

CACAO

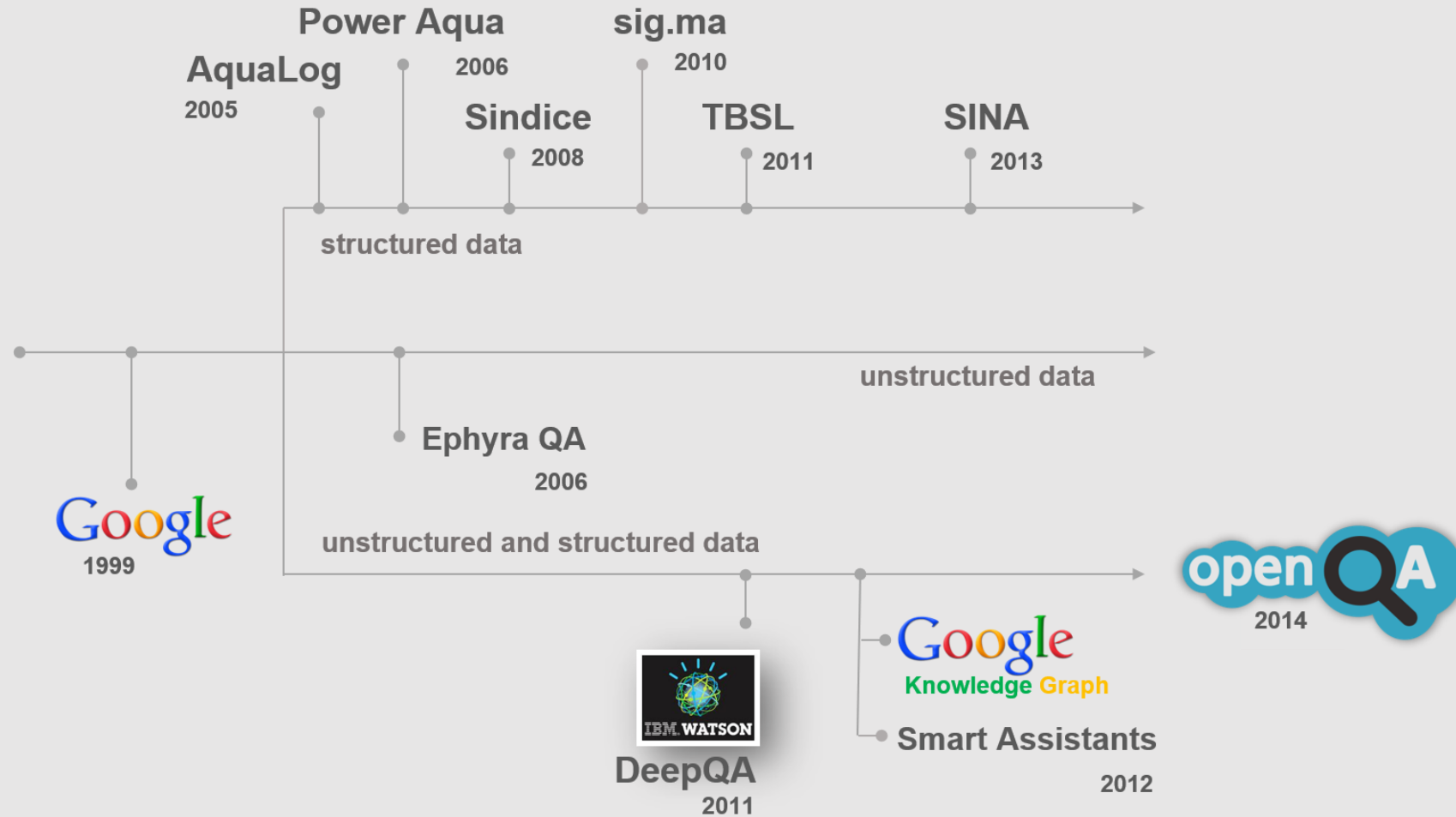
Conditional Spread Activation for Keyword Query Interpretation

