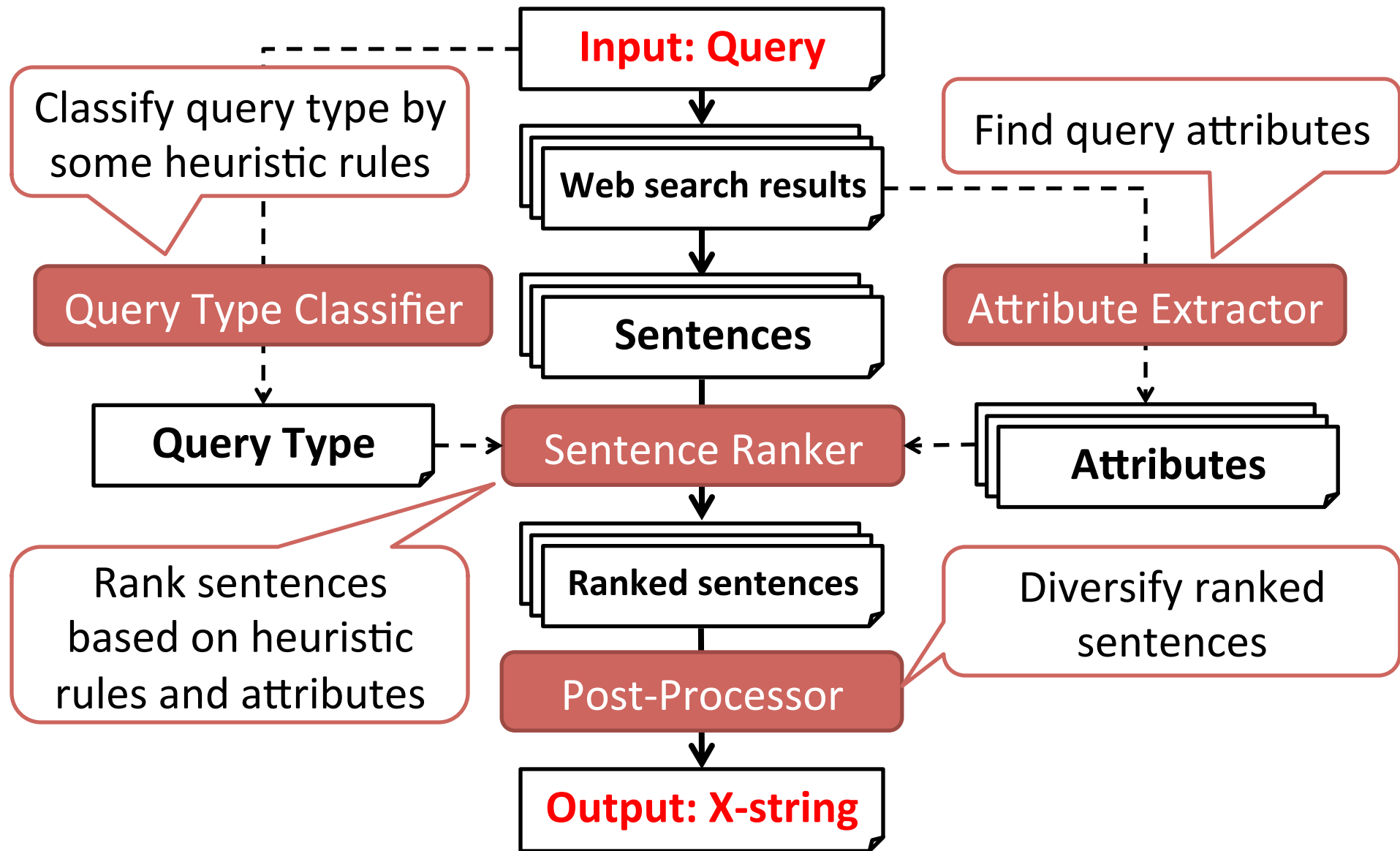


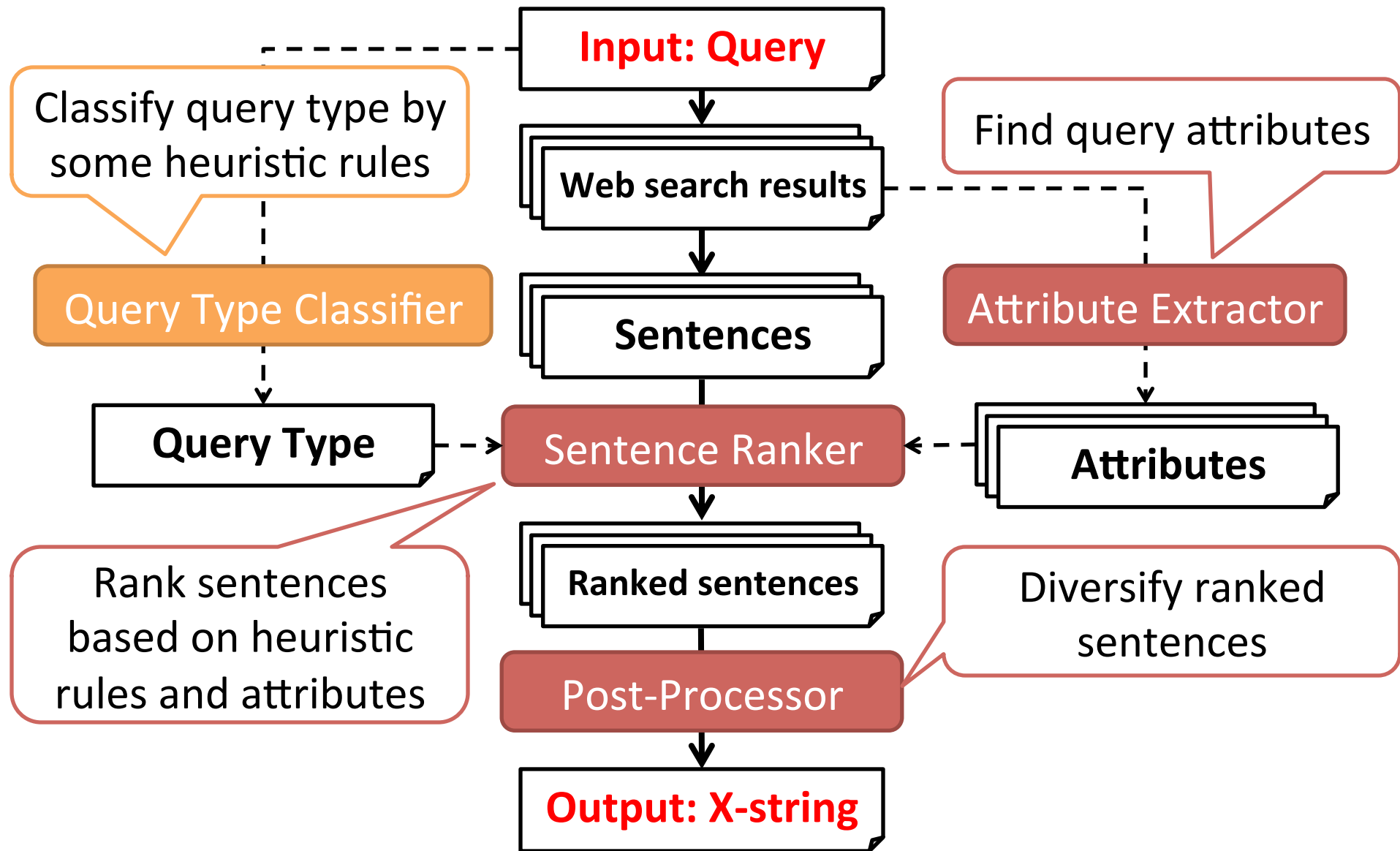
MSRA at NTCIR-10 1CLICK-2

Kazuya Narita (*Tohoku University, Japan*),
Tetsuya Sakai, Zhicheng Dou (*MSRA, China*),
Young-In Song (*NHN Corporation, Korea*)

Overview



Overview



Classification Rules

- Classification is based on some heuristic rules:

– E.g.

Query length

→ QA

“世界で初めてノーベル賞をとったのは誰か”

(Who took the Nobel Prize for the first time in the world?)

Long query may be QA

Classification Rules

- Classification is based on some heuristic rules:

– E.g.

