



Extracting Lack of Information on Wikipedia by Comparing Multilingual Articles

- ☆ Yuya Fujiwara Konan University(Japan)
- Yu Suzuki Nagoya University(Japan)
- Yukio Konishi Konan University(Japan)
- Akiyo Nadamoto Konan University(Japan)



Background 1

- Many people all over the world use Wikipedia on the Internet.
- An important policy of Wikipedia is that the contents of articles is the same for all language version
- Articles of the same topic of any language version are expected to have exactly identical contents except for language.

**This policy is not obeyed,
especially for culture-related topics.**



Background 2



The content of article about “Fish and Chips” is very rich in the English version, but poor in the Japanese version.

Because “Fish and Chips” is a very popular dish in the U.K., but not in Japan.

There are some lack of information on one language Wikipedia, however there may be rich information on other language Wikipedia.

Query: Fish and Chips

Japanese version

English version



Japanese user

目次	
1	概略
2	歴史
3	食べ方
4	関連項目

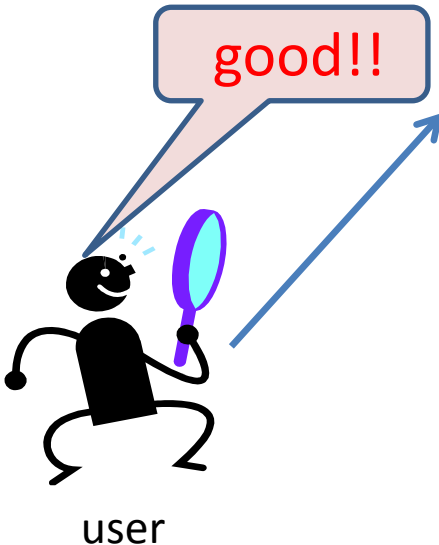
poor

Contents	
1	History
1.1	England
1.2	Scotland
1.3	Ireland
2	Composition
2.1	Cooking
2.2	Thickness
2.3	Batter
2.4	Choice of fish
2.5	Accompaniments
3	Vendors
4	Cultural impact
5	Environment
6	See also
7	Footnotes
8	External links

rich

Propose

- If there are lack of information in Wikipedia article, we complement it in native language version.



フィッシュ・アンド・チップス

この記事は検証可能な出典がまっ
出典を追加して記事の信頼性向上に

この項目では、料理について記述しています。1990年代後半

フィッシュ・アンド・チップス (英語: fish-and-chips または 英語:
るファーストフードの一つである手軽な食事。

目次 (非表示)

- 1 概略
- 2 歴史
- 3 食べ方
- 4 関連項目

概略 (編集)

タラやカレイ、オヒョウなどの白身魚の切り身に、小麦粉を卵や水ま
い棒状に切って油で揚げたチップスと合わせて供する。この場合
で言うフライドポテト (アメリカで言うフレンチフライ) のイギリスでの
の切り身小一切れにジャガイモ中一個分)で450キロカロリー程。

Ireland

Main article: Irish cuisine

In Ireland, the first fish and chips were sold by an Italian immigrant, Giuseppe C
started by selling fish and chips outside pubs from a handcart. He then found a
would ask customers "Uno di questa, uno di quella?" This phrase (meaning "on
which is still a way of referring to fish and chips in the city."^[5]

complement

Scotland

Lack of information

Ireland

Main article: Irish cuisine

In Ireland, the first fish and chips were sold by an Italian immigrant, Giuseppe C
started by selling fish and chips outside pubs from a handcart. He then found a
would ask customers "Uno di questa, uno di quella?" This phrase (meaning "on
which is still a way of referring to fish and chips in the city."^[5]

Composition

Cooking



Traditional frying uses *beef dripping*
of vendors in the north of England :
dish, but it has the side effect of m
museums, such as the *Black Count*