15th International Conference on Semantic Web
Karlsruhe 2019

Edgard Marx, Gustavo Publico & Thomas Riechert



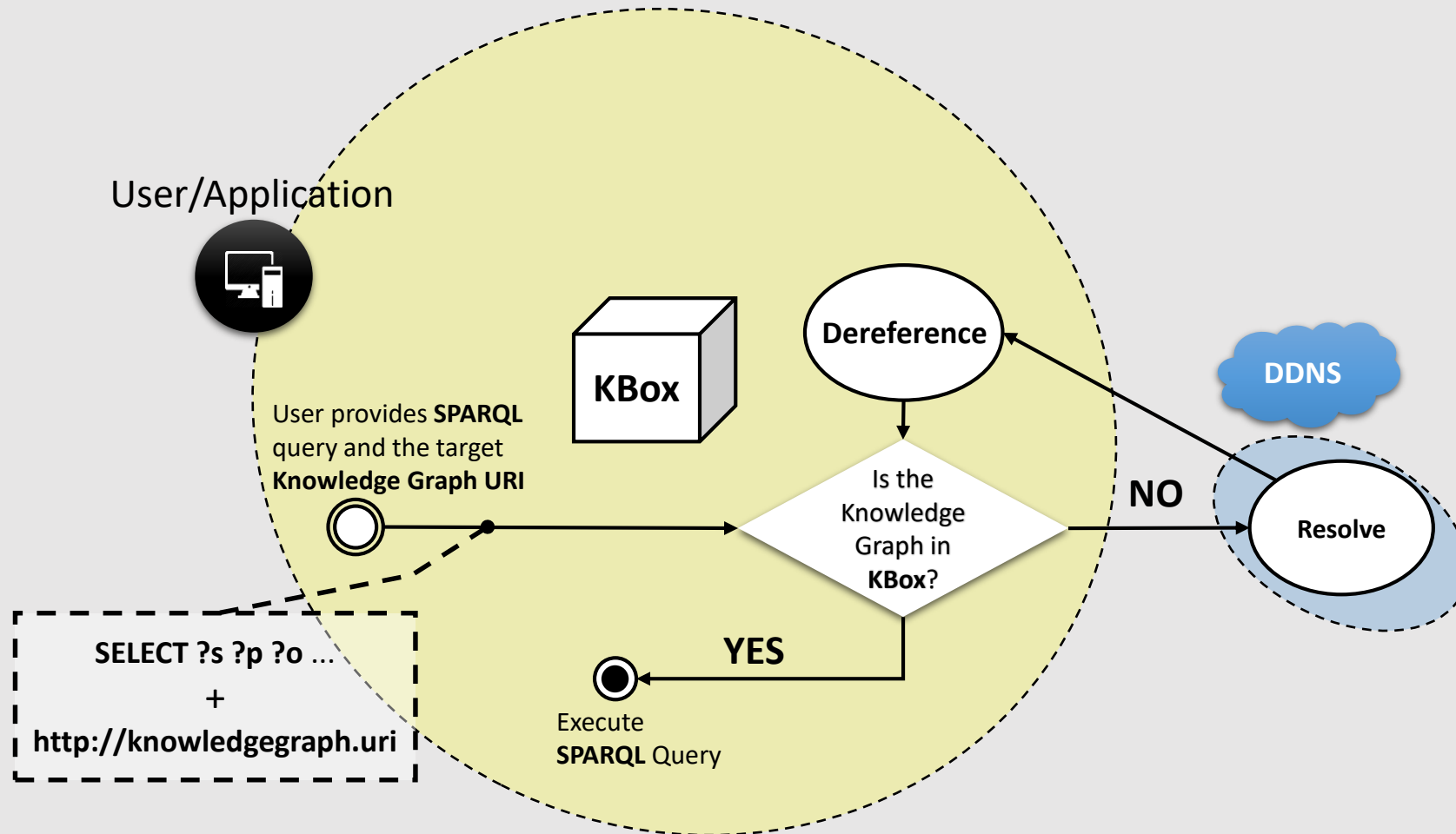
openQA



*<https://github.com/AKSW/openQA>

KBox

Previous Works



*<https://github.com/AKSW/KBox>

NSpM

Previous Works

Neural SPARQL Machines
@theLiberAI Follows you

Using Neural SPARQL Machines to translate natural language into machine language for data access. Part of @AKSWgroup. #DeepLearning #LinkedData

Leipzig, Germany
aksw.org/Projects/Neura...
Joined November 2017

Tweets 24 Following 9 Followers 39 Likes 36

Tweets Tweets & replies

Pinned Tweet
Neural SPARQL Machines @theLiberAI · Jul 10
We will be in Stockholm at the #ICML2018 workshop on Neural Abstract Machines & Program Induction to present our latest work. Check out our poster on July 15th! @akswgroup arxiv.org/abs/1806.10478

Neural SPARQL Machines @theLiberAI · Aug 31
We almost forgot to thank @near_ai for supporting the travel bursary which allowed us to present our poster. We are honored! #NAMPI #ICML2018

Neural SPARQL Machines Retweeted
DBpedia @dbpedia · Aug 8
Just recently, students of @edgardman integrated #DBpedia NSpM (Neural SPARQL Machines) with #Telegram and created an amazing #chatbot. Check it out tinyurl.com/ydb9w5nf

Who to follow Refresh · View all

- Knowledge Graphs and S...** Follow
- Piyush Chawla** @piyush_b... Follow
- Henriette Harmse** @Harm... Follow

Find people you know
Import your contacts from Gmail

Connect other address books

Trends for you Change

- #DiaMundialRBD14Años
Os fãs estão lembrando da 'era Rebelde'!
- #QuintaDetremuraSDV

DBNQA

Previous Works

Question Answering data [\[edit \]](#)

This section includes datasets that deals with structured data.