Query length

→ QA

GEO clue suffixes at middle

→ GEO

“博多駅 ホテル”

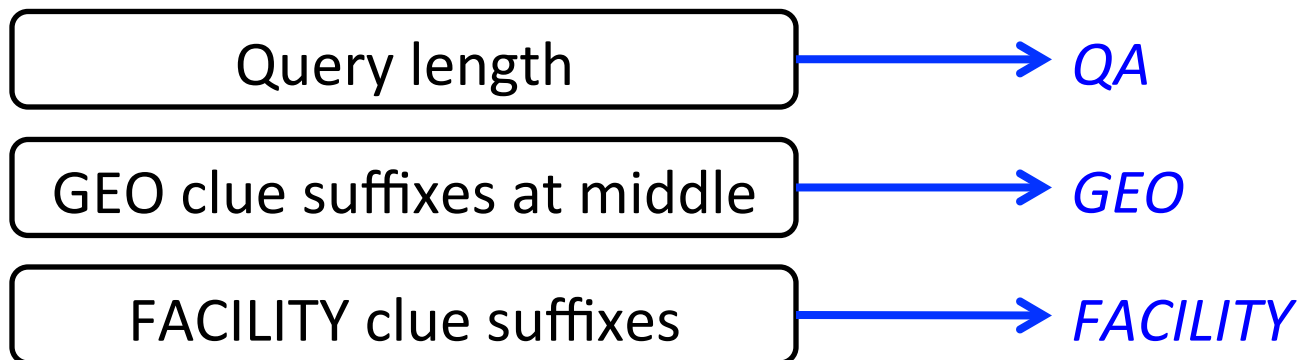
(*hotel near Hakata station*)

Other suffixes: “市” (*city*), “区” (*ward*), “町” (*town*) ...

Classification Rules

- Classification is based on some heuristic rules:

– E.g.



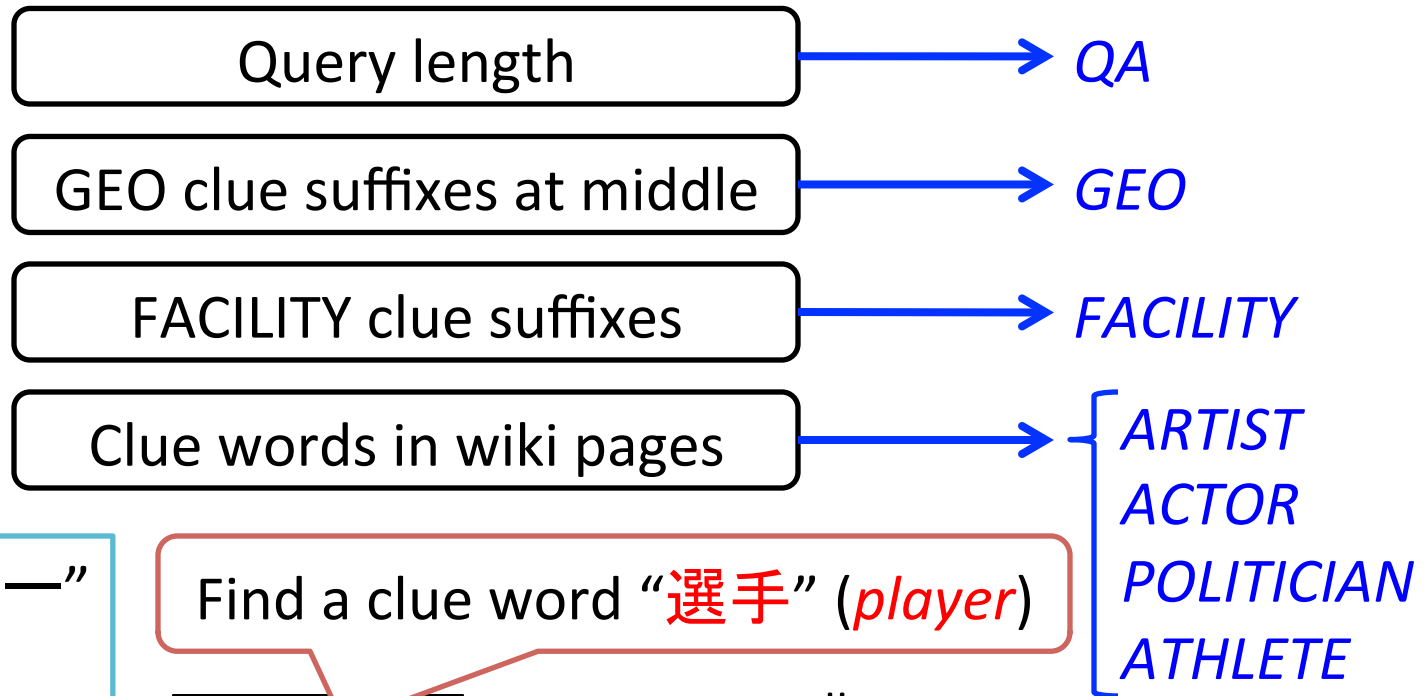
“須磨海浜水族園”
(Suma Aqualife *Park*)

Other suffixes: “学校” (*school*), “病院” (*hospital*) ...

Classification Rules

- Classification is based on some heuristic rules:

– E.g.



“イチロー”
(Ichiro)

Check
a wikipedia page



Find a clue word “選手” (*player*)

||

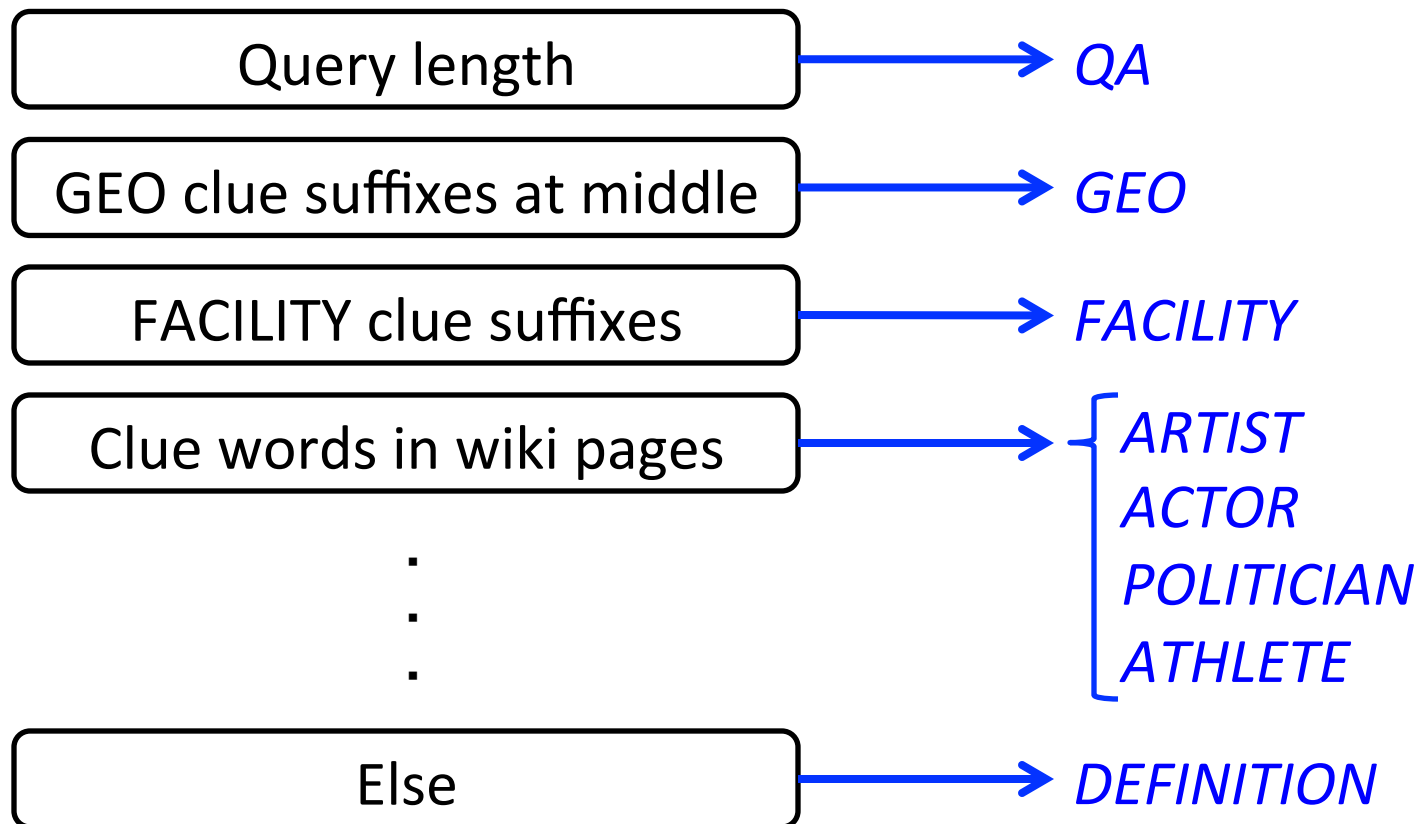
ATHLETE

Other clue words: “小説家” (*novelist*), “女優” (*actress*) ...

Classification Rules

- Classification is based on some heuristic rules:

– E.g.



