

ARC Day 2020

ARC所蔵資料データベースの評価 データセット作成プロジェクト

前田 亮(情報理工学部)
 LI Kangying(情報理工学研究科 D2)
 WANG Jiayun(情報理工学研究科 D3)
 SONG Yuting(情報理工学部 特任助教)
 BATJARGAL Biligsaikhan(衣笠総合研究機構 専門研究員)
 佐藤 英男(情報理工学研究科 M2)

1

研究の概要

- ARC所蔵の日本文化資源データベースのオープンデータ化および容易な多言語情報アクセスの実現に向けて研究を行っている



ARC浮世絵ポータルデータベース

ARC古典籍ポータルデータベース

2

現状の課題

- 情報系の研究では、提案手法の評価が重要
- 浮世絵や古典籍のデータについて、評価用のデータセットが存在しない
 - 現時点では、研究室内で小規模なデータセットを作成して使用
 - 小規模なデータセットでは、信頼性に欠ける
 - 研究室内で大規模なデータセットの構築は難しい
- クラウドソーシング、ゲーミフィケーションを用いた評価用データセットの作成手法を検討中

3

3

紹介する研究

- Multimodal Representation learning for Ukiyo-e records retrieval and analysis**
LI Kangying (情報理工学研究科 D2)
- ARC所蔵浮世絵資料のための作品推薦・同定手法**
WANG Jiayun (情報理工学研究科 D3)
- Finding Identical Ukiyo-e Prints across Databases in Japanese, English and Dutch**
SONG Yuting (情報理工学部 特任助教)

4

Exploring information of Ukiyo-e records by using Multi-source Data

Multimodal Representation Learning for Ukiyo-e Records Retrieval and Analysis

LI Kangying

5

Outline

- Motivation
- Some progress:
 - Pre-processing of titles
 - Compare some morphological analyzers
 - Build some user dictionaries
 - Token embedding from BERT's pre-trained model
 - Visualization of the results

6

1

Motivation

- If we only use a single database to search some items:



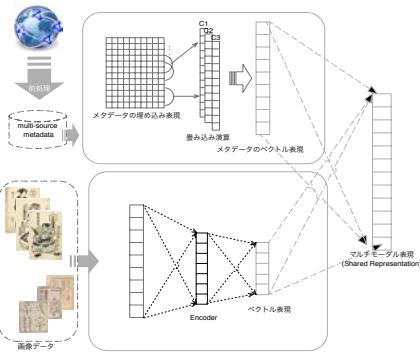
Sorry!
No image

~~metadata without image information~~

Motivation

- Now there are many multilingual open source databases, what we can do is...

- ARC database
 - WikiData
 - Wikipedia
 - WikiArt
 - ...

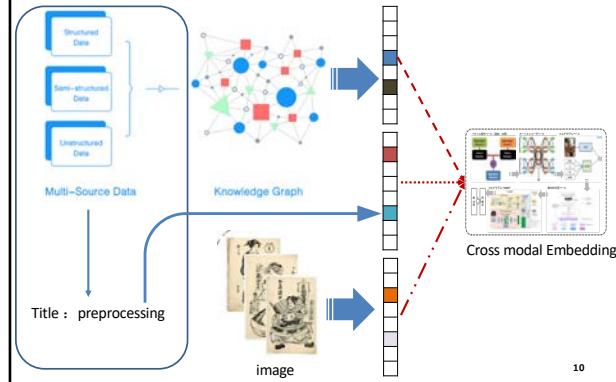


7

If we can build a system that...



Can we build it by?



10

Outline

- ## ● Motivation

- ## ● Some progress:

Pre-processing of titles

- Compare some morphological analyzers
 - Build some user dictionaries
 - Token embedding from BERT's pre-trained model
 - Visualization of the results

1

Pre-processing of titles: morphological analyzer (proper noun extraction)

- **Input sentence:**この時期は役者絵や合巻の挿絵などを描いていたが、あまり人気が出ず作品も僅かであった。また、勝川春亭にも学んでおり、さらに萬葉北斎の影響を受け、後に三代豊国等に学んで、雪舟とも昌した。