Dataset Name ↕	Brief description ↕	Preprocessing ↕	Instances ↕	Format ↕	Default Task ↕	Created (updated) ↕	Reference ↕	Creator ↕
DBpedia Neural Question Answering (DBNQA) Dataset	A large collection of Question to SPARQL specially design for Open Domain Neural Question Answering over DBpedia Knowledgebase.	This dataset contains a large collection of Open Neural SPARQL Templates and instances for training Neural SPARQL Machines; it was pre-processed by semi-automatic annotation tools as well as by three SPARQL experts.	894,499	Question-query pairs	Question Answering	2018	[402] [403]	Hartmann, Soru, and Marx et al.

https://wikipedia.org/wiki/List_of_datasets_for_machine-learning_research

*<https://github.com/AKSW/DBNQA>

Outline

- Motivation
- Problem Statement
- Background
- Approach
- Evaluation
- Conclusion & Future Works

Linked Data

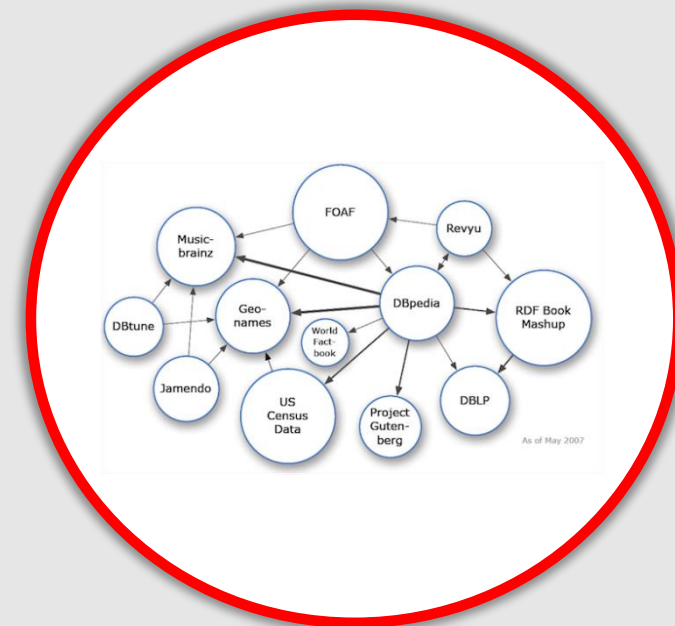
Unstructured



Structured

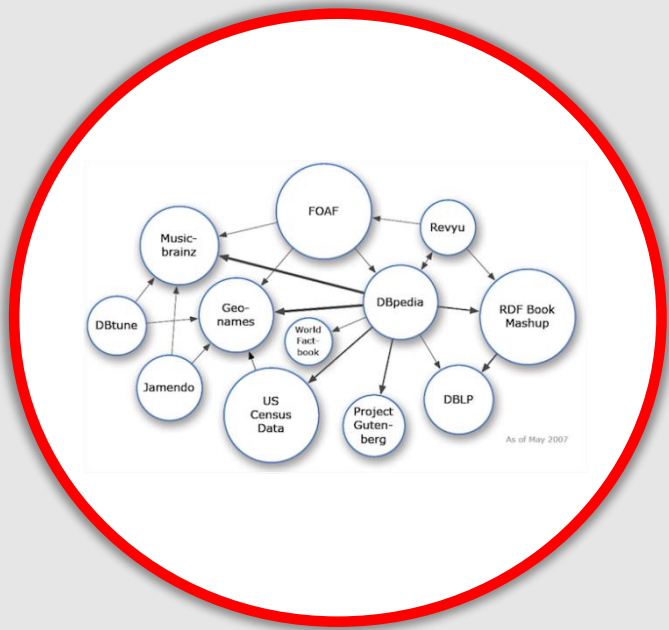


Crowd

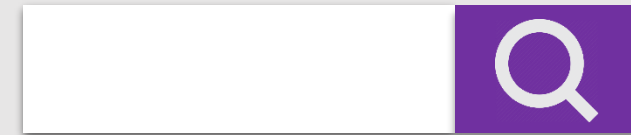


Linked Data

Access (the information needed)



Linked Data



- Entity Retrieval
- Question Answering
- Entity Linking

Entity Retrieval

Formally, a top-k entity retrieval takes a keyword query Q , an integer $0 < k$, a set of entities $E = \{e_1, e_2, \dots, e_{|E|}\}$, and returns the top-k entities based on a scoring function $S(Q, e)$

Entity Retrieval

```
Select * where {  
    ?s dbo:birthPlace dbpedia:Leipzig  
}
```

TF-IDF

$$w_{t,d} = tf_{t,d} \cdot \log \frac{|D|}{|\{d' \in D | t \in d'\}|}$$

- 1: Winter is coming.
2: Ours is the fury.
3: The choice is yours.



