# Extracting Difference Information from Multilingual Wikipedia

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

# Background 1

- Wikipedia, a large encyclopedia that is accessible using the Internet.

- Wikipedia has two characteristics.
  - Anyone can edit the content .
  - There are 250 over language versions.

- The content of articles differ among about respective language version.

**Native language version should have difference information.**

# 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.
In this way , **Native language version has difference information**.

Query: Fish and Chips

Japanese version

目次

1 概略
2 歴史
3 食べ方
4 関連項目

poor

English version

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

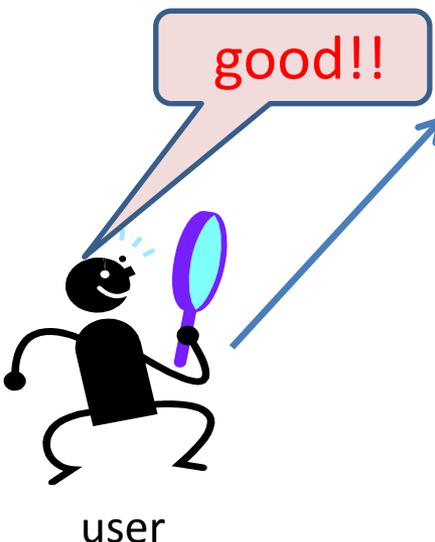rich

Japanese user

# Background 3

- Generally, users browse the Wikipedia of their own native language.

- Occasionally they would refer to the other language versions.

- It is difficult for them to understand the whole content which is written by non-native language.

It may be able to understand a passage of non-native language Wikipedia.

# Propose

- We consider that  if there are difference information, we add it in browsing language version.



**good!!**

user

**Difference information**

**add**

Extracting Difference information from Multilingual Wikipedia

# Naive Method

1. User inputs keyword in Japanese to the system.

2. The system retrieves one article related to keyword from the Japanese version of Wikipedia.

3. The system extracts the English article by using interlanguage link in Wikipedia.

4. The system compares Japanese article with English comparison target articles.

5. The system outputs Japanese article with sections of English articles that do not appear in the Japanese article.

# Naive Method

1. User inputs keyword in Japanese to the system.

2. The system retrieves one article related to keyword from the Japanese version of Wikipedia.

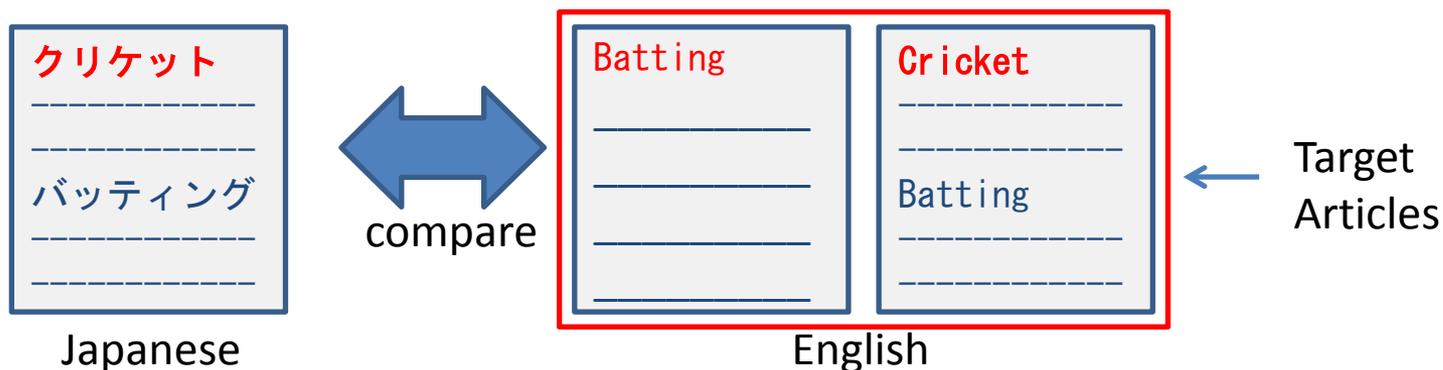3. The system extracts the English article by using interlanguage link in Wikipedia.← We focus on this step.

4. The system compares Japanese article with English

5. article.