この時期は役者絵や合巻の	連体形,*,*,*この,コノ,コノ 名詞副詞可能,*,*,*時期,ジキ,ジキ 助詞,系助詞,*,*,*は,ハ,フ 名詞一般,*,*,*役者,ヤクシャ,ヤクシャ 名詞接尾一般,*,*,*絵,エ 助詞,並立助詞,*,*,*や,ヤ,ヤ 名詞一般,*,*,*合ゴウ,ゴー 名詞一般,*,*,*巻,マキ,マキ 助詞,連体化,*,*,*の,ノ,ノ
--------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

名詞, 固有名詞, 地域, 一般, *, *, 勝川, カチガワ, カチガワ
名詞, 一般, *, *, *, 春, ハル, ハル
名詞, 接尾, 一般, *, *, 亭, テイ, テイ
助詞, 格助詞, 一般, *, *, に, ニ, ニ

By: mecab-unidic-neologd

12

11

12

Pre-processing of titles: morphological analyzer

この	コノ	この	連体詞
時期	ジキ	時期	名詞-副詞可能
は	ハ	は	助詞-係助詞
役者	ヤクシャ	役者	名詞-一般
絵	エ	絵	名詞-接尾-一般
や	ヤ	や	助詞-並立助詞
合	ゴウ	合	名詞-一般
巻	マキ	巻	名詞-一般
の	ノ	の	助詞-連体化
挿絵	サシエ	挿絵	名詞-一般
勝川	カチガワ	勝川	名詞-固有名詞-地域-一般
春亭	ハル	春亭	名詞-一般
に	ニ	に	助詞-格助詞-一般

By: 茶菴 Ochasan
13

13

Outline

● Motivation

● A book review

● Some progress:

Pre-processing of titles

– Compare some morphological analyzers

– Build some user dictionaries

– Token embedding from BERT's pre-trained model

– Visualization of the results

14

14

User dictionary

- Items, Events, Person's names and Stories Contained in Ukiyo-e:
人物画、風景画、伝説と神話、広告等の情報...

reference : <https://www.kumon-ukiyo.e.jp/history.html>



15

浮世絵用語辞典&日本人名辞書 (JBDB)

綱春水	らいしんすい	Rai Shunsui
阿波黒仁左衛門	あわくろいざえもん	Awaya Nizaemon
水戸笠崎庵	みとかさきあん	Nagatomo Dokushūan
福原玉五郎	ふくはらぎょく	Izumiya Shōsuke
伊達義重	いだぎじゅう	Taira Nizaemon
三浦屋仁左衛門	みうらやいざえもん	Shioya Nizaemon
北村義之	きたむらぎよし	Harmiya Gonbei
林喜吉	はやしうきよ	Hayashi Kōzai
留守道盛	るすたいどうせう	Rusu Taizō
林喜介	はやしうけいすけ	Hayashi Kiseiuke
愛宕里坊	あたごりぼう	Atago Bō
芥川龍之介	あくたがわりょうしけん	Akutagawa Yōken
高志義造	たかしよこう	Takashi Yōkō
丹後屋藤兵衛	たんごやとうべえ	Tangoya Tōbe'e
延徳	ちょうとく	Chōrin
青泉齋市	かいせんさいいち	Kaisen Kiichi
青石大膳	あおいだいじゅん	Aoshi Daishun
葛子	かつしん	Katsu Shinkin
木村萬蔵	きむらまんぞう	Kimura Kenkado
上林寅次郎	かんばらいんじじろう	Kambayashi Kichijirō
竜草屋	りゅうそうや	Ryō Sōro
無安和尚	むせんおしょう	Musen Oshō

16

16

歴史地名データ: 139,454 words

- For knowledge extraction: 地名一別名一属性一上位地名一上位地名属性一上位地名別名
- For user dictionary: 地名 + 別名 + 上位地名 + 上位地名属性 (.dic)



https://www.nihu.jp/ja/publication/source_map

17

Morphological analysis

• 「幡枝村円通寺」

幡
枝
村
圓
通
寺
名詞, 固有名詞-人名,姓,*,*幌,ハタ,ハタ
名詞, 固有名詞-地域-一般,*,*枝,エダ,エダ
名詞, 地域-一般,*,村,ムラ,ムラ
名詞, 固有名詞-地域-般,*,円通寺,エンツウジ,エンツウジ
mecab-unicid-neologd

