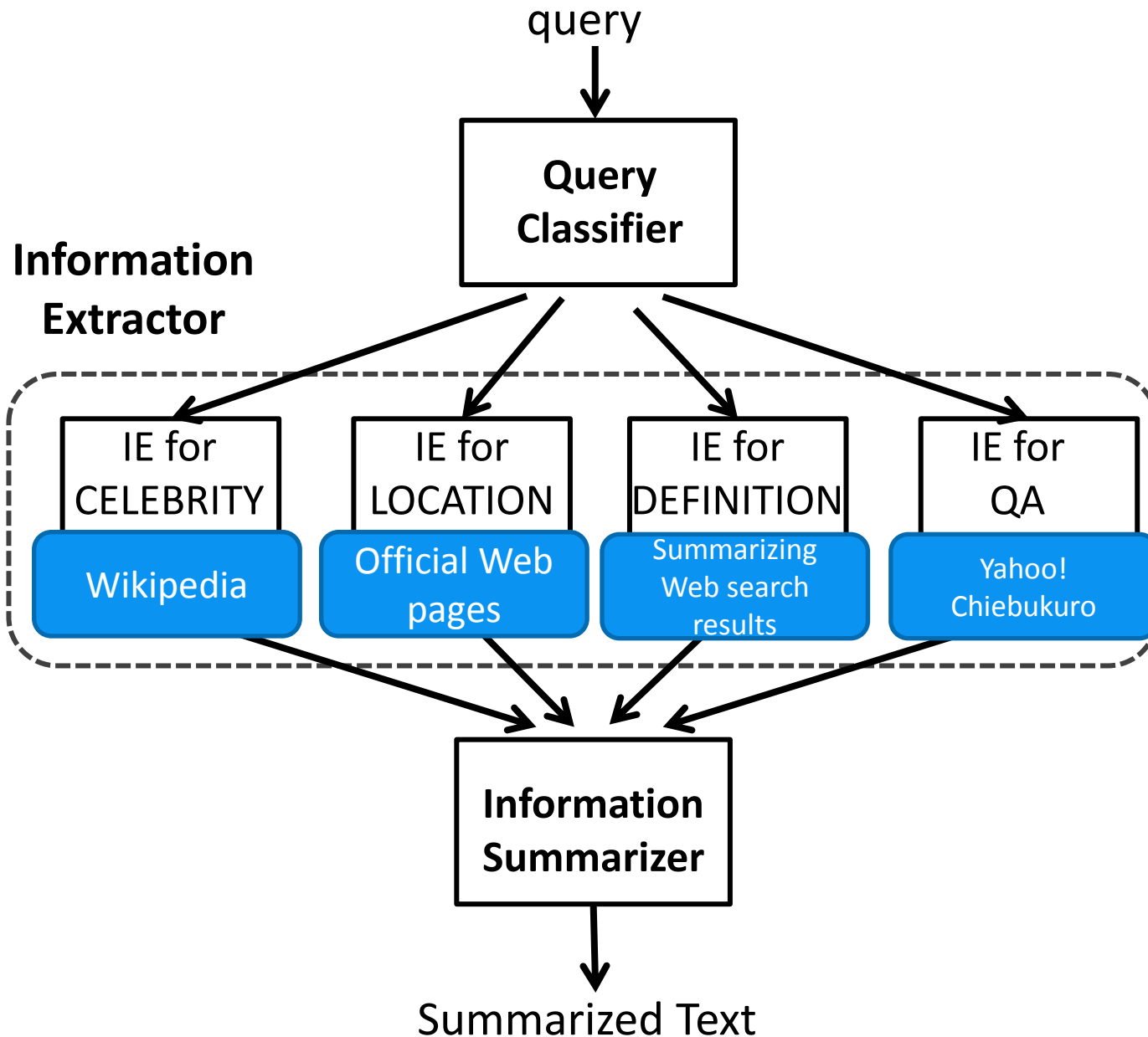# Information Extraction based Approach for the NTCIR-9 1CLICK Task