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.

| クリケット |
| --------- |
| ---------- |
| ---------- |
| バッティング |
| ---------- |
| ---------- |

compare

| Batting | Cricket |
| ------- | --------- |
| _____ | ---------- |
| _____ | ---------- |
| _____ | Batting |
| _____ | ---------- |
| _____ | ---------- |

← Target Articles

Japanese                English

**When we compare a Japanese Wikipedia with English Wikipedia, we have to consider multiple comparison English articles we call these articles "Target Article".**

# Our Flow

1.  Users input keywords in Japanese to the system.

2.  The system retrieves one article related to keywords from the Japanese version of Wikipedia.

3.  The system extracts the multiple English comparison target articles.

4.  The system compares sections in Japanese article and those in English articles, and detects which sections appeared in English articles but not in Japanese articles.

5.  The system outputs Japanese article with sections of English articles that do not appear in the Japanese article.

# Our Flow

1. Users input keywords in Japanese to the system.

2. The system retrieves one article related to keywords from the Japanese version of Wikipedia.

3. The system extracts the multiple English comparison target articles.
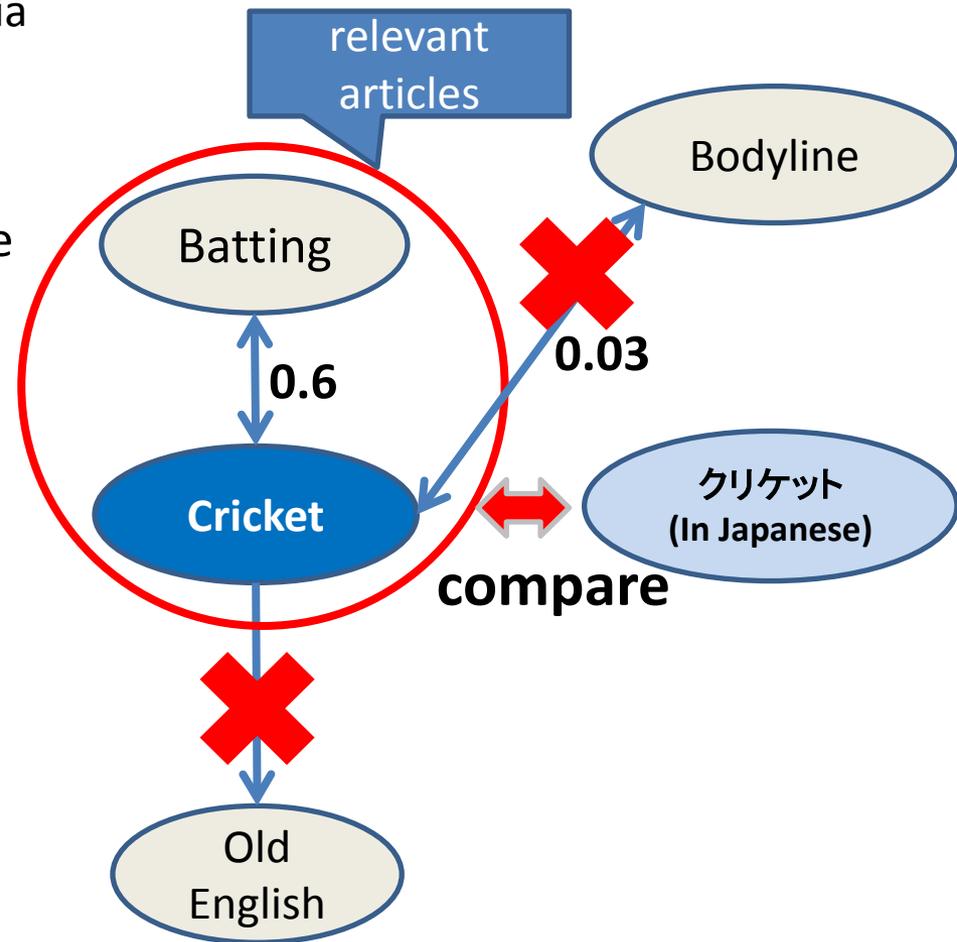
4. The system compares sections in Japanese article and ~~sections in English article~~ and ~~extracts sections~~

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

English articles that do not appear in the Japanese article.

# Extract comparison target articles

We create a link graph for English Wikipedia based on the user's input keyword.

1. The system extracts English articles that have the same title as the user's input and translated. We designate the English article as the basic article. We regard the basic article as the root node of the link graph.

2. The system extracts all interactive linked articles that are the subjects of link-out and link-in connections with the root node. Then it includes the articles as nodes and there by creates the link graph.