Juman++ Demo

幡枝村円通寺
名詞, 固有名詞-人名,姓,*,*幌,ハタ,ハタ
名詞, 固有名詞-地域-一般,*,*枝,エダ,エダ
名詞, 地域-一般,*,村,ムラ,ムラ
名詞, 固有名詞-地域-般,*,円通寺,エンツウジ,エンツウジ
EOS

chasen

幡
枝
村
圓
通
寺
名詞, 固有名詞-人名,姓,*,*幌,ハタ,ハタ
名詞, 固有名詞-地域-一般,*,*枝,エダ,エダ
名詞, 地域-一般,*,村,ムラ,ムラ
名詞, 固有名詞-地域-般,*,円通寺,エンツウジ,エンツウジ
EOS

mecab-unicid-neologd+userdic

18

17

3

Morphological analysis

- INPUT title: 仮名手本忠臣蔵 第十

-Ochasen	mecab-unidic-neologd	mecab-unidic-neologd+userdic
仮名 手本 忠臣蔵 記号:空白 第 ダイ 接頭詞:数跡続 +	カタカナ一般 手本 名詞一般 チラシシングラ 名詞:若有名詞一般 記号:空白 第 ダイ 接頭詞:数跡続 +	名詞一般,*****仮名:カメ 手本 名詞一般,*****手本子ホ 忠臣蔵,チラシシングラ 名詞固有名詞一般,***, 忠臣蔵,チラシシングラ 記号:空白***** 第 ダイ 名詞數,****,+,ジウ EOS
		仮名手本忠臣蔵 名詞 固有名詞一般 姓*** 仮名手本忠臣蔵,かなでほん ちゅうじんざく,かなでほんちゅうじん 記号:空白***** 第接頭詞数跡続****,*,第 名詞數****,+,ジウ EOS

19

Outline

- Motivation
- A book review
- Some progress:
 - Pre-processing of titles
 - Compare some morphological analyzers
 - Build some user dictionaries
 - Token embedding from BERT's pre-trained model
 - Visualization of the results

20

19

20

Token embedding from BERT's pre-trained model

MODEL (Google): BERT-Base, Multilingual Cased (New, recommended): 104 languages, 12-layer, 768-hidden, 12-heads, 110M parameters

First step: clean the data ,Because it contains a lot of symbols		
e.g. [○★●•[]「」?#-「」『』()>[]`]		「阿蘭々仙人」「☆豊次郎」 「耶敷陀羅姫」「羅☆羅太子」「車匿舍人」
Second step: Grouping the experimental subject:		
小倉擬百人一首 木津勘助 仮名手本忠臣蔵 古今忠孝伝 木津閑兵衛 中村駒之助 市川団十郎 八大伝大乃そうしの内 黒谷金戒光明寺 西山月輪寺 幡枝村円通寺	['小倉', '擬百人一首'], ['木津', '勘助'], ['仮名手本忠臣蔵'], ['古今', '忠孝', '伝'], ['木津', '閑兵衛', '中村', '駒之助'], ['市川', '団十郎'], ['八大', '伝', '大', '乃', 'そう', 'し', 'の', '内'], ['黒谷', '金戒', '光明', '寺'], ['西山', '月輪', '寺'], ['幡枝', '村', '円通', '寺']	['小', '倉', '擬', '百', '人', '一', '首'], ['木', '津', '勘', '助'], ['仮', '名', '手', '本', '忠', '臣', '藏'], ['古', '今', '忠', '孝', '伝'], ['木', '津', '閑', '兵', '衛', '中', '村', '駒', '之', '助'], ['市', '川', '团', '十', '郎'], ['八', '大', '伝', '大', '乃', 'そう', 'し', 'の', '内'], ['黒', '谷', '金', '戒', '光', '明', '寺'], ['西', '山', '月', '輪', '寺'], ['幡', '枝', '村', '円', '通', '寺']
Sentence	Words-level	Characters-level

21

Compareing the performance

Query: 幡枝村円通寺 Database: [小倉擬百人一首, 木津勘助, 仮名手本忠臣蔵, 古今忠孝伝, 木津閑兵衛, 中村駒之助, 市川団十郎, 八大伝大乃そうしの内, 黒谷金戒光明寺, 西山月輪寺, 幡枝村円通寺]