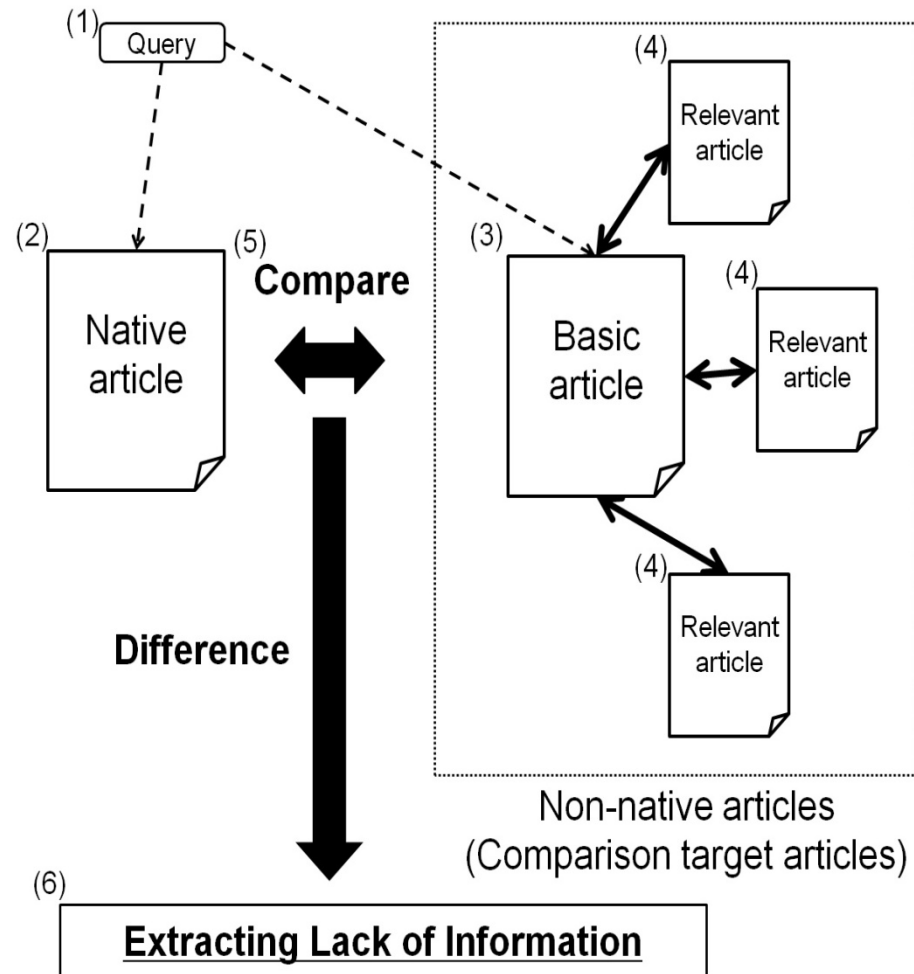
Thickness

British chips are usually significant
food chains, resulting in a lower fat

Extracting Lack of information on Wikipedia
by Comparing Multilingual Articles

Our Flow

1. Users input a query in their native language to the system.
2. The system retrieves one native article of which title is the same as the user's input query.
3. It translates the query to the non-native language using a language dictionary and retrieves a non-native article whose the title is the same as the user's input query.
4. It extracts comparison articles from the non-native articles using a Wikipedia link graph.
5. It compares a native article with non-native articles extracted in 4. and extracts lack of information.
6. It browses lack of information available on the web.