3. The system calculates the relevance degree between the root node and respective articles in the link graph.

4. When the relevance degree is greater than a threshold α value, then the system regards the articles as relevant articles.

relevant articles

Bodyline

Batting

0.6

0.03

Cricket

compare

クリケット
(In Japanese)

Old English

:root node

# Calculating Relevance Degree

- We proposed extraction of the relevance article using only cosine similarity between root node and the other nodes.

    ⇒The result of recall ratio is not good

- We propose Relevance Degree between root node and the other nodes.

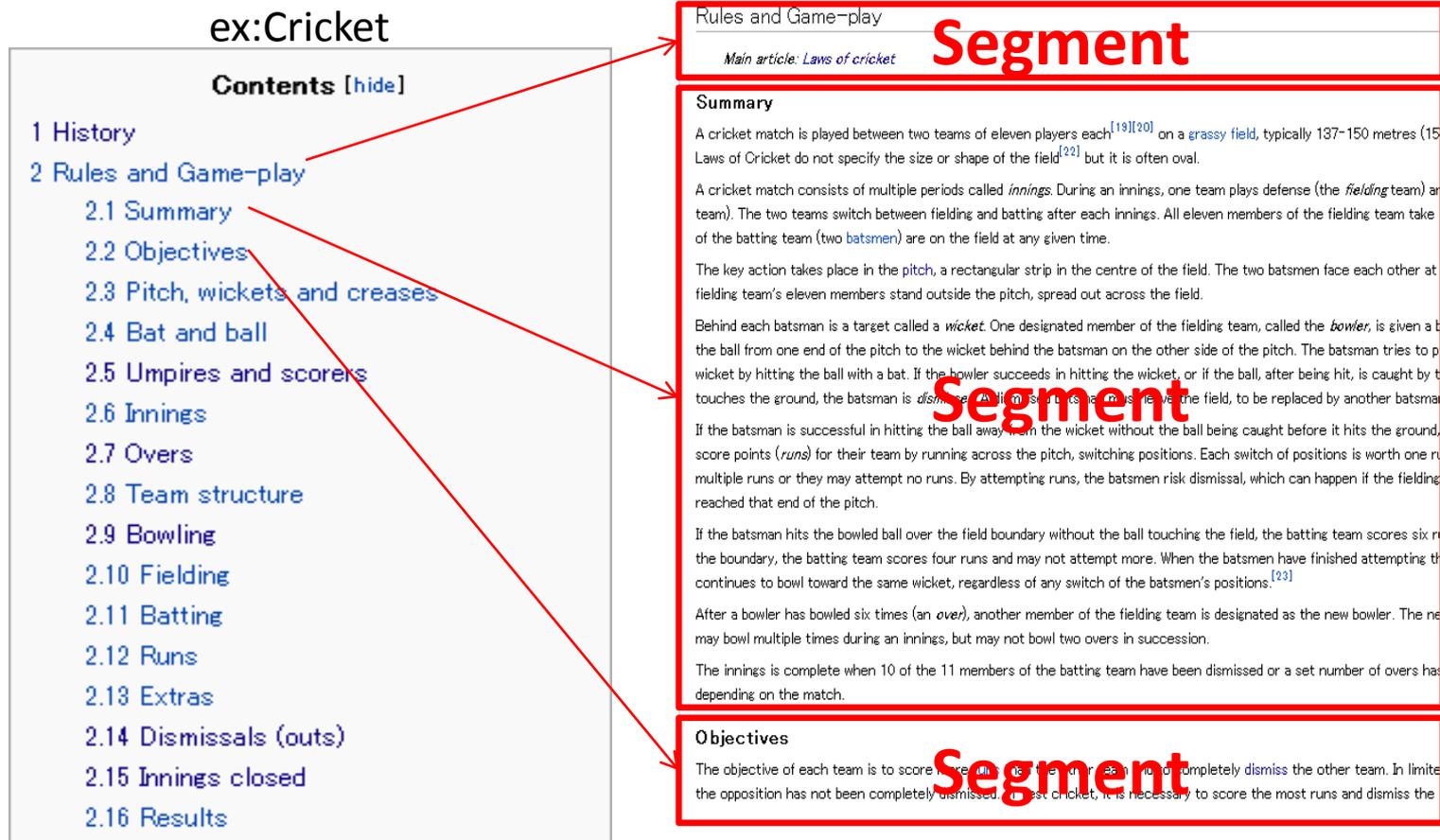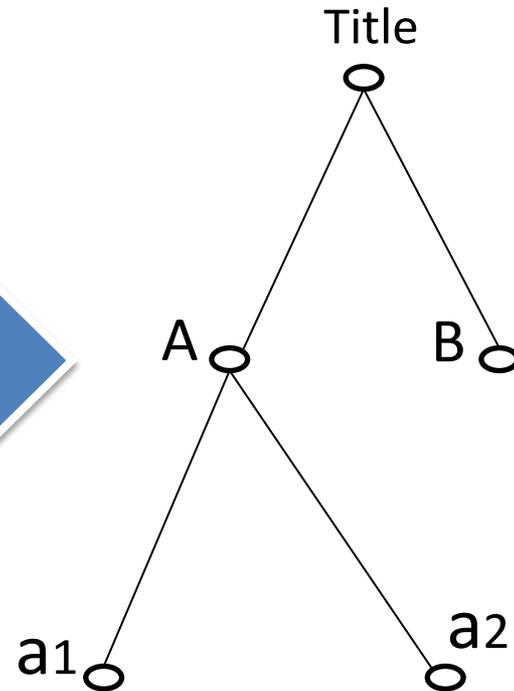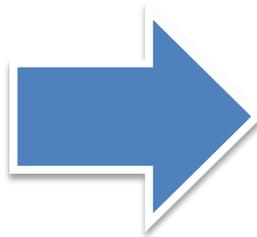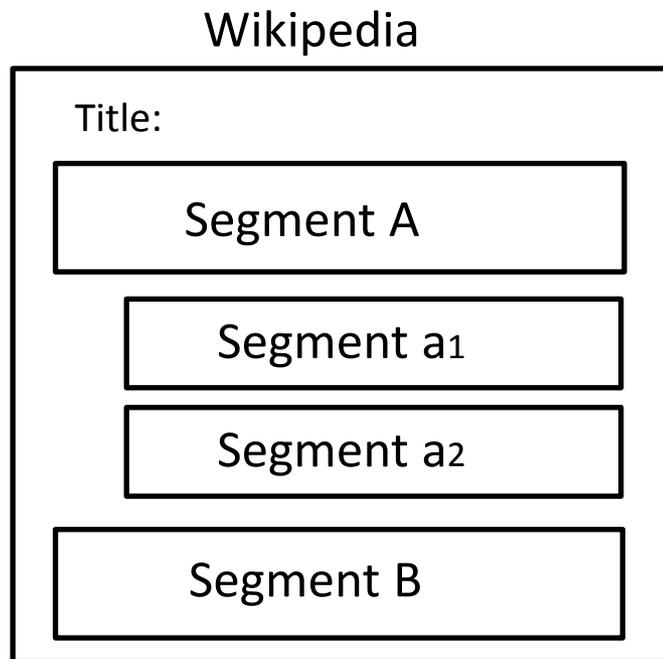| Position of the link anchor | → | The Important anchor appears in the upper area of the page and also appears in the section area more than in a subsection area. |
|---|---|---|
| Number of the link anchor | → | Important anchors related to the basic article appear many times in the basic article. |
| Similarity between the content | → | If articles are similar, relevance degree becomes high. |