Nearest points in the original space:	Nearest points in the original space:	Nearest points in the original space:	
黒谷金戒光明寺	0.669	黒谷金戒光明寺	0.655
西山月輪寺	0.728	古今忠孝伝	0.724
木津閑兵衛 中村駒之助	1.096	八大伝大乃そうしの内	0.785
八大伝大乃そうしの内	1.118	木津閑兵衛 中村駒之助	0.807
古今忠孝伝	1.193	八大伝大乃そうしの内	1.294
市川団十郎	1.227	西山月輪寺	1.294
木津勘助	1.271	古今忠孝伝	1.173
小倉擬百人一首	1.291	市川団十郎	1.246
仮名手本忠臣蔵	1.373	小倉擬百人一首	1.256
		木津閑兵衛	1.279
		西山月輪寺	1.437
		仮名手本忠臣蔵	1.389

Sentence Words-level Characters-level

22

compare the performance

Query: 木津勘助 Database: [小倉擬百人一首, 仮名手本忠臣蔵, 古今忠孝伝, 木津閑兵衛, 中村駒之助, 市川団十郎, 八大伝大乃そうしの内, 黒谷金戒光明寺, 西山月輪寺, 幡枝村円通寺]

Sentence	Words-level	Characters-level	
木津閑兵衛 中村駒之助	0.630	木津閑兵衛	0.608
八大伝大乃そうしの内	1.233	黒谷金戒光明寺	0.000
西山月輪寺	1.279	西山月輪寺	0.000
木津閑兵衛 中村駒之助	1.290	木津閑兵衛 中村駒之助	0.745
小倉擬百人一首	1.034	小倉擬百人一首	1.045
仮名手本忠臣蔵	1.074	古今忠孝伝	1.049
市川団十郎	1.298	市川団十郎	1.159
八大伝大乃そうしの内	1.370	木津閑兵衛	1.159
幡枝村円通寺		幡枝村円通寺	1.271
黒谷金戒光明寺		黒谷金戒光明寺	1.290
西山月輪寺		西山月輪寺	1.290
古今忠孝伝		古今忠孝伝	1.308
木津閑兵衛 中村駒之助		木津閑兵衛 中村駒之助	1.654
市川団十郎		市川団十郎	1.681
八大伝大乃そうしの内		古今忠孝伝	1.369

23

Outline

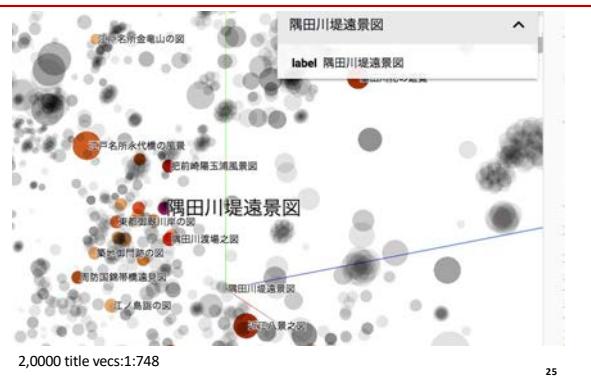
- Motivation
- A book review
- Some progress:
 - Pre-processing of titles
 - Compare some morphological analyzers
 - Build some user dictionaries
 - Token embedding from BERT's pre-trained model
 - Visualization of the results