Results of Query Type Classifier

- Accuracy: 83% (83/100)

Not good...

sys \ gold	ART	ACT	POL	ATH	FAC	GEO	DEF	QA
ARTIST	9	1	1		1	1		
ACTOR	1	9						
POLITICIAN			8					
ATHLETE				10			1	
FACILITY					7		1	
GEO					1	14	1	
DEFINITION			1		5		12	1
QA					1			14

Error Analysis for FACILITY

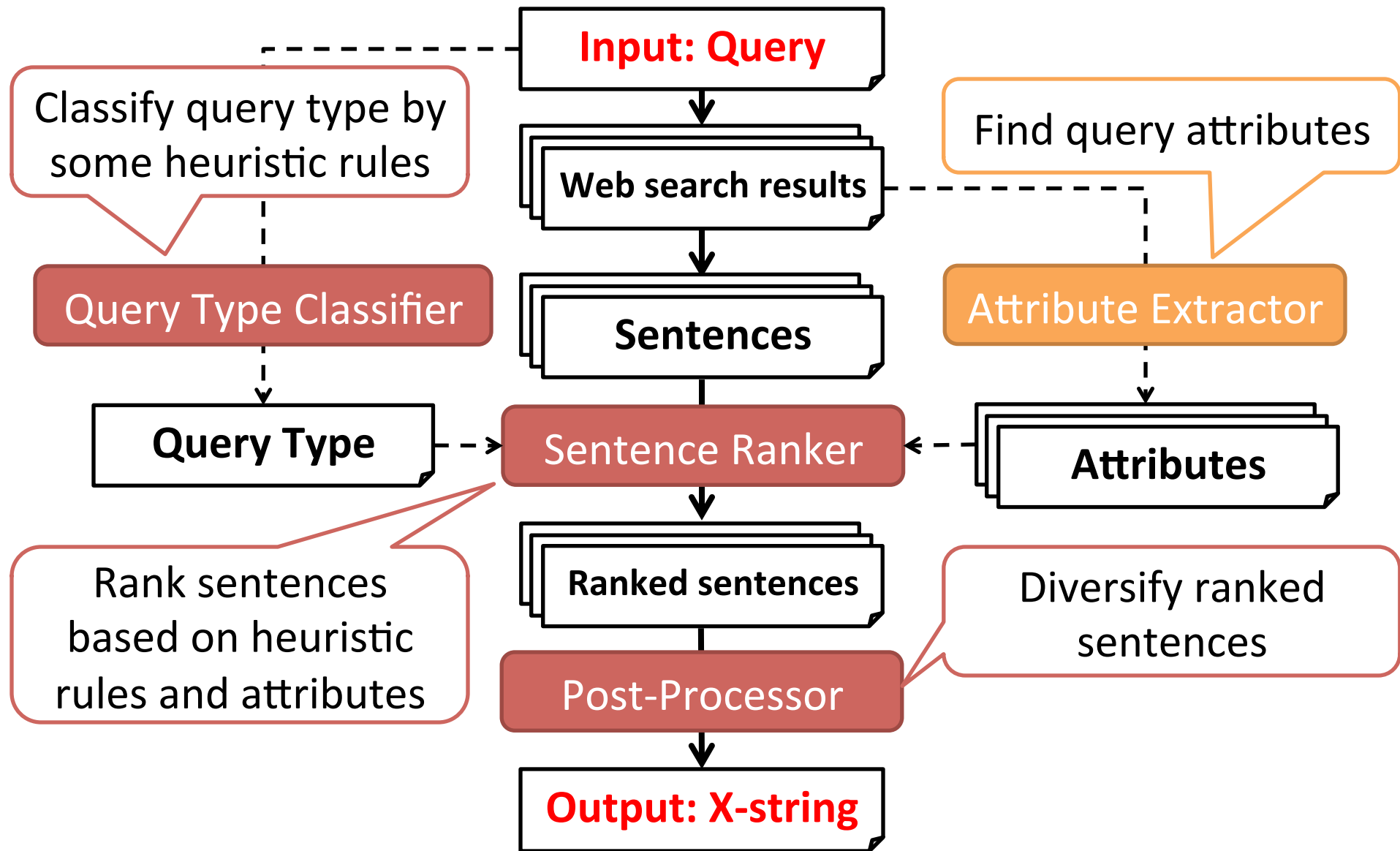
Correct Example

- Correct only by suffix rules
 - 須磨海浜水族園
(Suma Aqualife *Park*)
 - 横浜市役所
(Yokohama City *Hall*)
 - 小金井図書館
(Koganei *library*)
 - ハワイパシフィック大学
(Hawaii Pacific *University*)
 - あおやま矯正歯科医院
(Aoyama Orthodontic *Office*)

Wrong Example

- Insufficient suffixes
 - 京都真如堂
(Kyoto Shinnyo-do *temple*)
 - 南ヶ丘牧場
(Minamigaoka *dairy*)
- Difficult only by suffix rules
 - カーサ・ディ・ナポリ
(Casa di Napoli)
 - ザ・ペニンシュラ東京
(The Peninsula Tokyo)
 - らーめんてつや
(Ramen Tetsuya)

Overview



Query Attributes

Ichiro



Future Science Museum



Born

Height

Position

Address

Tel

Opening Hours

Query Attributes

Ichiro



Future Science Museum



These attributes may be useful when users are seeking for basic information of entities

Born

Height

Position

Address

Tel

Opening Hours

Attribute Extraction

Web search results

The screenshot shows a web browser displaying a Japanese Wikipedia article for Ichiro Suzuki. The page is titled "イチロー" (Ichiro) and is part of the "sportsnavi" website. The article text includes:

この項目では、野球選手について記述しています。その他の用法については「イチロー (曖昧さ回避)」をご覧ください。

イチロー（本名：鈴木 一朗（すずき いちろう）、1973年10月22日 - ）は、シアトル・マリナーズに所属するプロ野球選手（外野手）である。NPB・MLBの双方で活躍し、MLBのシーズン最多安打記録など多数の記録を保持している。夫人は元TBSアナウンサーの福島弓子。

日本での愛称は「イチ」。アメリカ合衆国での愛称は「魔法使い (Wizard)」、「安打製造機 (Hit Machine/Hitting Machine)」。

The article also features a table of contents and a detailed "選手情報" (Player Information) section:

選手情報	
国籍	● 日本
出身地	愛知県春日井市郷龜山町
生年月日	1973年10月22日 (38歳)
身長	5'11" = 約180.3 cm
体重	170 lb = 約77.1 kg
選手情報	
投球・打席	右投右打
ポジション	右翼手
プロ入り	1991年ドラフト4位
初出場	NPB / 1992年7月11日 MLB / 2001年4月2日
年俵	\$18,000,000 (2011年)
経歴 (所属内在職年)	
所属チーム	● 日本
WBC	2006年、2009年

Attribute Extraction

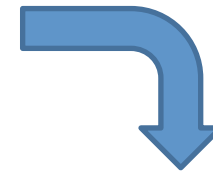
Web search results

The screenshot shows the Japanese Wikipedia page for Ichiro Suzuki. A blue circle highlights a table containing the following information:

国籍	日本
生年月日	1973年10月22日 (38歳)
身長	5'11" = 約180.3 cm
体重	170 lb = 約77.1 kg

A blue arrow points from this table towards the right side of the slide.

Extract tables or text containing “:”



Born	October 22, 1973 (age 38)
Height	5 ft 11 in (1.80m)
Position	Outfielder

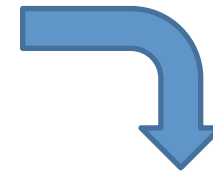
Born: October 22, 1973 (age 38)
 Height: 5 ft 11 in (1.80m)
 Position: Outfielder