12

# Calculating Relevance Degree

ex:Cricket



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

# Calculating Relevance Degree

Wikipedia

Title:

Segment A

Segment a₁

Segment a₂

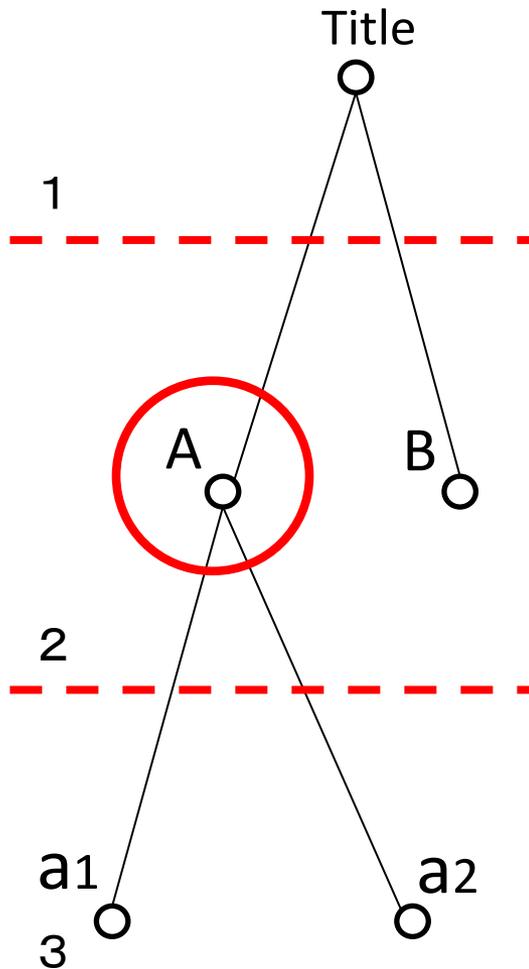Segment B

Title

A        B

a₁        a₂

The system creates segment tree for which the root node is the title name of the basic article, child nodes are segments.

In the segment tree, the child node of the root node is a section of the basic article; the grandchild node is a subsection of the parent node.

The left side nodes show younger segment numbers.
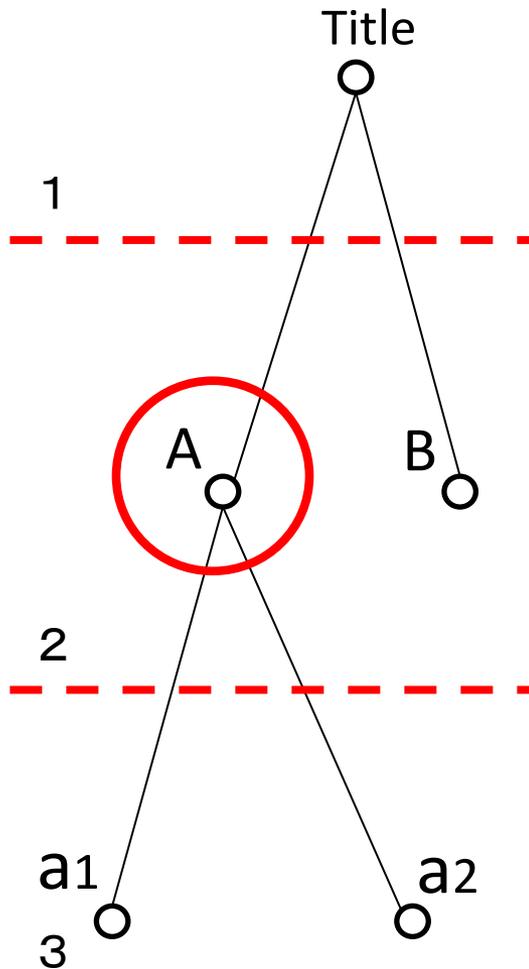
14

# Calculating Relevance Degree



- Our hypothesis is that important anchors related to the basic article appear many times in the basic article.

- They appear in upper areas of the basic article, and also appear in the section area more than subsection area.

- We infer that the relevance degree of the nodes of the low hierarchy and left side in the segment tree are higher than the nodes of the deep hierarchy and right side in the segment tree.

# Calculating Relevance Degree

Title

1 - - - - - - - - - - - - - - - -

A        B

2 - - - - - - - - - - - - - - - -

a1        a2

3

Position of the link anchor

Number of the link anchor

Similarity between the content

$$W_{kl} = af \times S_{kl} + \sum_{i=1}^{af} \left\{ \left(\frac{1}{d_i}\right)^{n_i} \times (n_i - o_i + 1) \right\} / \max(W_{km})$$

$W_{kl}$:Relevance Degree of article k and l

af:The number of the anchor

$d_i$:depth

$n_i$:sibling node of i

$o_i$:sequence number of l in the same depth

$S_{kl}$:The content is similarity between article of k and l

16

# Our Flow

1. Users input keywords in Japanese to the system.

2. The system retrieves one article related to keywords from the Japanese version of Wikipedia.

3. The system extracts the multiple English comparison target articles.

4. The system compares sections in Japanese article and those in English articles, and detects which sections appeared in English articles but not in Japanese articles.

5. The system outputs Japanese article with sections of English articles that do not appear in the Japanese article.

# Comparison between Japanese article and English 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 to all content, we extract the content as difference information.