<u>term</u>	<u>freq</u>	<u>documents</u>
choice	1	3
coming	1	1
fury	1	2
is	3	1, 2, 3
ours	1	2
the	2	2, 3
winter	1	1
yours	1	3

Dictionary Postings

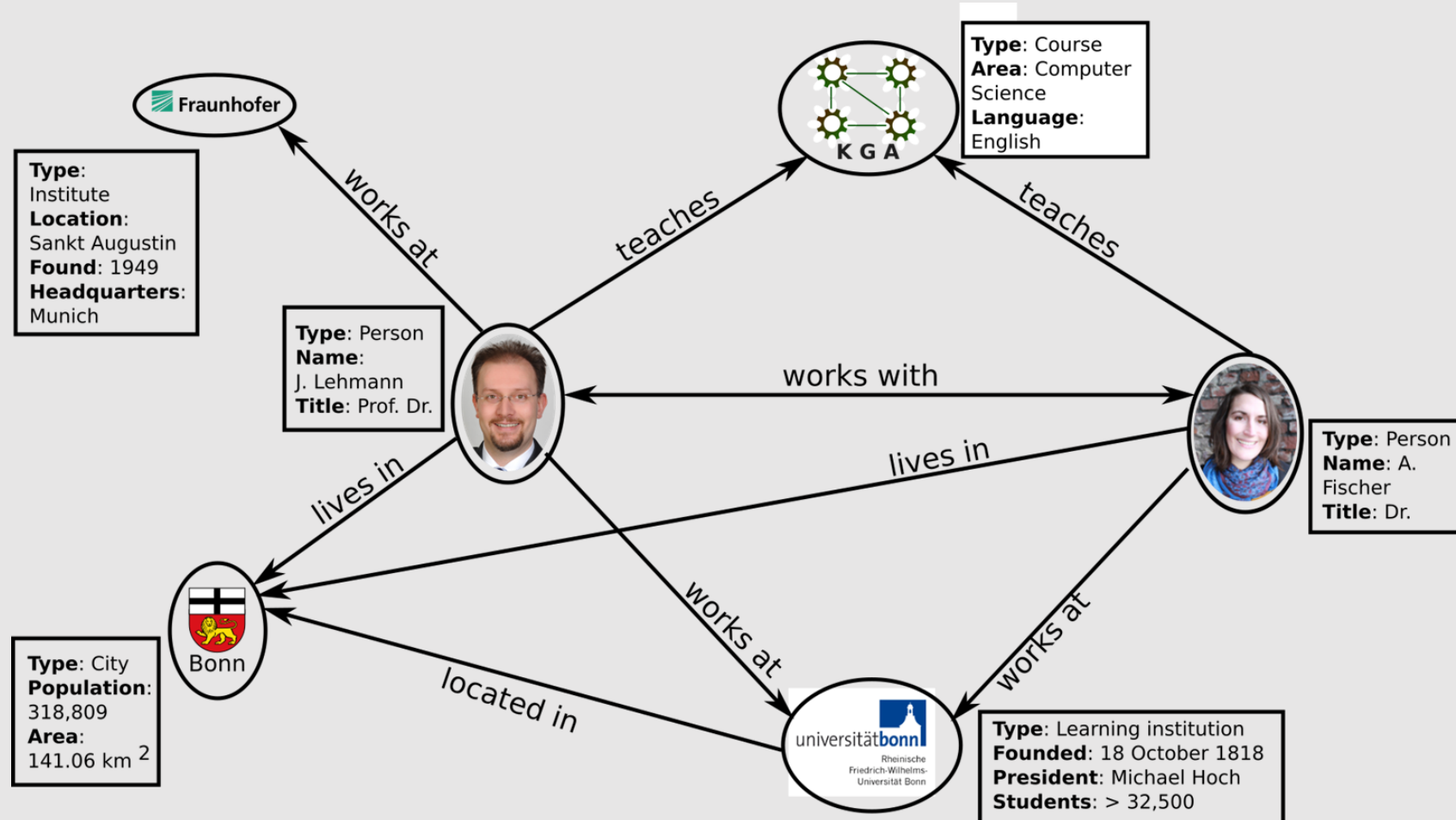
BM25

"This can opener can open any can that a can opener can and if this can opener cannot open any can that a can opener can, then this can opener is free"