Attribute Extraction

Web search results

Extract as query attributes

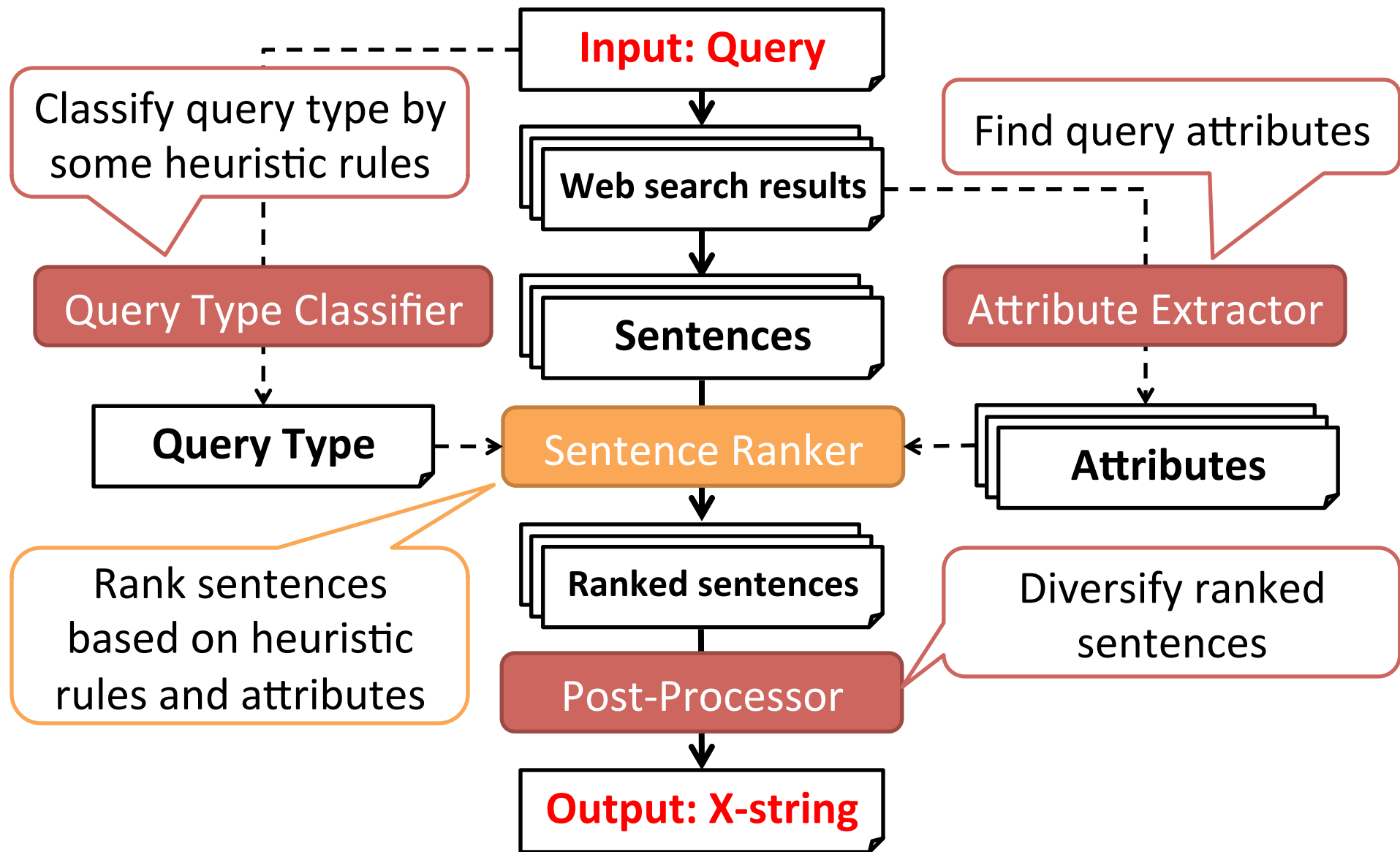
Extract tables or text containing “:”



Born	October 22, 1973 (age 38)
Height	5 ft 11 in (1.80m)
Position	Outfielder

Born: October 22, 1973 (age 38)
Height: 5 ft 11 in (1.80m)
Position: Outfielder

Overview



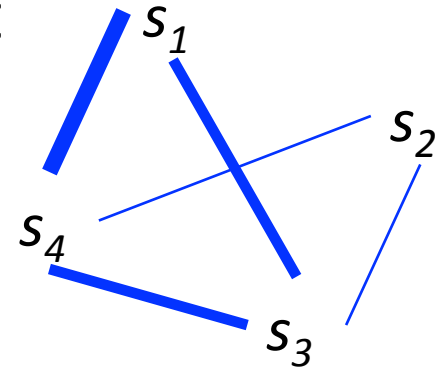
Sentence Ranking

- Scoring based on following components:
 1. Content similarity
 2. Query attributes
 3. Manual clue words and URLs
 4. Search Rank

Scoring

1. Content similarity (*LexRank* [Erkan and Radev 2004])

– Similar sentences are more relevant



Importance of each sentences
(similar to many other sentences -> High)

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

All elements are $1/n$
(n : the number of sentences)

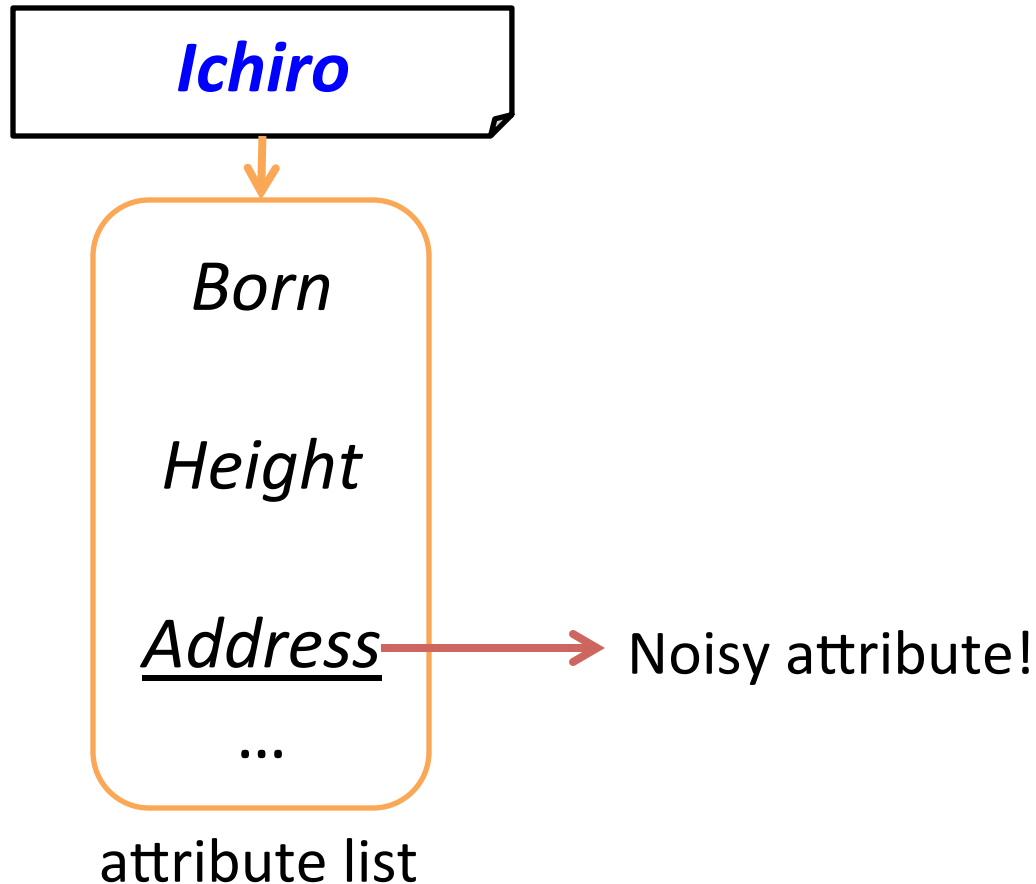
Adjacency matrix
of the cosine similarity