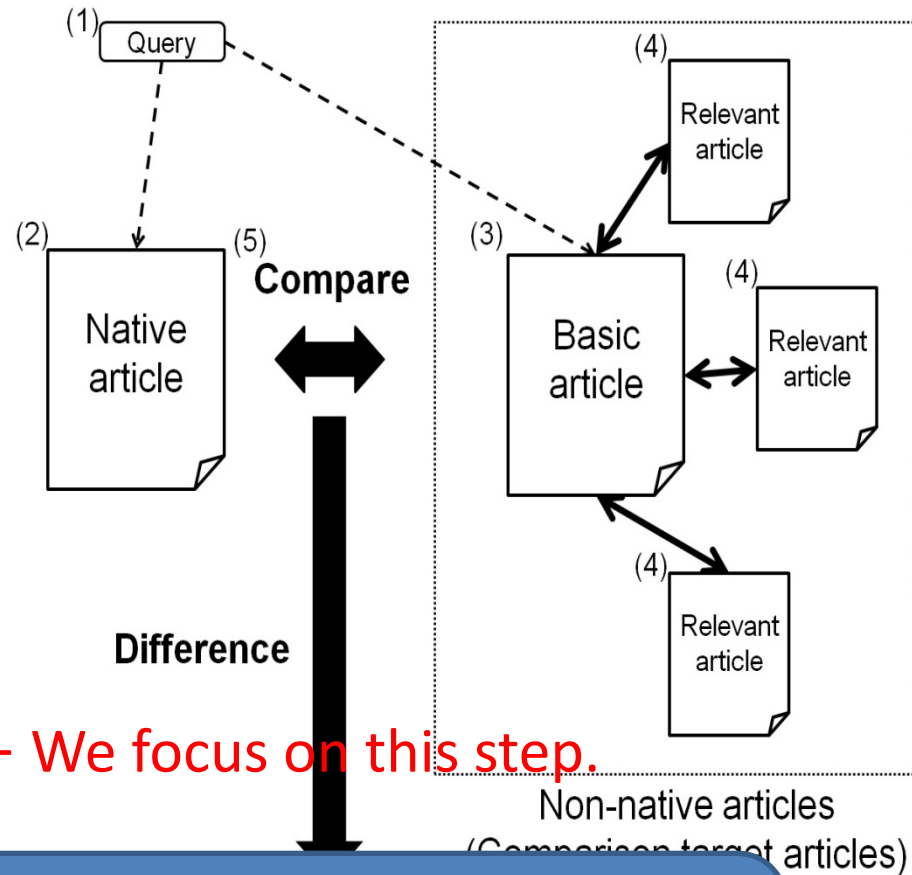
Our Flow

1. Users input a query in their native language to the system.
2. The system retrieves one native article of which title is the same as the user's input query.
3. It translates the query to the non-native language using a language dictionary and retrieves a non-native article of which the title is the same as the user's input query.

4. It extracts comparison articles from the non-native articles using a Wikipedia link graph.

5. It compares native article with

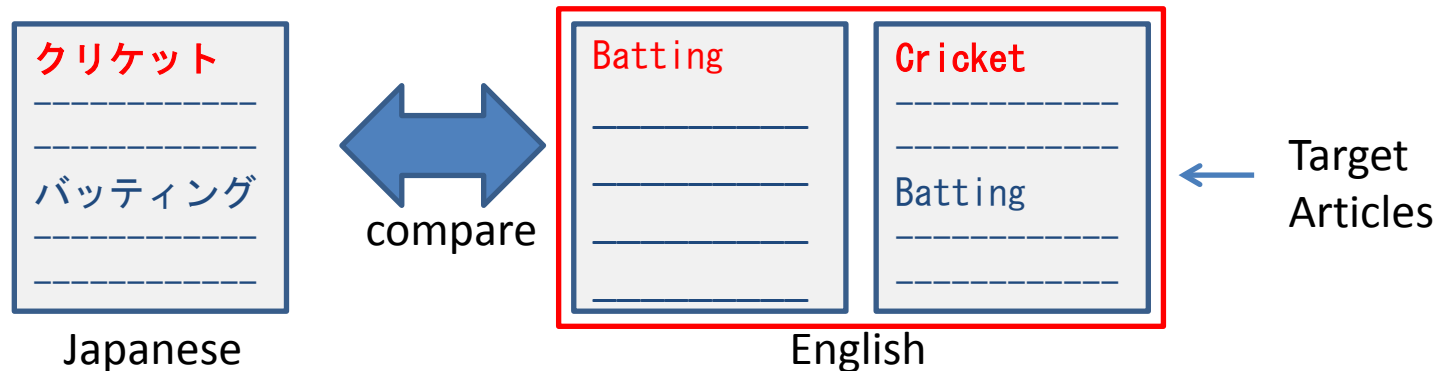
6. It av



Granularity of information differs between the languages in Wikipedia.

– For example:

- Article of "Cricket" is written about Batting of Cricket both Japanese Wikipedia and English Wikipedia.
- In English Wikipedia, there are another page about detail of Batting of Cricket.



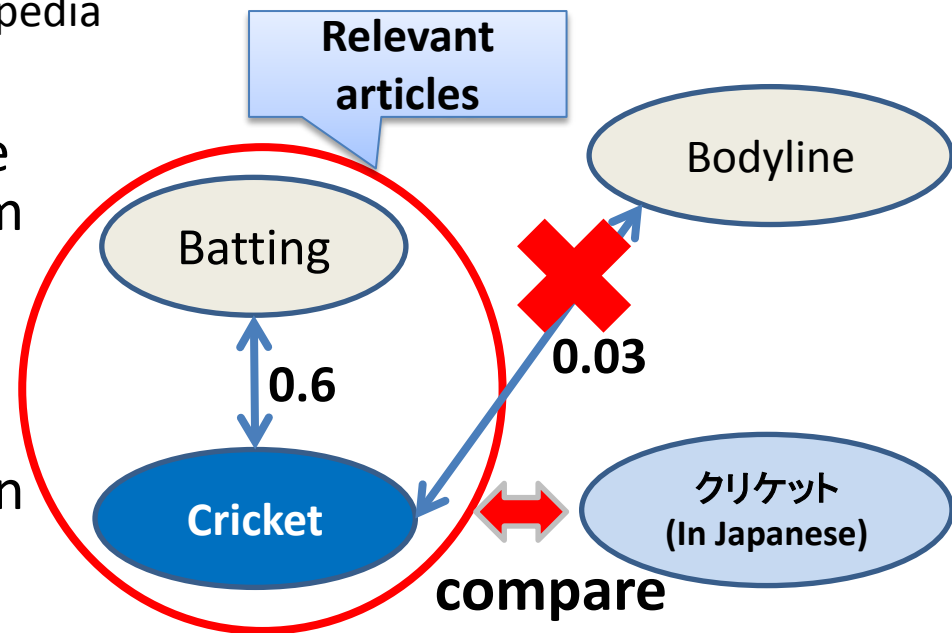
When we compare a native article with non-native articles, we have to consider multiple comparison non-native articles.