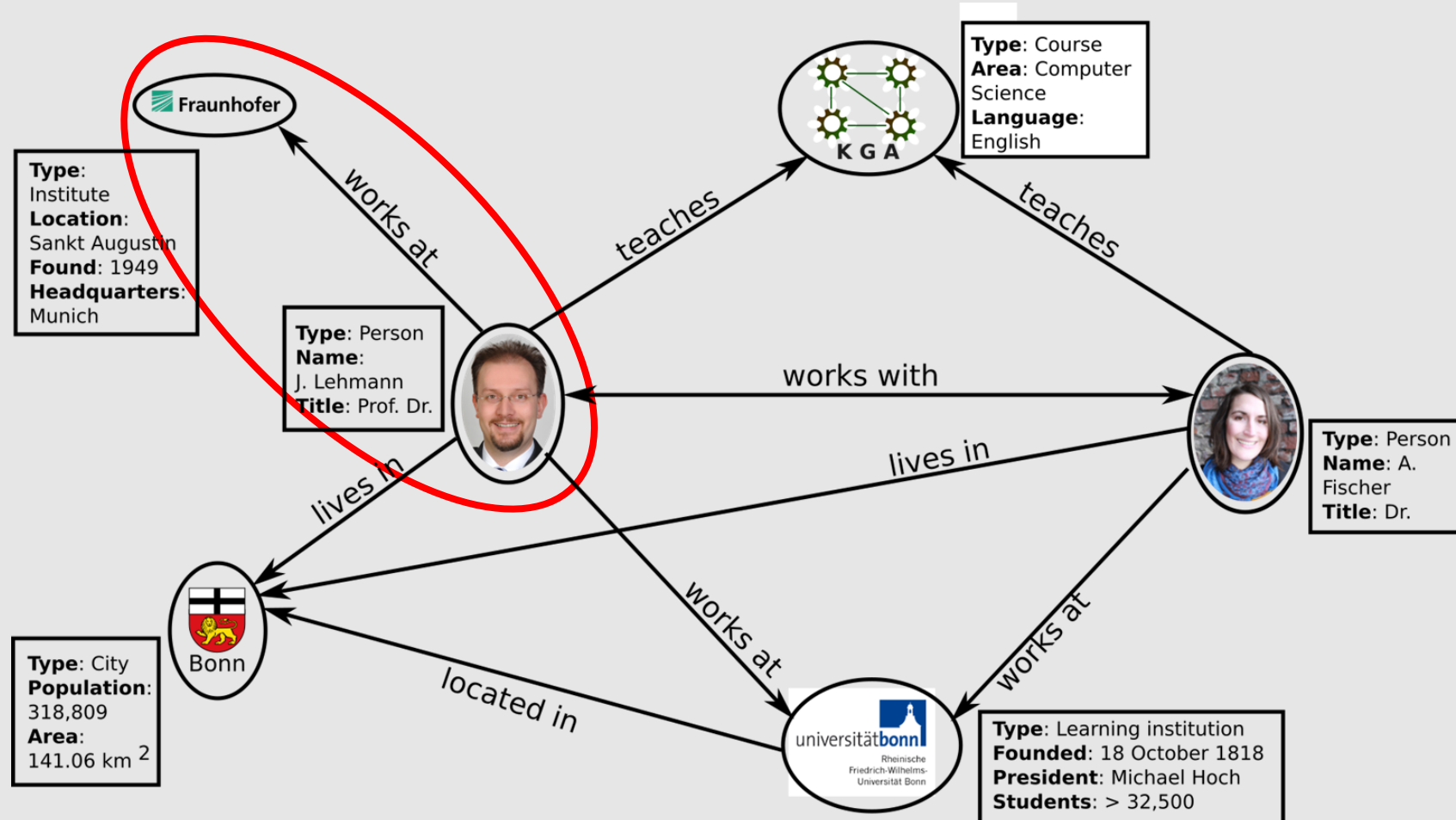


"If this can opener does not open any can, then you can take it for free"