24

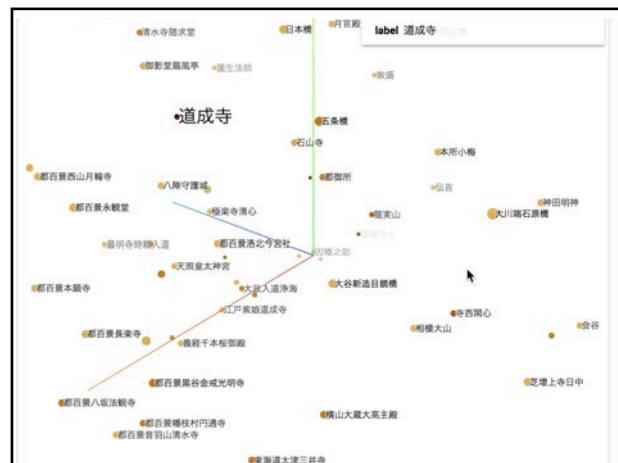
23

24

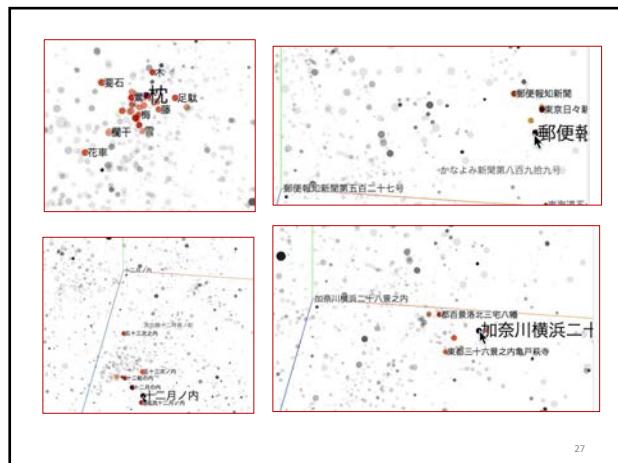
Visualization



25

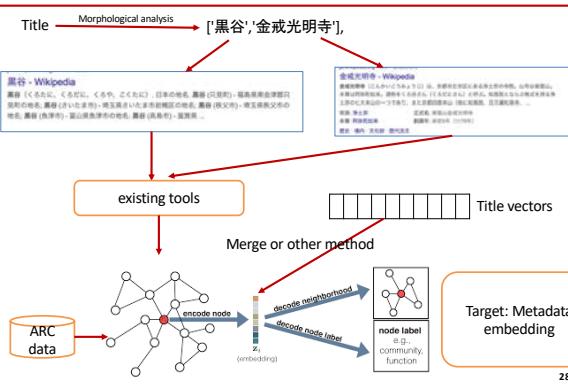


26



27

Plans



28

Plans

- Vector data analysis
- Extend user dictionary
- About the new research of seals retrieval: in the experimental stage..

29

ARC所蔵浮世絵資料のための 作品推薦・同定手法

WANG Jiayun



30

Finding Identical Ukiyo-e Prints across Databases in Japanese, English and Dutch

SONG Yuting



37

Ukiyo-e

- Japanese traditional woodblock print
- One of the popular arts of the Edo period (1603–1868)
- Many libraries, museums and galleries in Western countries have digitalized ukiyo-e woodblock prints



Title: Under the Wave off Kanagawa (Kanagawa-oki nami-ura), also known as The Great Wave. From the series Thirty-six Views of Mount Fuji (Fugaku sanjūrokkei).
Artist: Katsushika Hokusai (Japanese, Tokyo (Edo) 1760–1849 Tokyo (Edo))
Period: Edo period (1615–1868)
Date: ca. 1829–32
Culture: Japan

Image source: The Metropolitan Museum of Art, <https://www.metmuseum.org/art/collection/search/45434>

Japan

United States

United Kingdom

Netherlands

Edo-Tokyo Museum
Ritsumeikan University

...

Metropolitan Museum of Art
Museum of Fine Arts

...

British Museum
Ashmolean Museum

...

Rijksmuseum

38

38

Identical ukiyo-e prints in different languages

	Metadata		Language
	Title	Artist	
江戸東京博物館 (Japan)	雪月花 淀川	葛飾北斎	English
	凱風快晴	葛飾北斎	
Metropolitan Museum of Art (United States)	Moonlight on the Yodo River, from the series Snow, Moon, and Flowers	Katsushika Hokusai	English
	Morning Mist at Mishima	Utagawa Hiroshige (I)	
Rijksmuseum (Netherlands)	Helder weer en een zuidelijke wind	Katsushika Hokusai	English
	Mishima in ochtendmist	Hiroshige (I), Utagawa	

4

39

Motivation

- To find identical ukiyo-e prints in different languages
 - by comparing metadata

作品名	作者	Title	Artist
富嶽三十六景 神奈川沖浪裏	葛飾北斎	Under the Wave off Kanagawa, from the series Thirty-six Views of Mount Fuji	Katsushika Hokusai
富嶽三十六景 深川年輪橋下	葛飾北斎	Snow on the Sumida River, from the series, Snow, Moon, and Flowers	Katsushika Hokusai
日本橋 朝之景	歌川広重(初代)		
雪月花 濱田	葛飾北斎	Morning View of Nihonbashi	Utagawa Hiroshige