**Meng Zhao**, Kosetsu Tsukuda, Yoshiyuki Shoji, Makoto P. Kato, Takehiro Yamamoto, Hiroaki Ohshima, Katsumi Tanaka

*Graduate School of Informatics, Kyoto University*

NTCIR-9 Workshop Meeting, 2011/12/9, NII Tokyo, Japan

# Our Framework

query

Query Classifier

**Information Extractor**

IE for CELEBRITY — Wikipedia

IE for LOCATION — Official Web pages

IE for DEFINITION — Summarizing Web search results

IE for QA — Yahoo! Chiebukuro

Information Summarizer

Summarized Text

2

# Query Classifier

Classify four types of query by multi class SVM

185 Features

| Feature | # of features |
|---|---|
| Has Wikipedia article | 1 |
| Frequency of parts-of-speech | 44 |
| Query unigram | 85 |
| Sentence pattern | 2 |
| Number of documents containing expanded query | 15 |
| Has travel services | 1 |
| Number of search results | 1 |
| Terms in search results | 39 |
| **Total** | **185** |

Search with 15 expanded query such as "query-san(さん)" and "query-senshu(選手)" to distinguish name of CELEBRITY.

Count selected 39 terms such as "profile," "born at" and "Chome (丁目)" in snippets of Yahoo Japan search result

# Information Summarizer

- Eliminates overlap from the extracted sentences.

Information Extractor

Because it is **too similar** to the already selected sentences

*Sentence 1*
*Sentence 2*
*Sentence 3*
*...*
*Sentence N*

MMR-based Information Summarization Method

Because of the output length **limitation**

*Sentence 1*
*Sentence 3*
*...*
*Sentence K*

Shorter than the limitation

Output

$$MMR = \operatorname*{argmax}_{d_i \in D/S} [\lambda(Score(d_i)) - (1-\lambda)\max_{d_j \in S}\operatorname{sim}(d_i, d_j)]$$

Score of the sentence

Similarity to the output sentences

# IE for CELEBLITY (Attribute-based)

- Extract pairs of an attribute name and its value from *Infobox* in Wikipedia



attribute name    value

| Country | 🇺🇸 USA |
| Residence | Las Vegas, Nevada, U.S. |

http://en.wikipedia.org/wiki/Andre_Agassi

"Andre Agassi" ➡️ 🖥️ ➡️ "Country is USA, residence is Las Vegas, Nevada, U.S., height is …., and highest ranking is No.1."

# IE for CELEBLITY (Sentence-based)

- Extract important sentences from Wikipedia

- 12 features employed



query: Hayao Miyazaki

(Yahoo! Japan Web search API)

ContainsQuery

Hayao Miyazaki

RelativePosition

The position of a sentence in a section

Summaries(Top n)

MaxCos, MinCos, MaxInnerProd, AvgCos, AvgInnerProd

Calculate cosine, inner-product similarity with each summary

http://en.wikipedia.org/wiki/Hayao_Miyazaki

- Train a regressor to predict the importance score of each sentence (using pre-distributed queries)

# IE for LOCATION

Collect
official Web pages

*query*: **Kyoto Royal Hotel**

**www.royalhotel.com AND ( address OR access OR …)**

(Yahoo! Japan Web search API)

Extract attributes and their values

## Regular Expression based extraction

〒 [0-9]{3} -[0-9]{4}

→ **Address:〒606-8501**

[0-9]{3} -[0-9]{3} -[0-9]{4}

→ **Tel: 075-753-5385**

[a-Z]*@[.a-Z]*

→ **Mail: hotel@kyoto.com**

## Sentence based extraction

10 minutes walk
from Kyoto station

**access information?
(SVM classifier)**

## Table based extraction

| check in | 15:00 |
|----------|-------|
| check out | 10:00 |

Check in  :  15:00
Check out: 10:00

# IE for DEFINITION

- Query : *x*

- Search with a new query "*x towa*（とは）"