Damping factor

Scoring

2. Query attributes

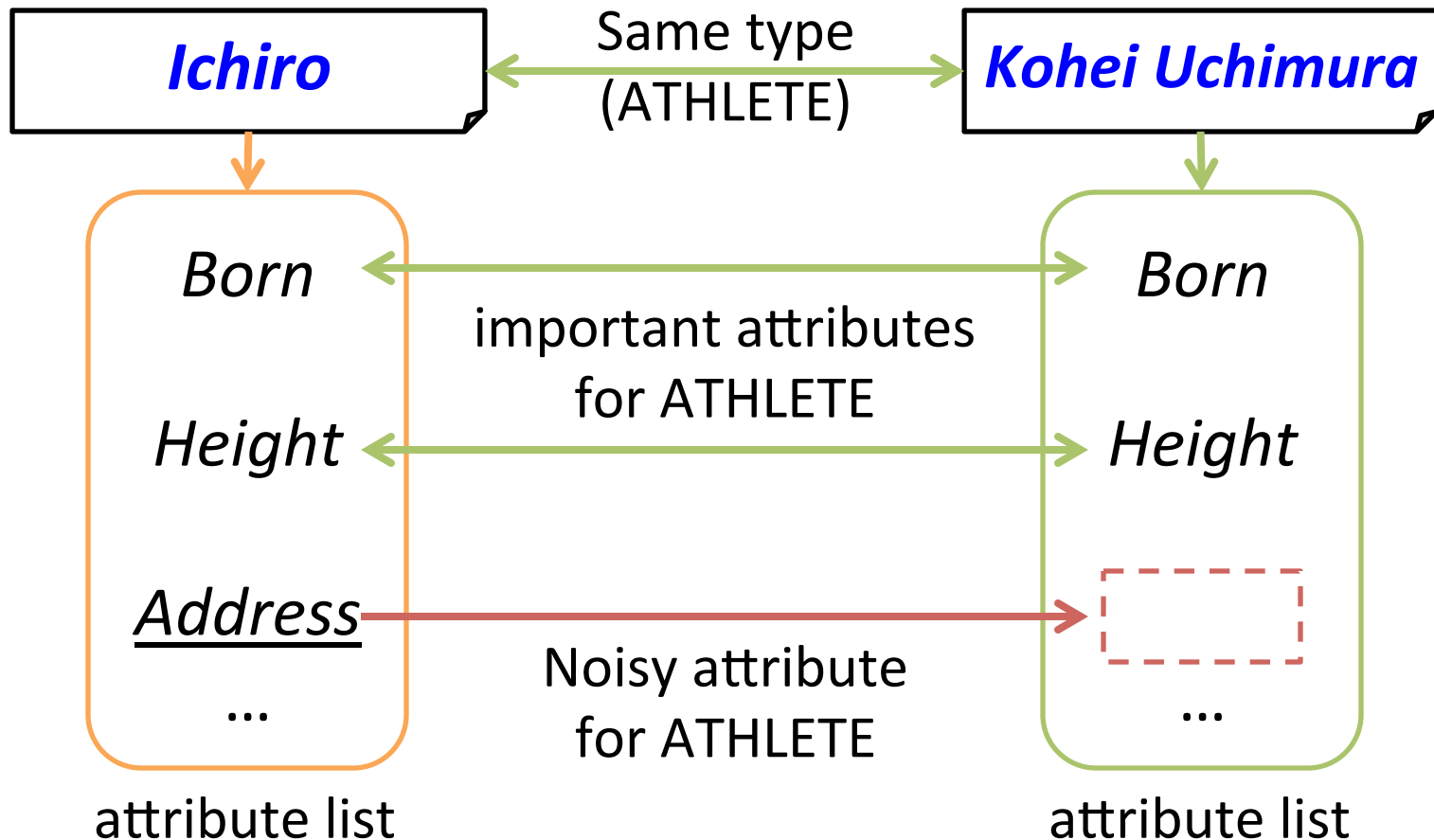
- Attributes extracted automatically are *noisy*



Scoring

2. Query attributes

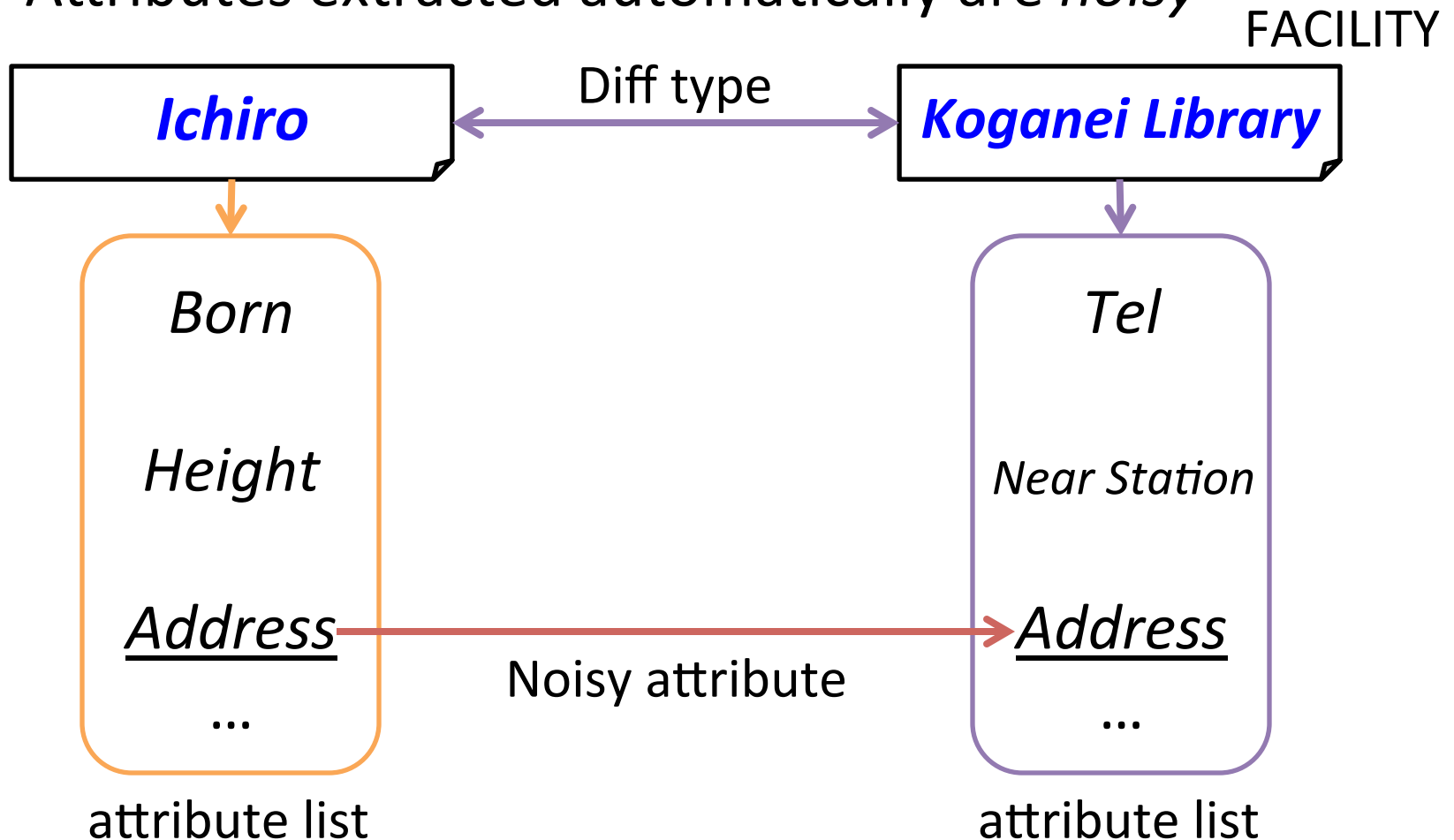
- Attributes extracted automatically are *noisy*



Scoring

2. Query attributes

- Attributes extracted automatically are *noisy*



Scoring

2. Query attributes

- Weight each attribute as importance like TF-IDF

Attribute list of the query

Same type attributes

A B

A

A

A

Attribute list of other query

Diff type attributes

B

B

B

We can give high priority to characteristic attributes

Diff type attributes

Weight of A > Weight of B

A B

B

B

A

Scoring

3. Manual clue words and URLs

- “**birthday**” or “**birth town**”: more relevant for celebrity
- “**talent.yahoo.co.jp**”: more relevant for celebrity
- “**address**” or “**phone number**”: more relevant for facility

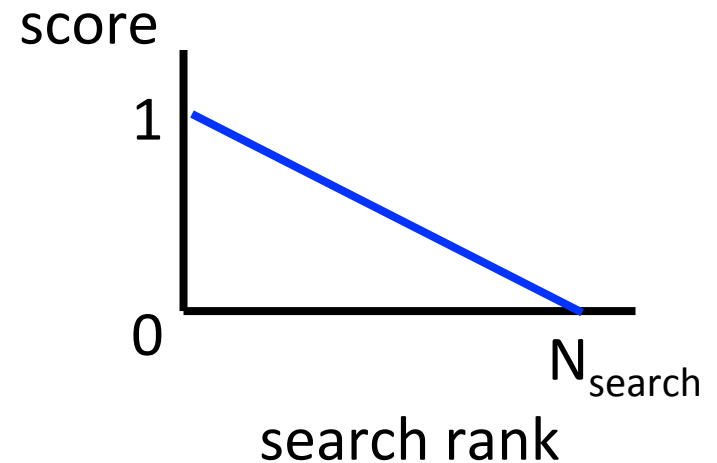
- “**posted by**” or “**answerer**”: less relevant for all
- “**amazon.co.jp**”: less important for all

