of the most widely-used languages for Web content
     one of the most linguistically distant languages from English

#### Outline

# **Building a DocRE Dataset w. xLingual Transfer**

- Background & Motivation
  - task definition: DocRE
  - existing DocRE datasets

### Approaches

- automatic annotation: Re-DocRED<sup>ja</sup>
- semi-automatic annotation: JacRED
- Dataset Analysis & Experiments
  - superiority of our annotation approach
  - usefulness of our collected dataset

#### Approach: Automatic Annotation

## **Starting from English language resources**

### ■ **Re-DocRED**<sup>ja</sup>:Translate Re-DocRED<sub>(Tan+, 2022)</sub> into Japanese

<e0:LOC> Morogoro Region </e0:LOC> is one of <e1:LOC> Tanzania </e1:LOC> 's <e2:NUM> 31 </e2:NUM> administrative regions .

The regional capital is the municipality of <e3:LOC> Morogoro </e3:LOC> .

According to the <e4:MISC> 2012 national census </e4:MISC> , the region had a population of <e5:NUM> 2,218,492 </e5:NUM> , which was higher than the pre - census projection of <e6:NUM> 2,209,072 </e6:NUM>.



#### Approach: Automatic Annotation

### Train a DocRE model on Re-DocRED<sup>ja</sup>

■ The model fails to extract many relations on Japanese Wikipedia

Topic Shift of Contents	<b>JA: <u>堀 直宥</u>(ほり なおさだ、寛文5年11月17日(1665年12月23日) - 正徳元年6月8日</b> (1711年7月23日))は、江戸時代前期から中期の大名で、 <mark>上総八幡藩</mark> 第2代 <mark>藩主</mark> 。
Little / No contents about Japanese bistory / figures /	EN: <u>Naosada Hori</u> (December 23, 1665 - July 23, 1711) was a feudal lord of the early to mid-Edo period, the second lord of the <u>Joso Hachiman domain</u> .
architecture	missed relation: (Naosada Hori, head of government, Joso Hachiman domain)
Gap of Surface Structures	JA: <u>ザカリアーシュ・ヨージェフ</u> (1924年3月25日 - 1971年11月22日)は、ハンガ リー出身のサッカー選手、サッカー指導者。1954年のFIFAワールドカップでは決 勝戦を除く4試合にフル出場し準優勝に貢献した。

### Approach: Semi-Automatic Annotation

### **Pipeline of Human Annotation**

Annotator's job: edit recommendations from models



![](_page_14_Figure_6.jpeg)

#### Semi-Automatic Annotation: Proposal 1

# **Recommend Relations with Model Predictions**

Utilize model trained on translated dataset to recommend relations Prior Works

![](_page_15_Figure_5.jpeg)

### Semi-Automatic Annotation: Proposal 2 Refine Relation Label Set

- Merge labels based on:
  - Frequency: select most-frequent relation labels
  - Hierarchy:
    - Merge sub-properties into super-properties
      - E.g. author -> creator
    - Merge inverse properties
      - E.g. has part(s) -> part of
  - **Similarity**: pretrained graph embedding: <u>GraphVite</u>(Zhu+, 2019)
- Reduced relation labels from 96 to 35
  - While keeping >88% relation instances in Re-DocRED

	author			
S	subproperty of		e creator	
			▼ 0 references	
	has part(s)			
	inverse property	e part of		
			▼ 0 references	

### Semi-Automatic Annotation

### **Collected Dataset: JacRED**

■ <u>Japanese Document-level Relation Extraction Dataset</u>

Dataset	Language	# Instances	# Docs	Evidence
DocRED <sub>(Yao+, 2019)</sub>	en.	50,503	4,051	Y
Re-DocRED <sub>(Tan+, 2022)</sub>	en.	120,664	4,053	Ν
HacRED <sub>(Cheng;, 2021)</sub>	zh.	56,798	7,731	Ν
HistRED <sub>(Yang+, 2023)</sub>	kr.	9,965	5,816	Y
JacRED	ja.	42,241	2,000	Y

#### Outline

# **Building a DocRE Dataset w. xLingual Transfer**

### Background & Motivation

- task definition: DocRE
- existing DocRE datasets
- Approaches
  - automatic annotation: Re-DocRED<sup>ja</sup>
  - semi-automatic annotation: JacRED