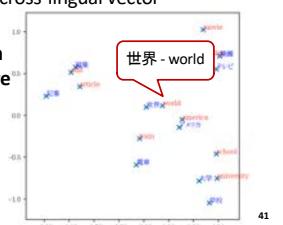
Challenges

- Metadata values are in different languages
 - Language barriers: similarity measures cannot be employed directly

40

Our proposed method

- To calculate metadata similarity in different languages by using bilingual word embeddings
 - Bilingual word embeddings
 - Word embeddings: dense, low-dimensional and real-valued vectors for representing words
 - Bilingual word embeddings: cross-lingual vector representations of words
 - Advantage: similar words in different languages can have similar vector representations



41

41

Metadata similarity calculation (1/2)

- We measure the metadata similarities through word-to-word matching between metadata in different languages

Notation

- Let b be the vector space of Japanese-English bilingual word embeddings
- Each word w_m^J in Japanese metadata is represented as $b(w_m^J)$ by using the bilingual word embeddings.
- Similarly, each word w_n^E in English metadata can be represented as $b(w_n^E)$

Word similarity

- $roma(w_m^J)$ is the romanization of the Japanese word w_m^J
- $Similarity(w_m^J, w_n^E) = \max[\cos_sim(b(w_m^J), b(w_n^E)), \cos_sim(b(roma(w_m^J)), b(w_n^E))]$

The cosine similarity between w_m^J and w_n^E

The cosine similarity between the romanization of w_m^J and w_n^E

42

42

Metadata similarity calculation (2/2)

- We measure the metadata similarities through word-to-word matching between metadata in different languages

- Metadata similarity

- Each word w_m^J in the Japanese metadata is used to compare with all the words in English metadata
- The maximum similarity score is used as the contribution of w_m^J to the metadata similarity

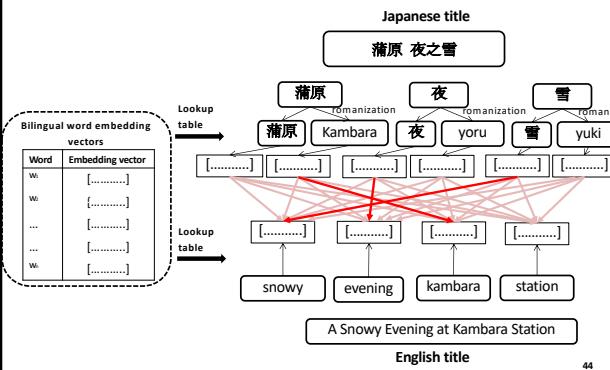
$$\text{Metadata Similarity} = \sum_{m=1}^{N^J} \max_{n \in \{1, 2, \dots, N^E\}} [\text{Similarity}(w_m^J, w_n^E)]$$

All the words in an English metadata

The word similarity between w_m^J and w_n^E

43

Example



A Snowy Evening at Kambara Station

44

43

44

Experiments

- Finding identical metadata records of ukiyo-e prints between Japanese, English and Dutch databases

- Experimental data

- Ukiyo-e prints dataset

- Metadata: title, artist name

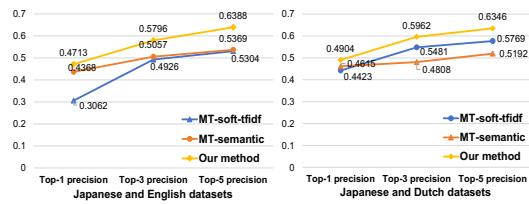


Language	Data source	Number of Records
Japanese	Edo-Tokyo Museum	203
English	Metropolitan Museum of Art	3,398
Japanese	Edo-Tokyo Museum	53
Dutch	Rijksmuseum	614

- Each Japanese ukiyo-e metadata record has at least one corresponding ukiyo-e metadata records in the English or Dutch dataset

45

Experimental results



- Compared with the methods that relies on machine translation(MT)

- Soft-tfidf similarity (MT-soft-tfidf)

- A string-based similarity metric and showed the best performance in title matching against 20 other commonly used string-based similarity metrics

- Semantic matching (MT-semantic)