Scoring

4. Search rank

- Sentences of high rank webpage is more important

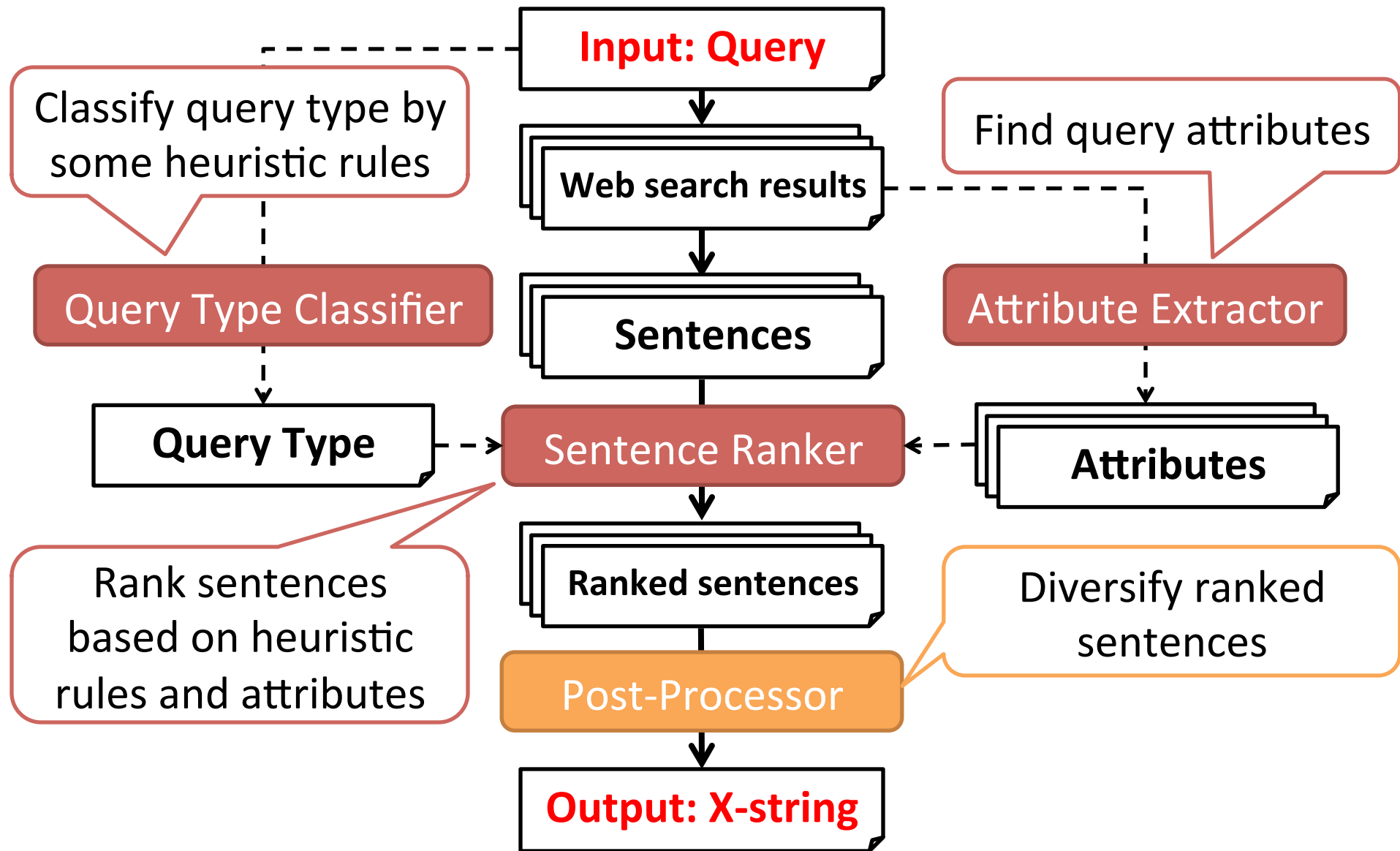
$$S_{SRank}(s) = 1 - \frac{rank(s)}{N_{search}}$$



Sentence Ranking

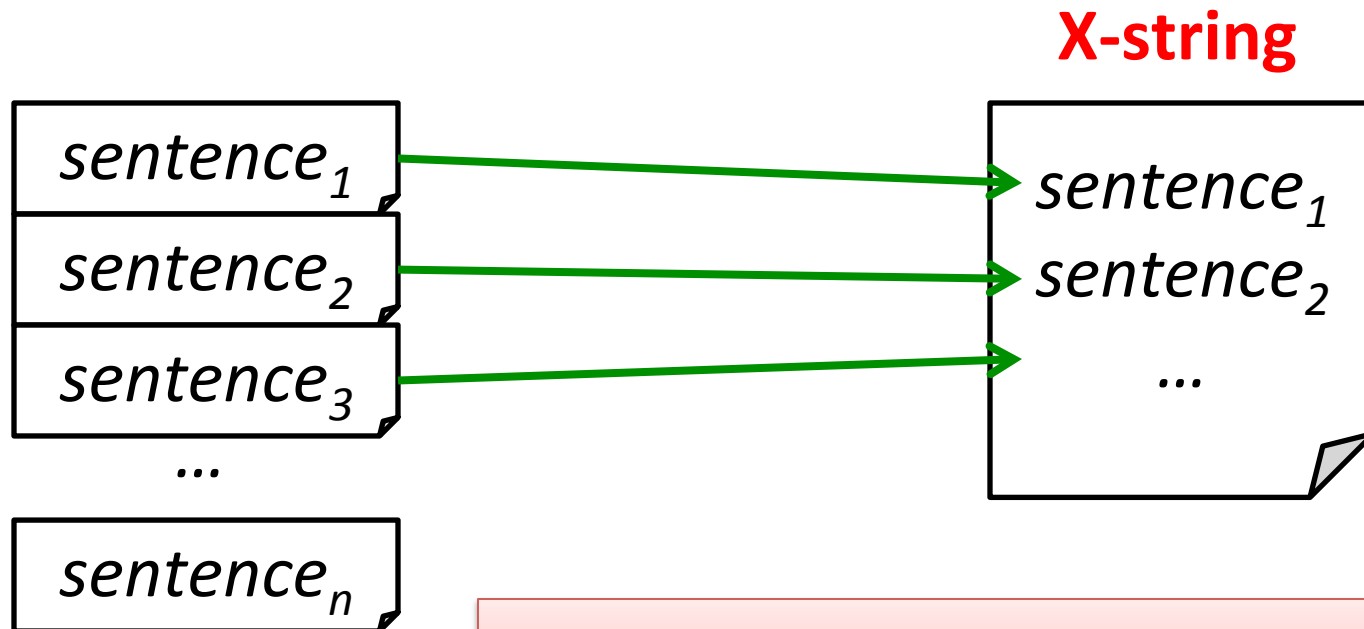
- Scoring based on following components:
 1. Content similarity
 - Similar sentences are more relevant
 2. Query attributes
 - Weight for each attribute as importance like TF-IDF
 3. Manual clue words and URLs
 4. Search Rank
 - Sentences of high rank webpage is more important

Overview



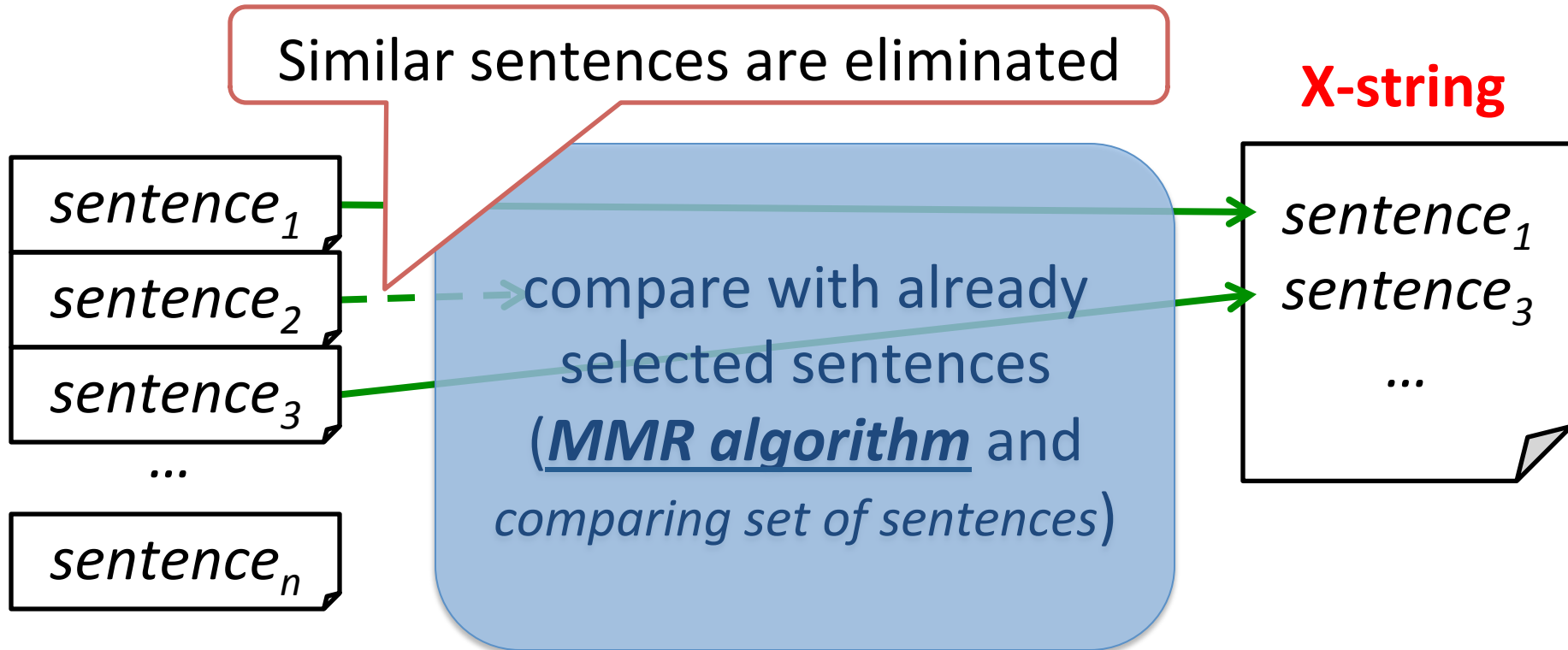
Sentence diversity

- Output top N sentences as X-string
 - But similar sentences tend to be output
 - Redundancy is undesirable for 1CLICK