- Extract a sentence including the phrase of "*x towa*" from each retrieved Web page

- Apply the *LexRank*[1] algorithm to the set of sentences for estimating the sentence importance

$$\mathbf{p} = \left[ d\mathbf{U} + (1-d)\mathbf{B} \right]^{\mathrm{T}} \mathbf{p}$$

**p** : a vector representing the importance degree of each sentence
**U** : $n \times n$ square matrix whose elements are $1/n$
**B** : adjacency matrix of the cosine similarity between two sentences
*d* : damping factor

[1]G. Erkan and R.D. Radev. LexRank: Graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research*, 22(1), pages 457-479, 2004.

# IE for QA

- Query: Which is taller, Tsutenkaku or Ustunomiya tower?

Calculate *n*-gram similarity

Questions

Yahoo! Chiebukuro
(Yahoo! Answers in Japan)

most similar $q$

$q_1$  $q_2$  $q_3$  $q_4$

$a_1$  $a_2$  $a_3$  $a_4$  $a_5$

Answers

- Extract answers for most similar questions to the given query from <u>Yahoo! Chiebukuro</u>

- Return <u>the best answer</u>
  - if no best answer is given, return the most similar answer to the query

9

- The accuracy of our query classifier: **0.93**

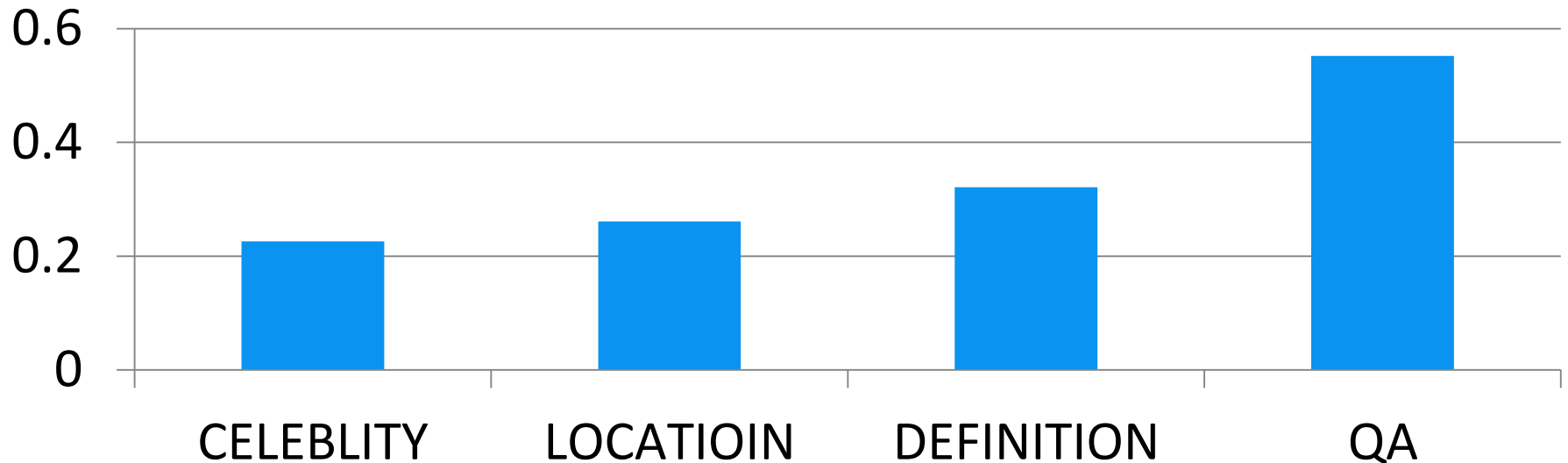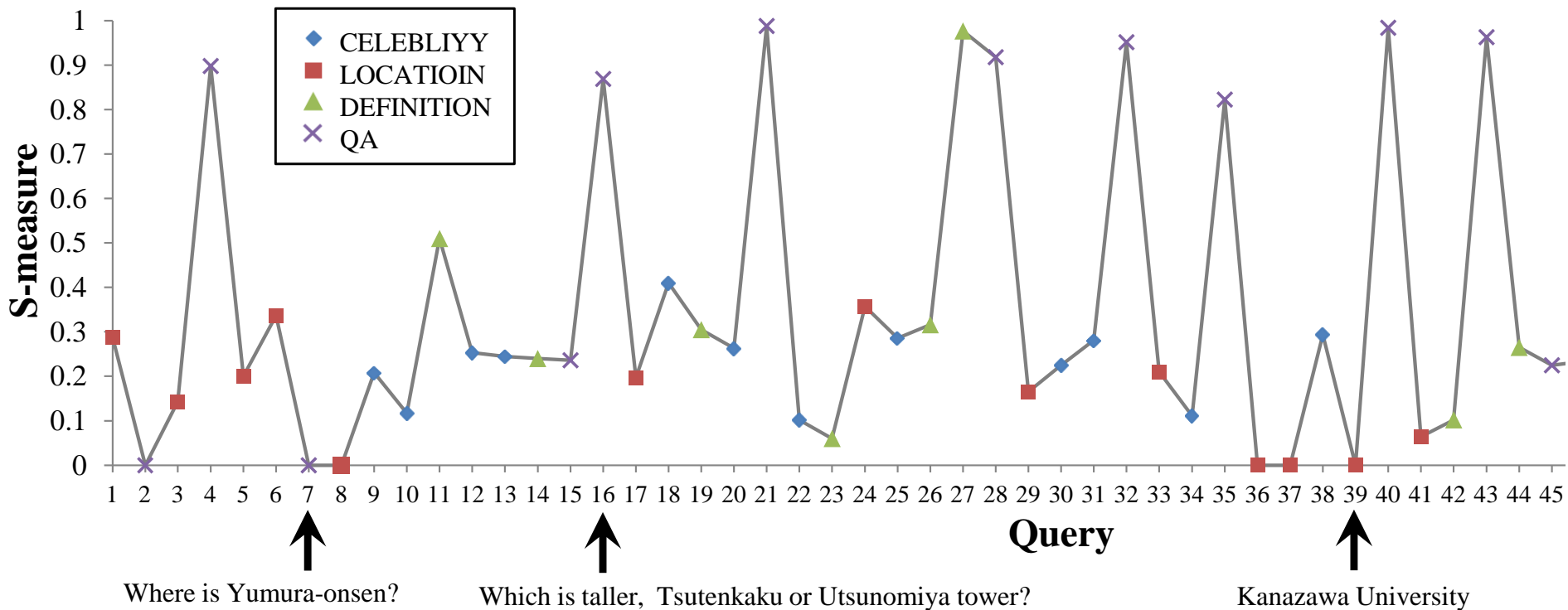- **Highest S-measures for CELEBLITY and QA queries among participants**
  - though the difference is not significant

**S-measure**



10

- For the 1CLICK task, the object identification problem should be tackled (see Kanazawa Univ.)

- QA queries obtained high S-measures when exactly the same questions are available (see the two questions)

11

# Conclusions

- Information extraction framework for 1CLICK
  - Query classifier (**multiclass SVM**)

  - Four types of information extractor
    - CELEBRITY (**attribute and sentence extraction from Wikipedia**)
    - LOCATION (**postal and access information extraction**)
    - DEFINITION (**summarization of Web search results**)
    - QA (**QA pair extraction from Yahoo! Chiebukuro**)

  - Information summarizer (**MMR**)

*Thanks! ntcir-9@dl.kuis.kyoto-u.ac.jp*