We extract target articles based on the Wikipedia link graph and our proposed relevance degree.

Extract comparison target articles

We create a link graph for non-native Wikipedia based on the user's input query.

1. We extract articles having the same title as the user's input from the comparison Wikipedia.
Basic article → root node
2. We extract all interactive linked articles and they become nodes in link graph.
3. We calculate the relevance degree between the root node and the other nodes in the link graph.
4. When the relevance degree is greater than a threshold β value, then we regard the articles as relevant articles.



:root node(basic article)
:other node

Calculating Relevance Degree

- Extraction of the relevance article using only cosine similarity between root node and the other nodes.
 - ⇒ The result of recall ratio is **not good**
- **Relevance Degree** between root node and the other nodes.

Position of the link anchor



The Important anchor appears in the summery area in Wikipedia.

Number of the link anchor



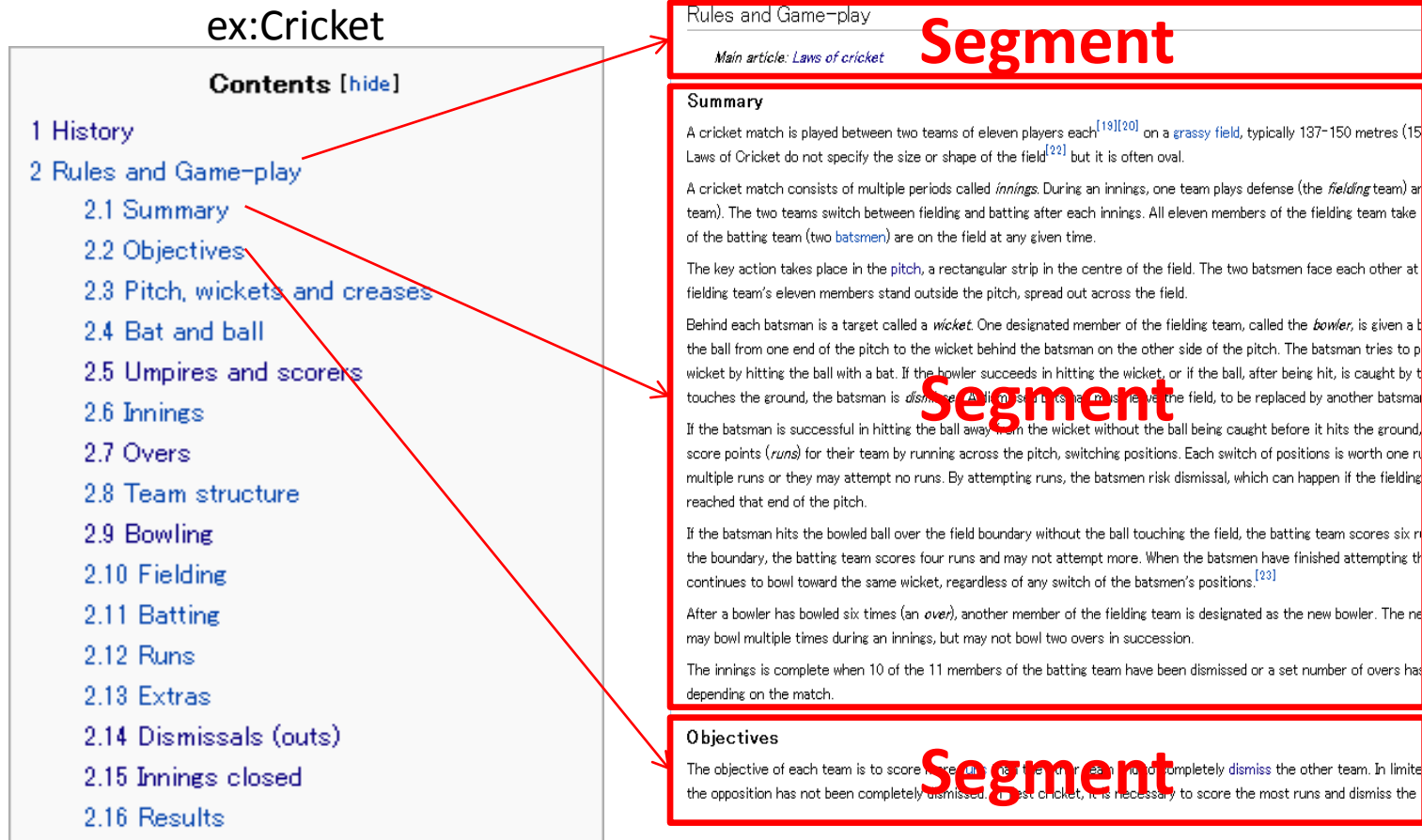
Important anchors related to the basic article appear many times in the basic article.