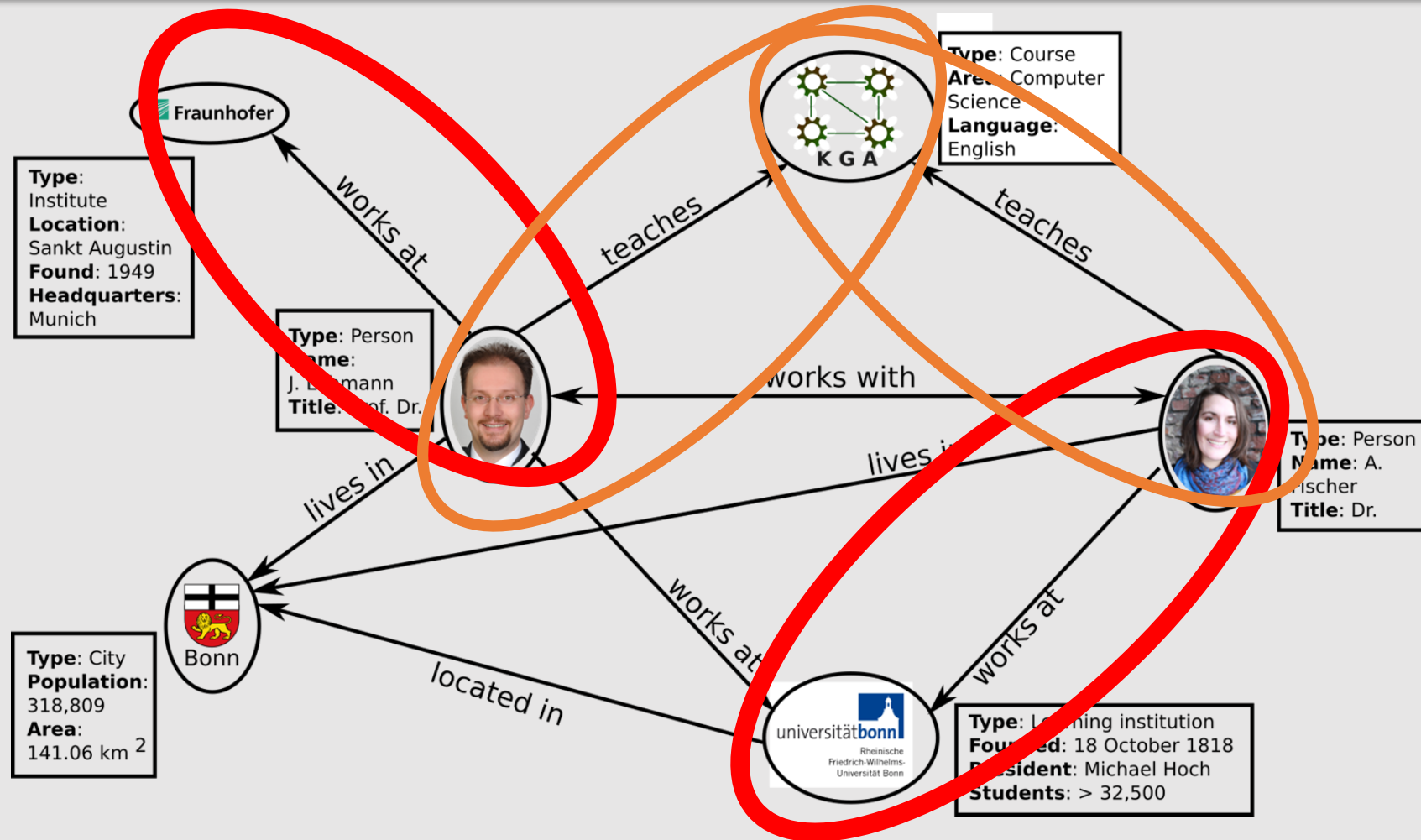
BM25



BM25F

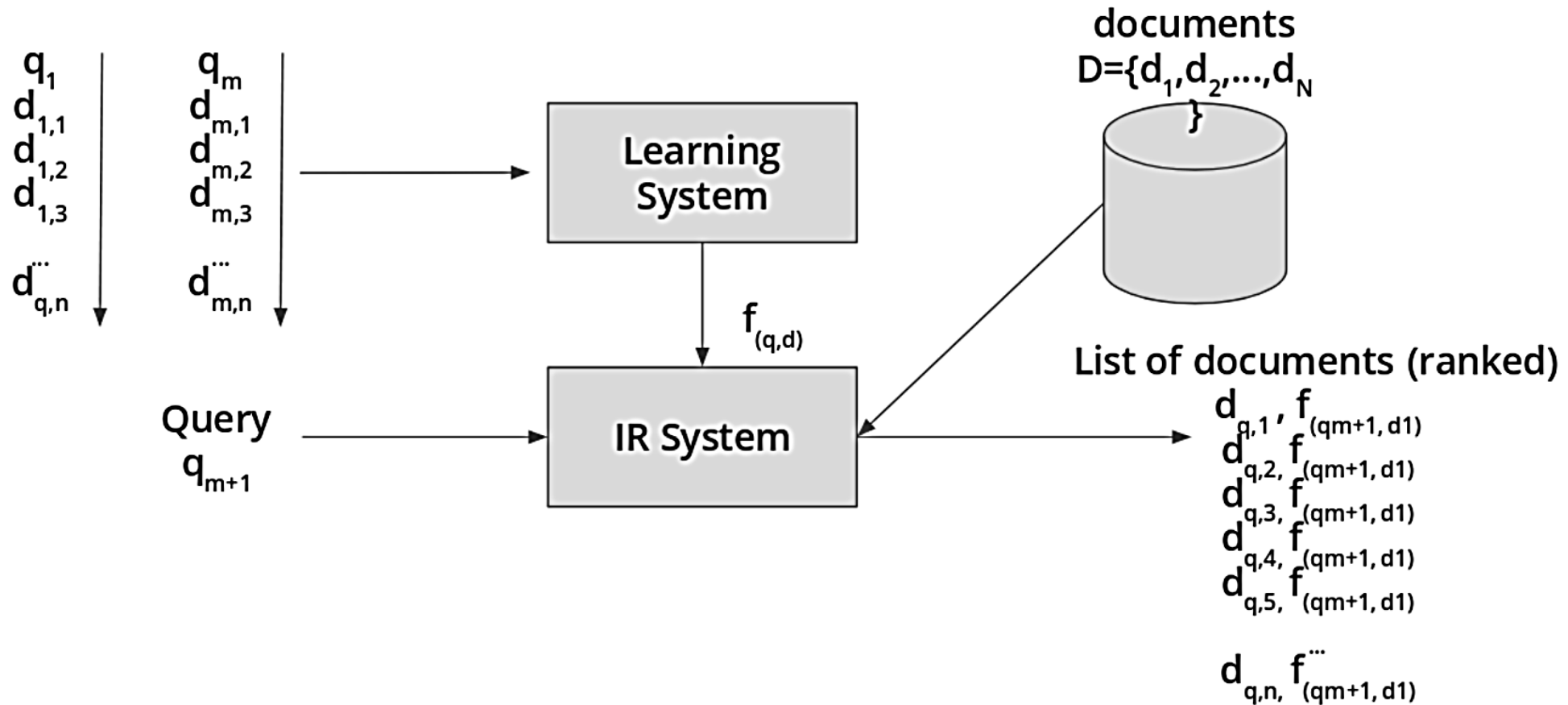


Field Weighted

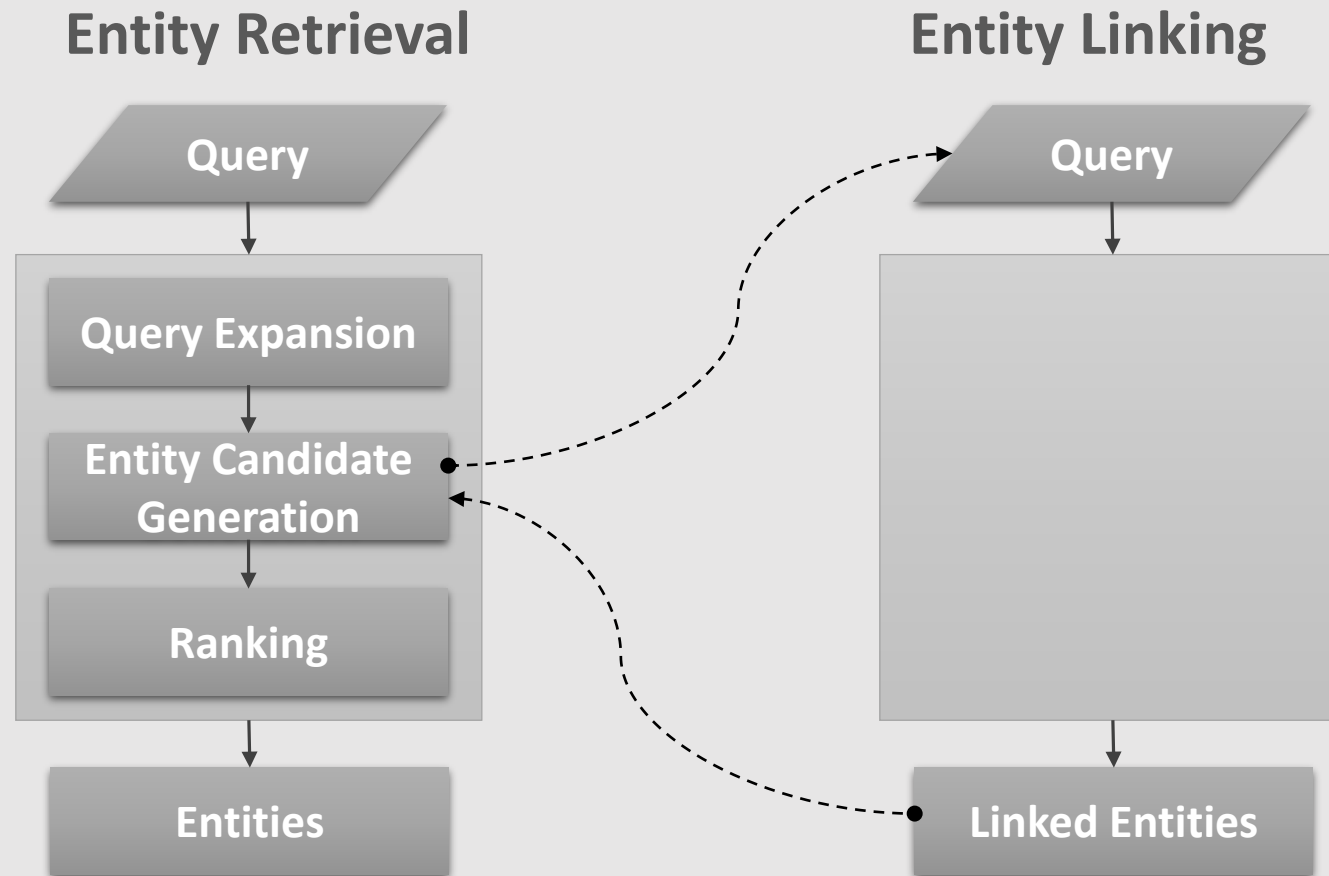