We should diversify sentences

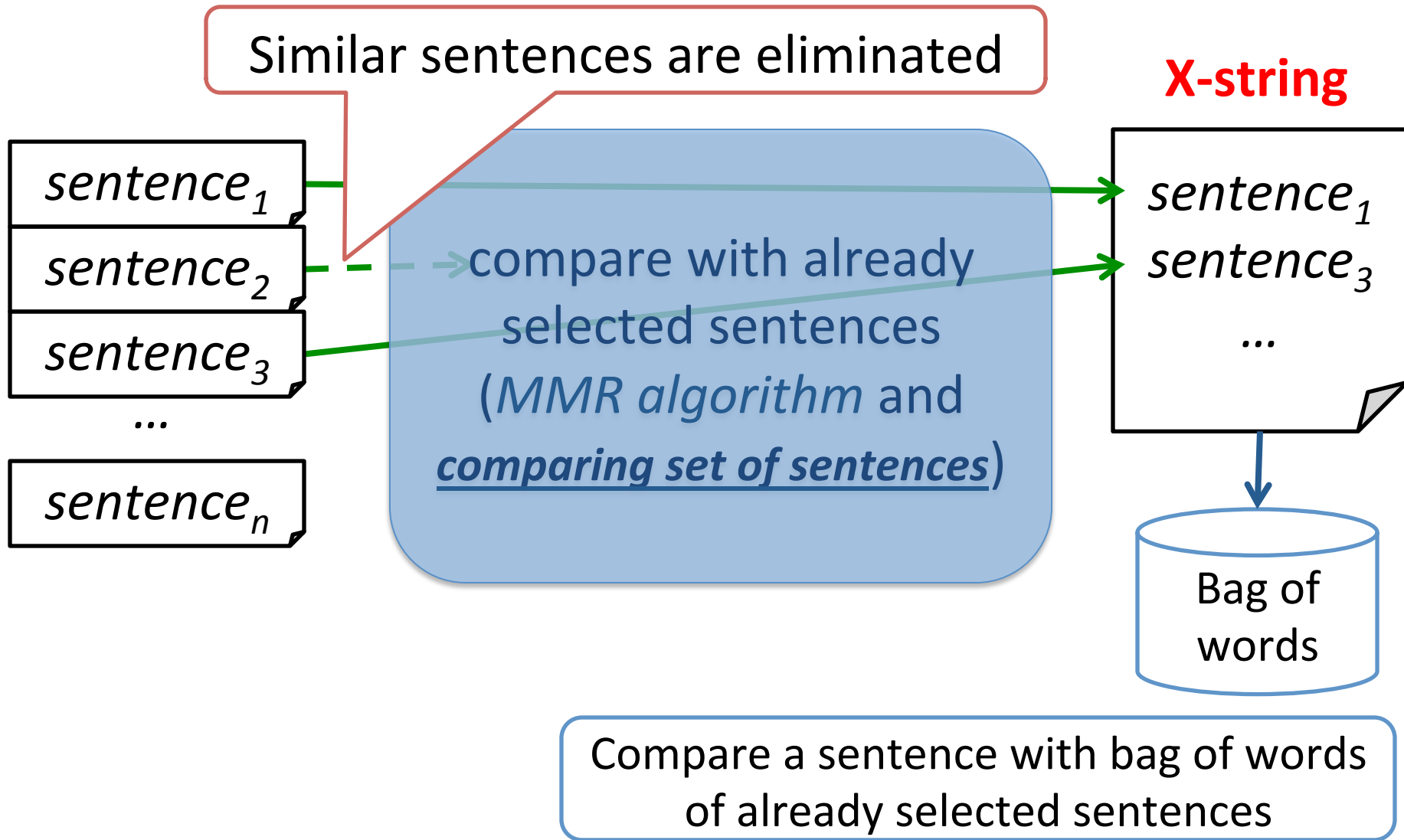
To consider sentence diversity



$$MMR = \arg \max_{s_i \in R \setminus S} \left[\lambda \text{Score}(s_i) - (1 - \lambda) \max_{s_j \in S} \text{Sim}(s_i, s_j) \right]$$