Similarity between
the content



If articles are similar, relevance degree becomes high.

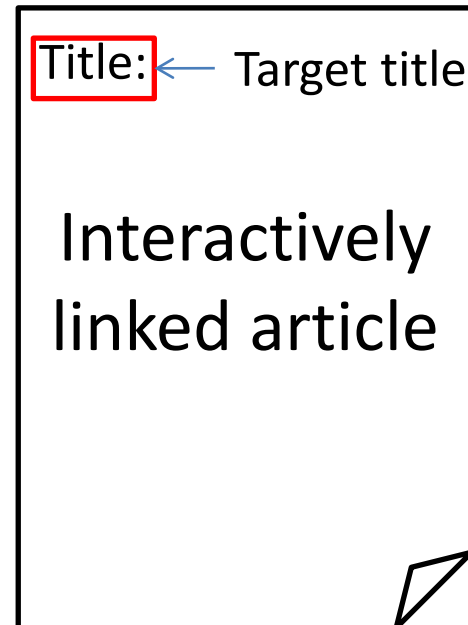
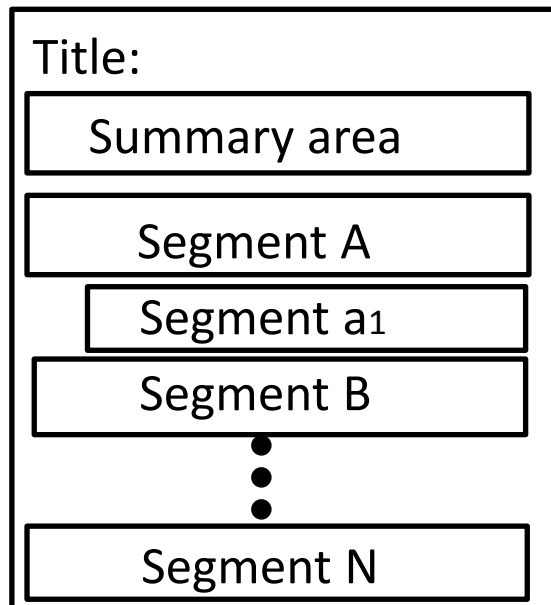
Calculating Relevance Degree



The system divides the basic article according to the structure of the table of contents of the basic article.
The divided parts → segments.

Calculating Relevance Degree

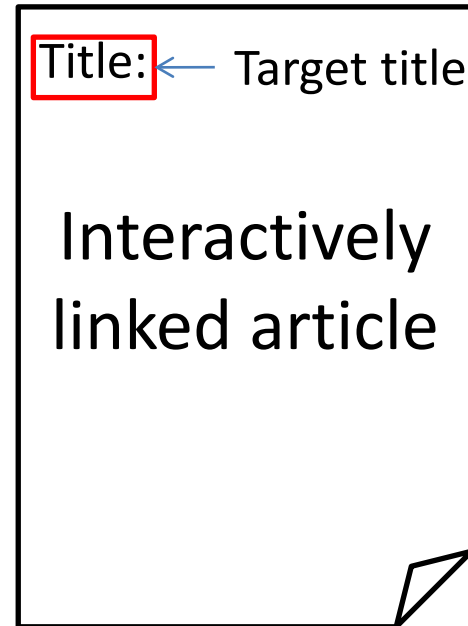
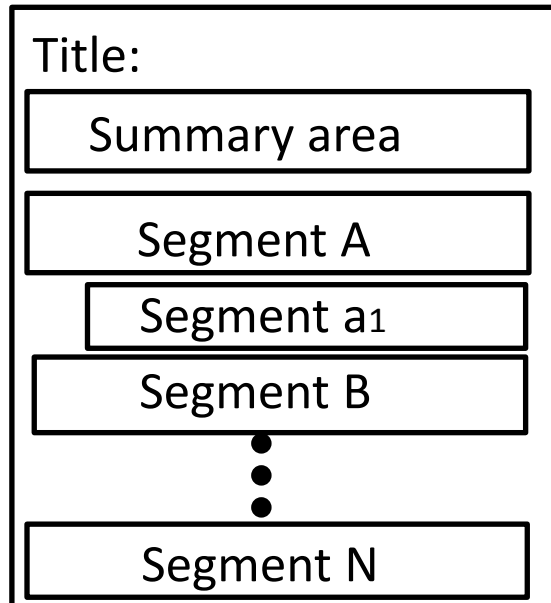
Basic article