$\text{weight}(p1) > \text{weight}(p2) > \text{weight}(p3) \dots$

Learning to Rank

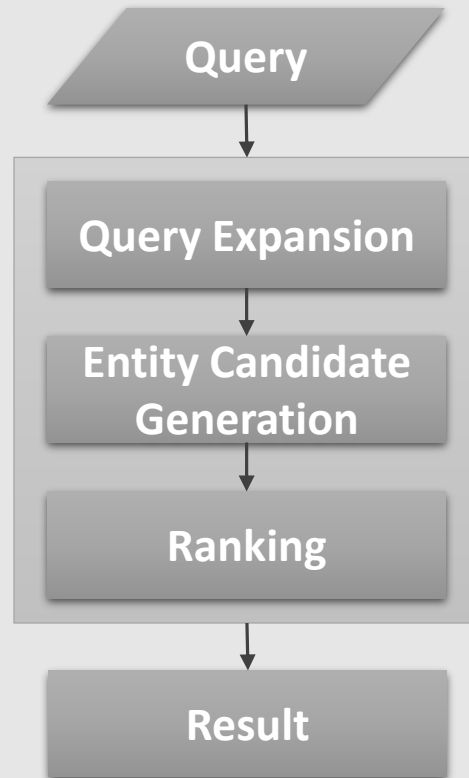