Ex:Fish and chips



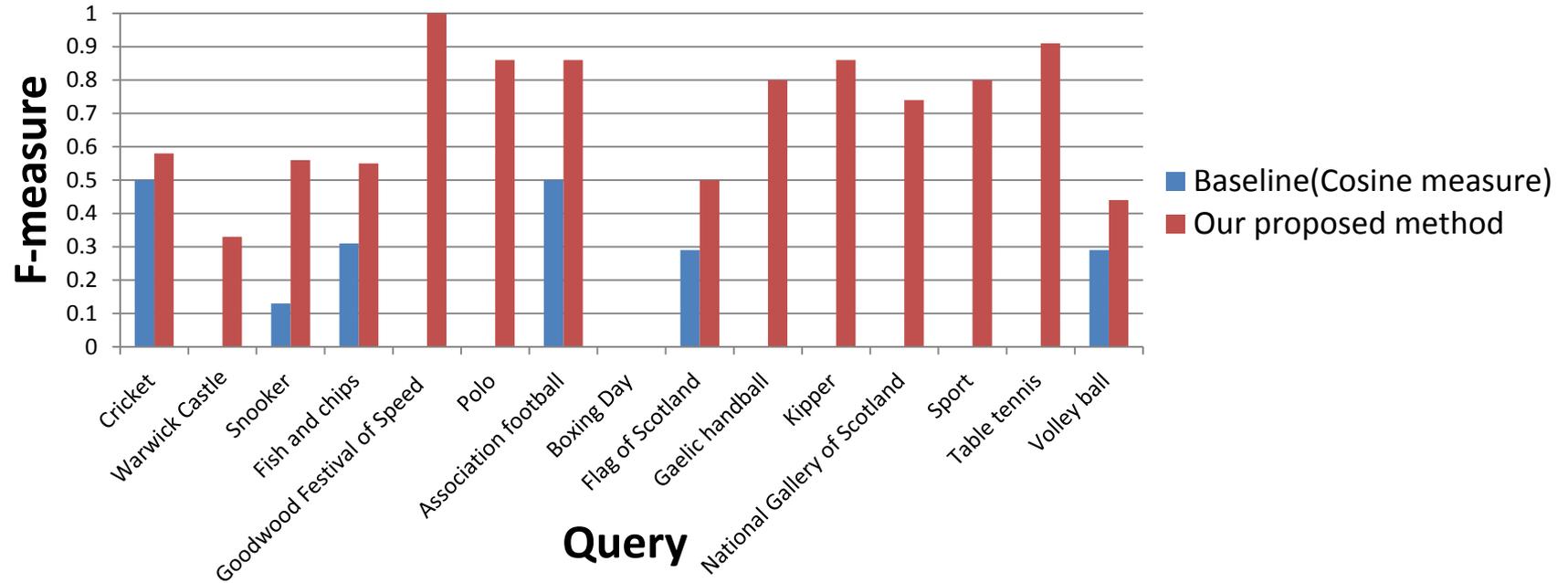**compare**

**Difference information**

# Experiment

- We confirmed the accuracy of difference information extraction methods.
  - Using precision, recall and F-measure by comparison of our proposed method with the baseline, which is only the cosine similarity between basic articles and other linked articles.
  - Compared Wikipedia
    - We used Japanese and English Wikipedia for experiment
  - We set the threshold of relevance degree to 0.2 based on the result of our pre-experiment.

# Experiment



When we compare the results of baseline with our proposed method,
we can observe that the average of precision ratio improves from 0.37 to 0.68,
and that the average of recall ratio improves from 0.1 to 0.61.

Moreover, the average of *F-measure improves* from 0.13 to 0.61.

From these improvements, we confirmed that our proposed method can extract appropriate
parts of articles from English Wikipedia articles.

# Result and discussion

- Average of precision ratio: 0.37⇒0.68
- Average of recall ratio: 0.1⇒0.61
- Average of F-measure: 0.13⇒0.61

- Discussion
  - when we use the queries "Association Football" and "Table Tennis," which are generally used terms, our proposed system is more effective than when we use terms for specific fields, because when we use general terms, the relevant terms are numerous. Also, the links to Wikipedia articles are numerous.
  - Therefore, the relevance of correct sections is large.

# Conclusion

- We proposed a method for extracting difference information.


- Two points
  - Examine the link graph of Wikipedia and structure of and article of  Wikipedia.
    - Extract comparison target articles of Wikipedia using our proposed degree of relevance.
  - Compare between Japanese article and English articles.

# Future work

- Calculating credibility of Wikipedia articles.
  - Wikipedia credibility is not always good.
  - For that reason, we must assess Wikipedia credibility in future studies.

- Considering word sense disambiguation.
  - Many instances of word sense disambiguation exist, but we ignored such cases in this time.
  - We intend to consider word sense disambiguation in later investigations.