- The system extracts a title from a node, which is the interactively linked article.
- The title becomes a keyword used for extraction of the anchor text from the basic article.
- The title → target title.

Calculating Relevance Degree

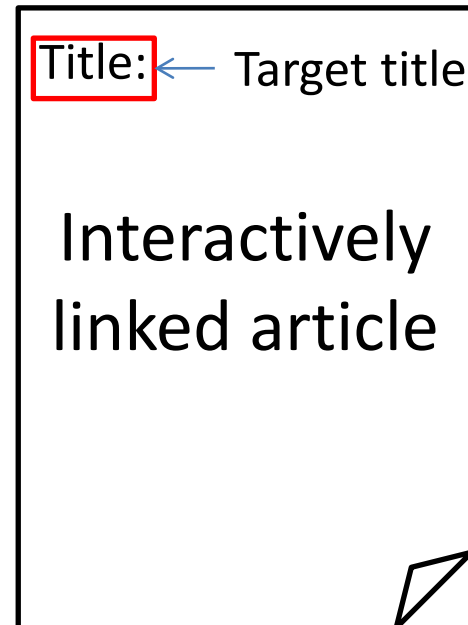
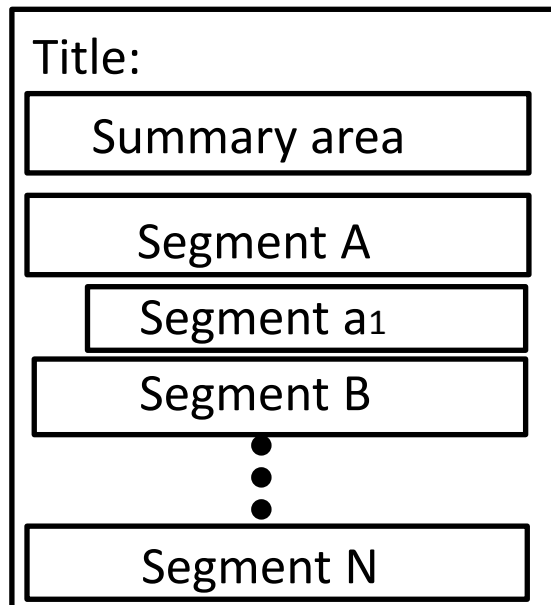
Basic article



- The system counts the anchor text of the target title in the summary area of the basic article.
- It also counts the anchors in each segment of the basic article.

Calculating Relevance Degree

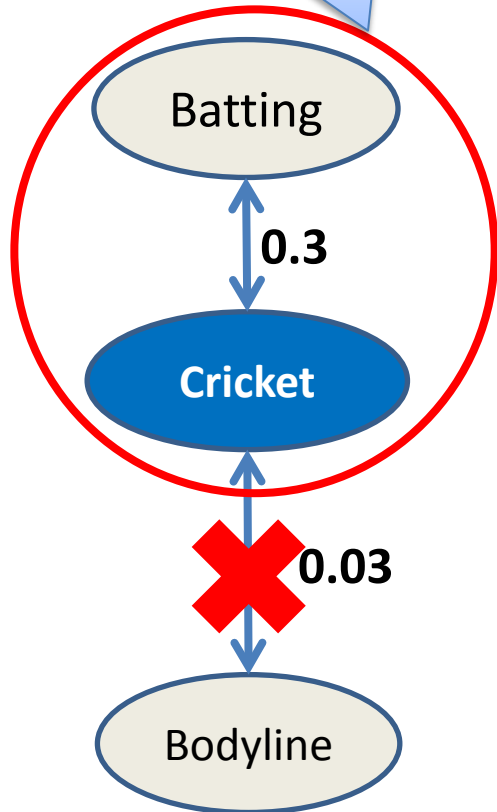
Basic article



- The system calculates the similarity between interactively linked article and the summary area of basic article.
- The system calculates the similarity between interactively linked article and the segment area.