Give a penalty depending on similarity to already selected sentences

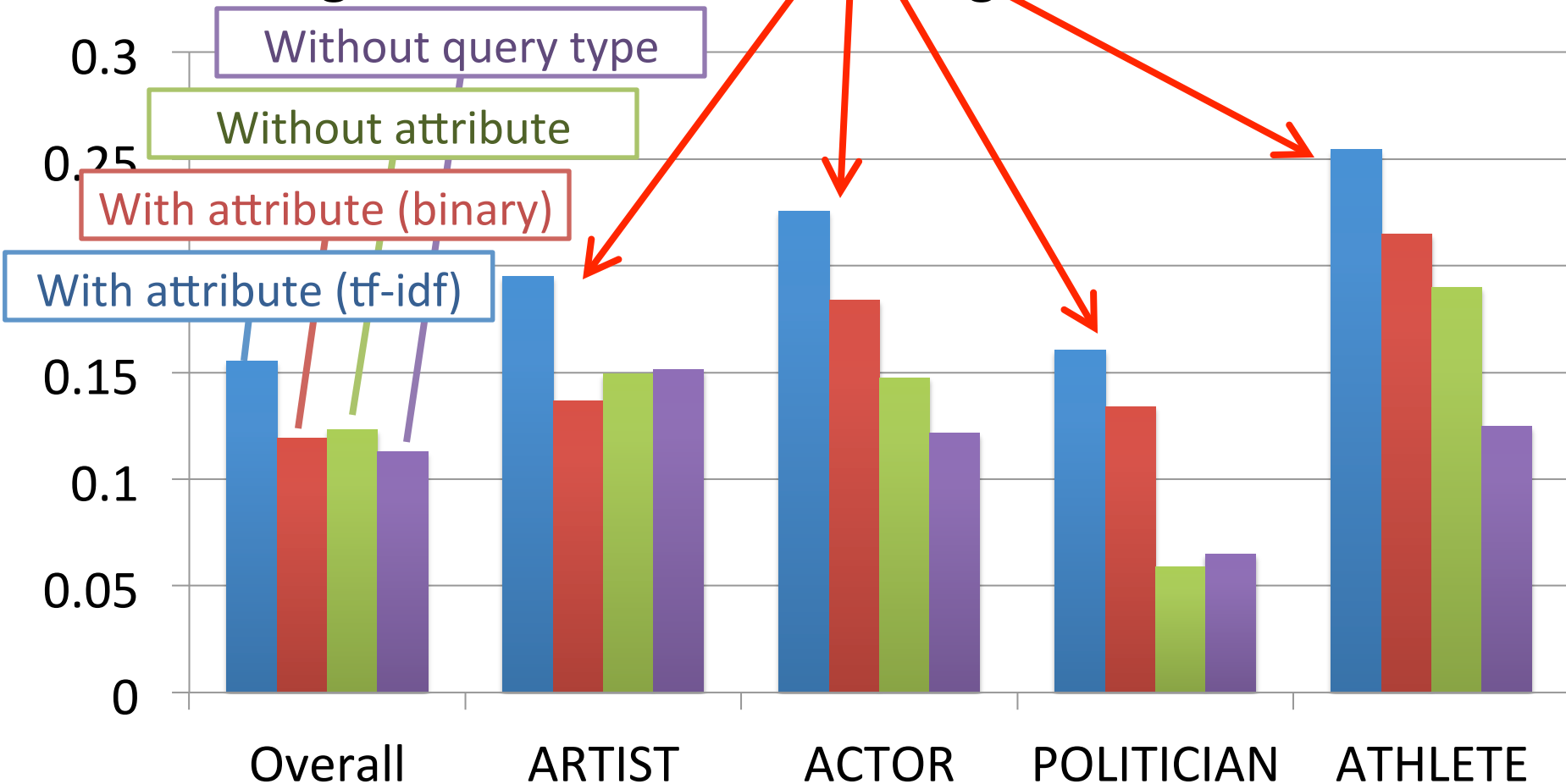
To consider sentence diversity



EVALUATION RESULTS

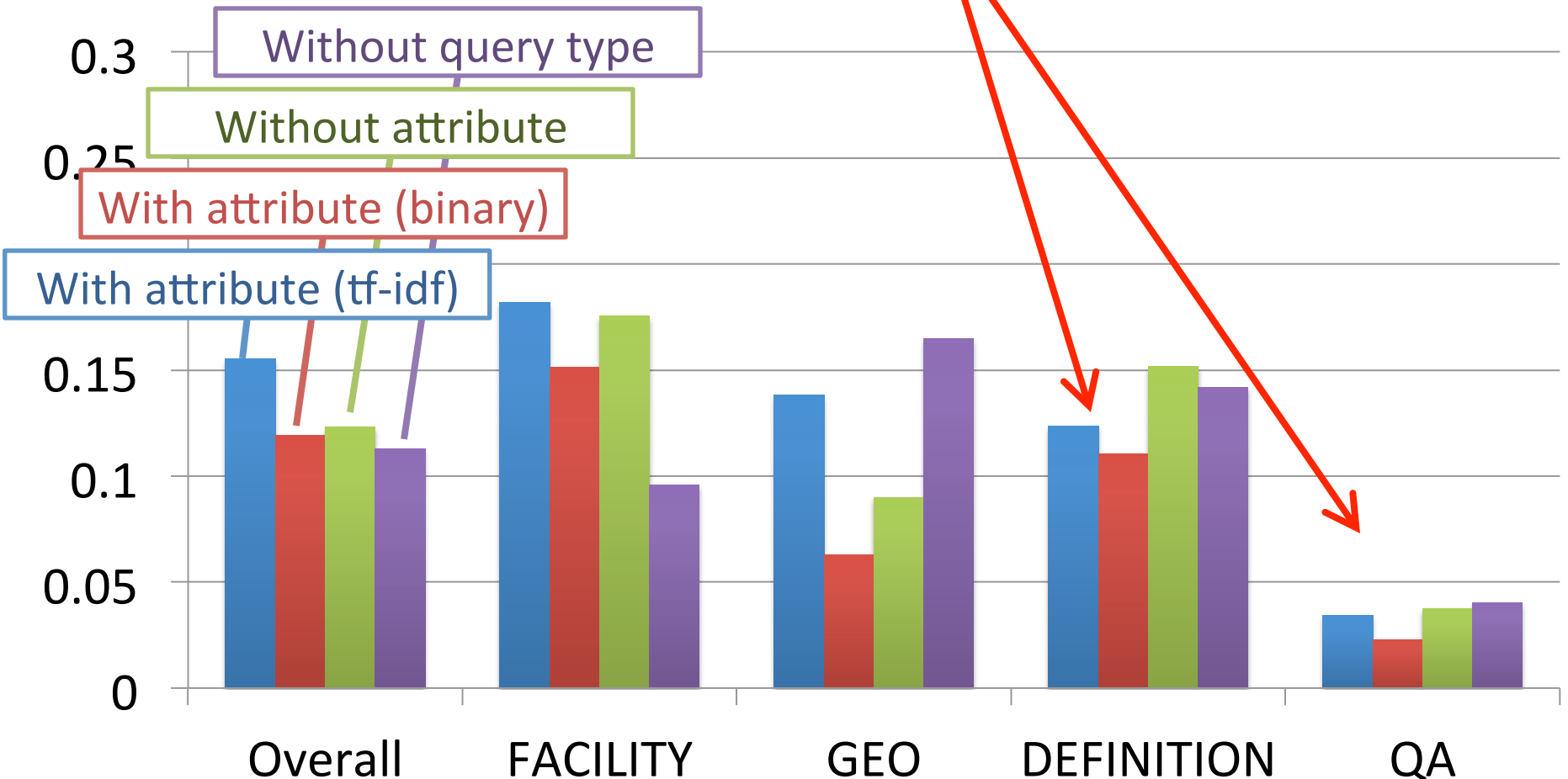
Results (mean S#-measures)

- Query attributes are **effective for celebrity**
 - Though the difference is not significant



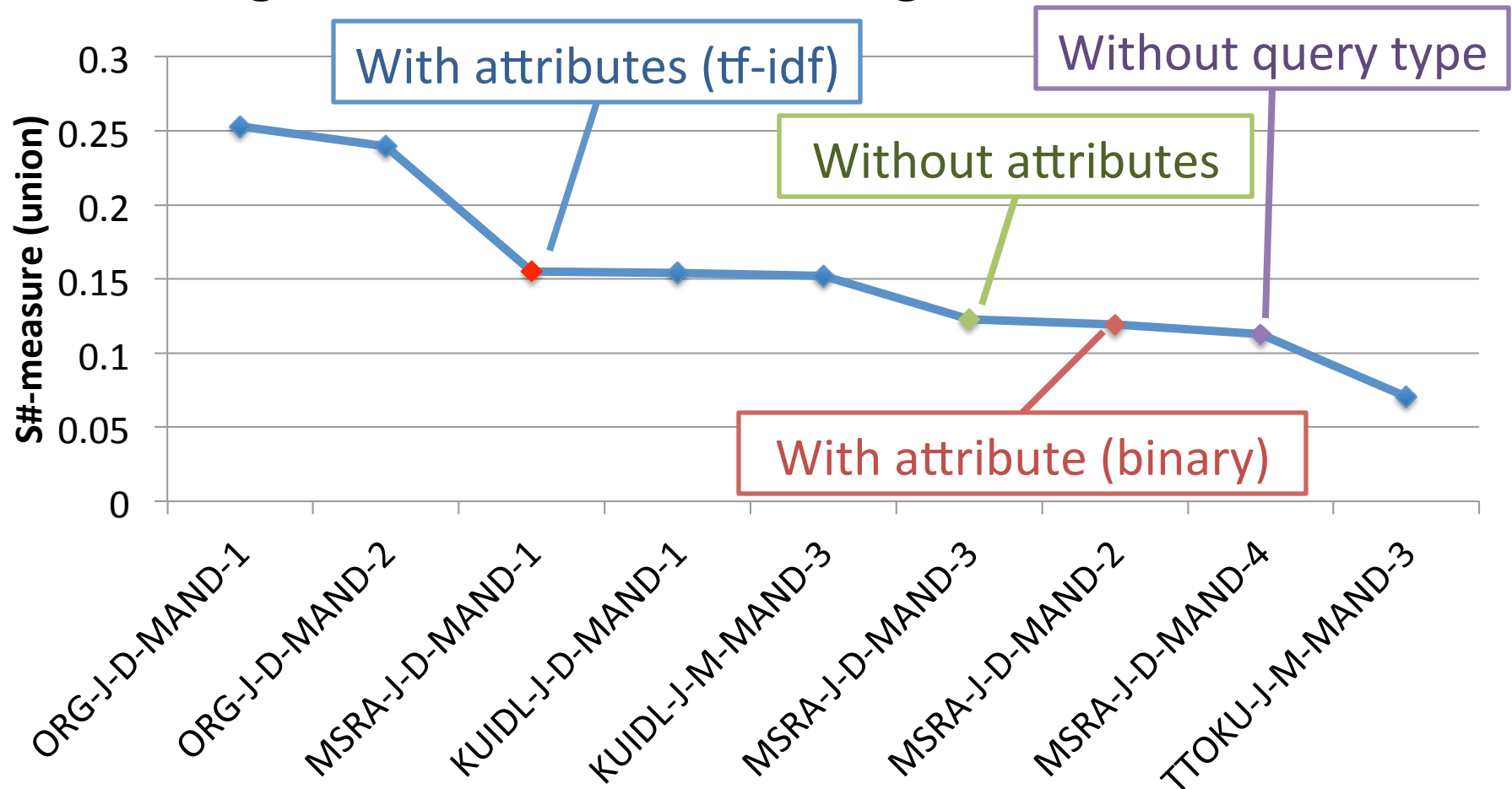
Results (mean S#-measures)

- Query attributes do **not work well for DEFINITION and QA queries**



Comparison with other participants

- Top performance among MANDATORY runs
 - Though the difference is not significant



Conclusions

- Automatic query attribute extraction is effective especially for celebrity queries
 - tf-idf weight is effective
- Future work:
 - Automatic extraction of clue words for query type classification
 - New framework instead of query attributes for DEFINITION and QA queries