- This method is based on monolingual word embeddings

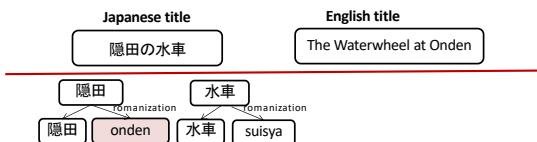
46

45

46

Discussion (1/2)

- Example: bilingual word embedding is better than MT



- MT-based method

- MT results of 「隠田の水車」
 - Hidden Water Wheel (Microsoft Translator)
 - Hidden Waterwheel (Google Translator)
 - water wheel in a hidden land (DeepL Translator)
 - The Water Wheel of Hidden Field (みらい翻訳)

47

Discussion (2/2)

- Limitation of our approaches

- The corresponding English title is inadequate translation



Utagawa Hiroshige, 1857

48

47

48

Translation patterns of Japanese titles

- 1) **Transliteration:** all the words in an original Japanese title are converted to the other language one by one via transliteration
 - Example: Okabe; Utsu no Yama (岡部 宇津之山)
- 2) **Literal translation:** all the words in Japanese title are converted to their equivalents in the other language, which remains the content of original Japanese title as much as possible
 - Example: Surugadai in Edo (東都駿台)
 - The Waterwheel at Onden (隱田の水車)
- 3) **Free translation:** The titles in the target language are created based on the content of Japanese objects without sticking to the original Japanese titles
 - Adding some information
 - Example: Snow on the Sumida River (隅田)
 - The Japanese title and corresponding English title have no meaning overlap

49

Translation patterns (cont')

- 1) **Transliteration: 64**
 - **Adequate transliteration:** 49
 - Example: Okabe; Utsu no Yama (岡部 宇津之山)
 - **Inadequate transliteration:** 15
 - Example: Yanagishima no Zu (柳嶋之図 楊本)
- 2) **Literal translation: 135**
 - **Adequate iteration translation:** 119
 - Example: The Waterwheel at Onden (隱田の水車)
 - **Inadequate iteration translation:** 16
 - Example: The Tanabata Festival (市中祭七夕祭)

50

Experiments

- **Datasets**

Language	Databases	Number of ukiyo-e prints
Japanese	Edo-Tokyo Museum	203
English	Metropolitan Museum of Art	3,398
- **Baseline methods**
 - Transliteration + softtfidf (new)
 - MT + softtfidf
 - MT + semantic
 - MT + mix
 - BiWE
- **Proposed method: BiWE+Roma**

51

Experimental setup

- **Transliteration + softtfidf**
 - MeCab + IPADic
 - Hepburn romanization system
 - **Softtfidf**
 - Secondary similarity function: edit distance
 - Similarity threshold: 0.9
- **MT + softtfidf**
 - Microsoft Translator Text API (translation results on 2020/06/28)
 - Statistical MT model (st)
 - Neural network MT model (nn)
- **MT + semantic / MT+ mix**
 - English word embeddings (Wikipedia articles + word2vec)
- **BiWE / Wib**
 - Japanese-English word embeddings
 - Monolingual word embeddings (Wikipedia articles + word2vec)
 - 9000 Japanese-English word pairs

52

Experimental results (MAP)

Table 1: Results (MAP)

	All	Transliteration		Literal translation		Free translation
		Adequate	Inadequate	Adequate	Inadequate	
Transliteration	0.4555	0.5818	0.3434	0.4518	0.375	0.1723
MT+softtfidf	0.4949	0.3931	0.4413	0.5474	0.3877	0.3147
MT+semantic	0.4696	0.375	0.37	0.5431	0.3085	0.2887
MT+mix	0.4693	0.4114	0.3618	0.5488	0.3177	0.2045
Bi-WE	0.2698	0.1009	0.0097	0.353	0.1719	0.2483
Bi-WE+Roma	0.5338	0.5671	0.1351	0.6046	0.2867	0.3288

- All datasets, adequate literal translation, free translation
 - Proposed method is better than baseline methods
- Adequate transliteration
 - Transliteration+softtfidf performs best

53

Conclusion

- We presented a metadata similarity calculation method by using bilingual word embeddings
- We showed our method's ability to find identical metadata records of ukiyo-e prints across Japanese, English and Dutch databases

Future work

- To further improve the performance by using other metadata
- To apply our methods on other humanities databases

54