Calculating Relevance Degree

Relevant articles



●:root node
○:other node

Position of the link anchor

Number of the link anchor

Similarity between the content

$$R_i = \{ \alpha \cdot (TF_{sum_i} \cdot S_{sum_i}) + \sum_{k=1}^n (TF_{ik} \cdot S_{ik}) \} / \max(R_{im})$$

i : the identification number of the interactively linked article

R_i : Relevance Degree of article i

TF_{sum_i} : number of the anchor in summary area

S_{sum_i} : the similarity between i and the summary area

TF_{ik} : number of the anchor in the segment k

S_{ik} : the similarity between i and the segment k

K : the segment number

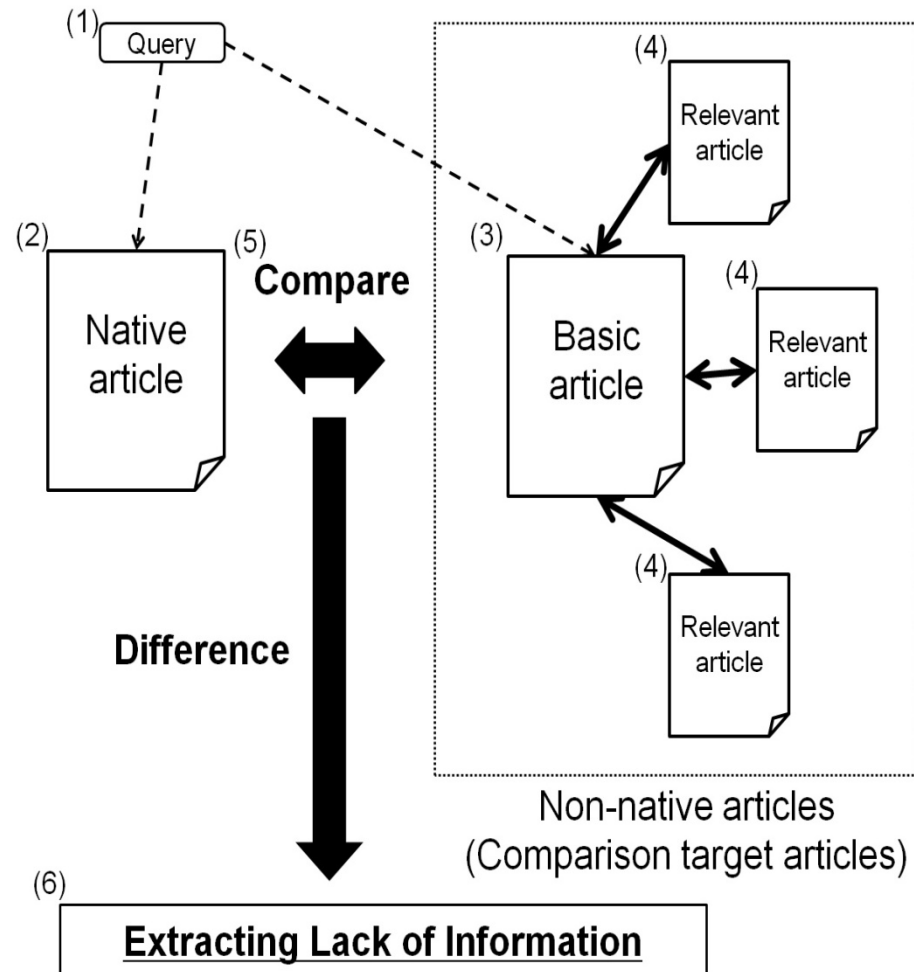
N : the number of segments in the basic article

$\max(R_{im})$: the maximum value in all R_i

$\alpha \rightarrow 3.0$ $\beta \rightarrow 0.2$

Our Flow

1. Users input a query in their native language to the system.
2. The system retrieves one native article of which title is the same as the user's input query.
3. It translates the query to the non-native language using a language dictionary and retrieves a non-native article of which the title is the same as the user's input query.
4. It extracts comparison articles from the non-native articles using a Wikipedia link graph.
5. **It compares a native article with non-native articles extracted in 4. and extracts lack of information.**
6. It browses lack of information available on the web.



Comparison between native article and non-native articles

- Almost all Wikipedia articles are divided into segments based on the table of contents meaning that the segments are divided semantically.
- When comparing the similarity of multilingual Wikipedia, we examine the segment of the table of contents of Wikipedia.
- If the similarity of a content is lower than all content, we extract the content as lack of information.