### Dataset Analysis & Experiments

- superiority of our annotation approach
- usefulness of our collected dataset

#### **Dataset Analysis**

# **Statistics: Comparison with Existing Datasets**

- JacRED has its advantages over existing datasets as a general language resource
  - Even without the distinctiveness of language

# Sentences, # Entities and # Relations averaged over documents
# Evidences averaged over relation instances

	# Sentences	# Entities	# Relations	# Evidences
DocRED <sub>(Yao+, 2019)</sub>	7.98	19.51	12.45	1.60
Re-DocRED <sub>(Tan+, 2022)</sub>	7.98	19.45	29.77	0.88
JacRED	8.39	17.87	17.87	1.67
	more rel	ation instances	more evide	nce instances

#### **Dataset Analysis**

### **Statistics: Number of Human Edits**

- Sample 400 documents from JacRED and calculate the number of edit steps before reaching at final annotations
  - Starting from model predictions reduces the number of human edit steps

![](_page_20_Figure_6.jpeg)

# **Training on Translated Dataset**

- Evaluate DREEAM<sub>(Ma+, 2023)</sub> trained with Re-DocRED<sup>ja</sup> on JacRED
  - Number of documents used during training indicated in parenthesis

![](_page_21_Figure_6.jpeg)

#### Experiments

### **Training on Translated Dataset**

- Evaluate DREEAM<sub>(Ma+, 2023)</sub> trained with Re-DocRED<sup>ja</sup> on JacRED
  - Number of documents used during training indicated in parenthesis

![](_page_22_Figure_6.jpeg)

#### Summary

# Building a DocRE Dataset w. xLingual Transfer

### Purpose

- explore how to build Document-level Relation Extraction (DocRE) datasets with minimal human efforts
- Method
  - automatic annotation using cross-lingual transfer -> 12
  - human annotation assisted by cross-lingual transfer ->
- Contributions & Findings
  - collected the first Japanese DocRE dataset
  - showed that although the automatic annotation is not ready for use on its own, it serves as a good start point for human annotation

![](_page_23_Picture_12.jpeg)

# References

- Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin, Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou, and Maosong Sun. DocRED: A Large-Scale Document-Level Relation Extraction Dataset. ACL 2019.
- Qiao Cheng, Juntao Liu, Xiaoye Qu, Jin Zhao, Jiaqing Liang, Zhefeng Wang, Baoxing Huai, Nicholas Jing Yuan, and Yanghua Xiao. HacRED: A Large-Scale Relation Extraction Dataset Toward Hard Cases in Practical Applications. ACL-IJCNLP 2021 Findings.
- Soyoung Yang, Minseok Choi, Youngwoo Cho, and Jaegul Choo. HistRED: A Historical Document-Level Relation Extraction Dataset. ACL 2023.
- Leonhard Hennig, Philippe Thomas, and Sebastian Möller. MultiTACRED: A Multilingual Version of the TAC Relation Extraction Dataset. ACL 2023.
- Qingyu Tan, Lu Xu, Lidong Bing, Hwee Tou Ng, Sharifah Mahani Aljunied. Revisiting DocRED -- Addressing the False Negative Problem in Relation Extraction. EMNLP 2022.
- Youmi Ma, An Wang, Naoaki Okazaki. DREEAM: Guiding Attention with Evidence for Improving Document-Level Relation Extraction. EACL 2023.
- Zhaocheng Zhu, Shizhen Xu, Meng Qu, Jian Tang. GraphVite: A High-Performance CPU-GPU Hybrid System for Node Embedding. WWWC 2019.
- Qingyu Tan, Ruidan He, Lidong Bing, and Hwee Tou Ng. Document-Level Relation Extraction with Adaptive Focal Loss and Knowledge Distillation. ACL 2022 Findings.
- Wenxuan Zhou, Kevin Huang, Tengyu Ma, and Jing Huang. Document-level relation extraction with adaptive thresholding and localized context pooling. AAAI 2021.
- Ningyu Zhang, Xiang Chen, Xin Xie, Shumin Deng, Chuanqi Tan, Mosha Chen, Fei Huang, Luo Si, and Huajun Chen. Document- level relation extraction as semantic segmentation. IJCAI 2021.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. 2022.