Ex: Fish and chips

概略 [編集]

タラやカレイ、オヒョウなどの白身魚の切り身に、小麦粉を卵や水または棒状に切った油揚げのチップと合わせて供する。この場合のチで言うフライドポテト(アメリカで言うフレンチフライ)のイギリスでの呼び名の切り身小一切れにジャガイモ中一個分)で450キロカロリー程。

歴史 [編集]

白身魚の切り身を揚げた料理は、少なくとも中世ヨーロッパに存在してロップ各地でジャガイモを揚げた料理も作られるようになった。両者はなつたかは諸説あり争われている。記録に残る限りでは、1860年にロンドンが最古のものである。19世紀後半に広がり網魚の技術革新が起こりチップスは労働者階級の日常食になった。第二次世界大戦下のイギリス、フィッシュ・アンド・チップスであった。戦後もフィッシュ・アンド・チップス

食べ方 [編集]

モルトビネガー(麦芽を原料とする穀物酢)と食塩をかけてマッシュイビー一般的だが、ヨーグルト、ケチャップ、マスタード、ソース、マヨネーズなど好みによ、多様な味付けが行われてよい。次の店内では皿に芋のように、紙袋に入れて、円錐型に丸めた新聞紙に包まれて渡される店もある。ファストフード店では、フィッシュをパンズに挟み、チップス

compare

History

Main article: British cuisine

Fish and chips became a stock meal among the working classes in Great Britain as a cities during the second half of the 19th century.^[2] In 1860, the first fish and chip sh Deep-fried chips (slices or pieces of potato) as a dish may have first appeared in Brita earliest usage of "chips" in this sense the mention in Dickens' *A Tale of Two Cities* ("drops of oil"). (Note that Belgian tradition as recorded in a manuscript of 1781, dates 1699.)^[4]

Lack of information

England

The dish became popular in wider circles (Charles Dickens mentions a "fried fish w England a trade in deep-fried chipped pote Tommyfield Market.^[3] It remains unclean and-chip shop industry we know today, dc London in 1860 or in 1865, while a Mr Le

Segment

Segment

Segment

Segment



Experiment 1

- We confirmed the availability of extracting relevant articles in non-native articles.
 - We compare our method with baseline.
 - The baseline is the cosine similarity.
 - using precision, recall, and F-measure by comparing our proposed method with the baseline.

Result of Experiment 1

Query	#	Baseline			Proposed		
		Precision(%)	Recall(%)	F-measure	Precision(%)	Recall(%)	F-measure
Bannock (food)	2	0	0	0	20	50	28
Warwick Castle	2	15	100	27	25	100	40
Black dog (ghost)	7	67	29	40	100	29	44
Fish and chips	4	40	50	44	50	75	60
Goodwood Festival of Speed	2	0	0	0	50	50	50
Bowls	2	33	100	50	10	50	14
Burlesque	3	60	50	55	100	67	80
Flag of Scotland	6	50	50	50	67	33	44
Gaelic handball	4	25	25	25	80	100	89
Kipper	3	67	67	67	100	67	80
National Gallery of Scotland	12	72	67	70	75	50	60
Lipton	1	0	0	0	25	100	40
Average	-	37	45	36	59	64	52

#:Number of correct results



Experiment 2

- We confirmed the accuracy of extracting lack of information.
 - We use English Wikipedia as native article and Japanese Wikipedia as non-native articles.
 - The correct answer was judged by a bilingual person.

Result of Experiment 2

Query	#	Precision(%)	Recall(%)	F-measure
Bannock (food)	2	33	50	40
Warwick Castle	12	79	92	85
Black dog (ghost)	32	89	78	83
Fish and chips	11	45	82	58
Goodwood Festival of Speed	10	60	60	60
Bowls	9	50	100	67
Burlesque	22	71	45	56
Flag of Scotland	56	98	88	92
Gaelic handball	16	68	94	79
Kipper	16	88	94	91
National Gallery of Scotland	4	57	100	72
Lipton	8	71	63	67
Average	–	67	79	71

#:Number of correct results



Discussion of Experiment 2

- They are almost good result.
- Bad result case
 - When we target on tea brand of “Lipton”, we extract “Thomas Lipton” as a relevant article. He created the Lipton tea brand. He is also famous for sportsman. It is not related to the tea brand of “Lipton”. But we extract it as a lack of information for the tea brand of “Lipton”.
- Other case is attributable to a translation problem.



Conclusion and Future work

- We proposed a method for extracting information that exists in one language version, but which does not exist in another language version.
- Two points
 - Examine the link graph of Wikipedia and structure of an article of Wikipedia.
 - Extract comparison target articles of Wikipedia using our proposed degree of relevance.
 - Compare between native article and non-native articles.
- Future work
 - Considering word sense disambiguation.
 - Comparing other languages (ex. Chinese, Korean etc...)