Entity Link in Queries

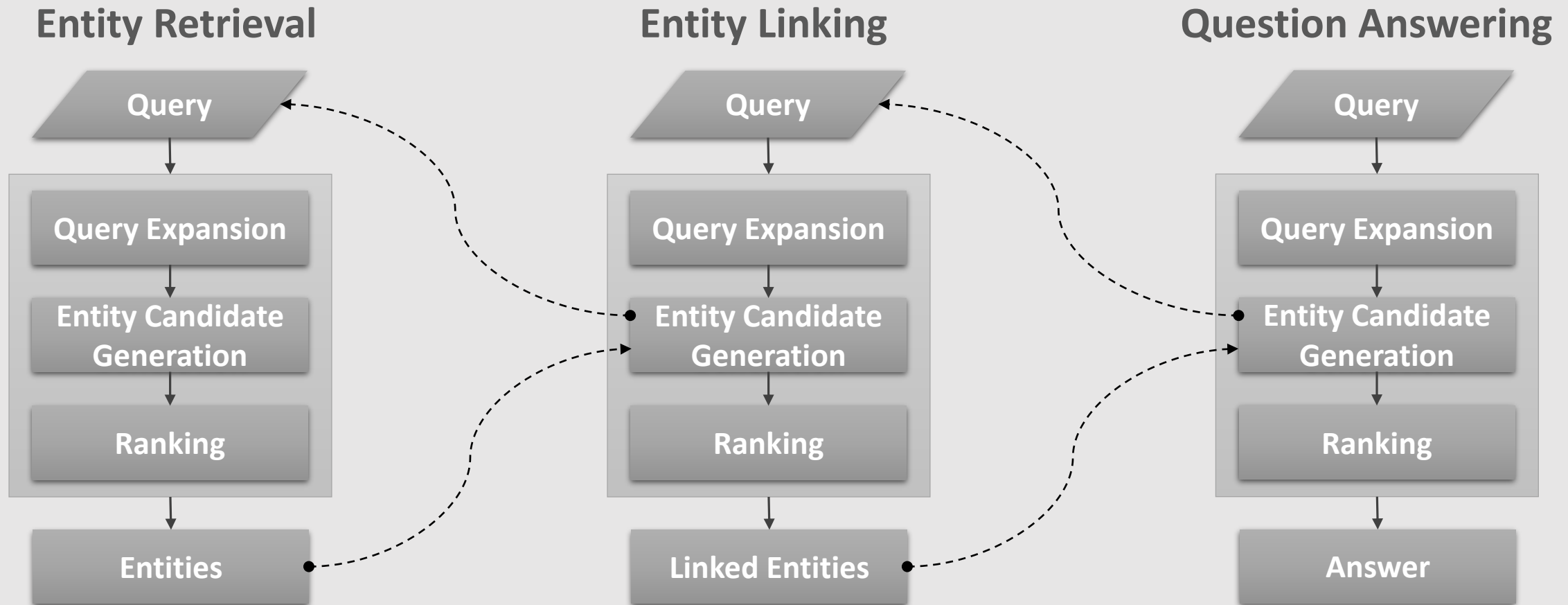


General IR Architecture

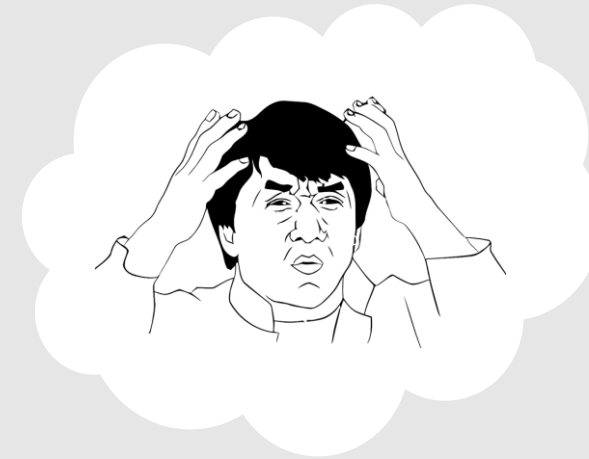
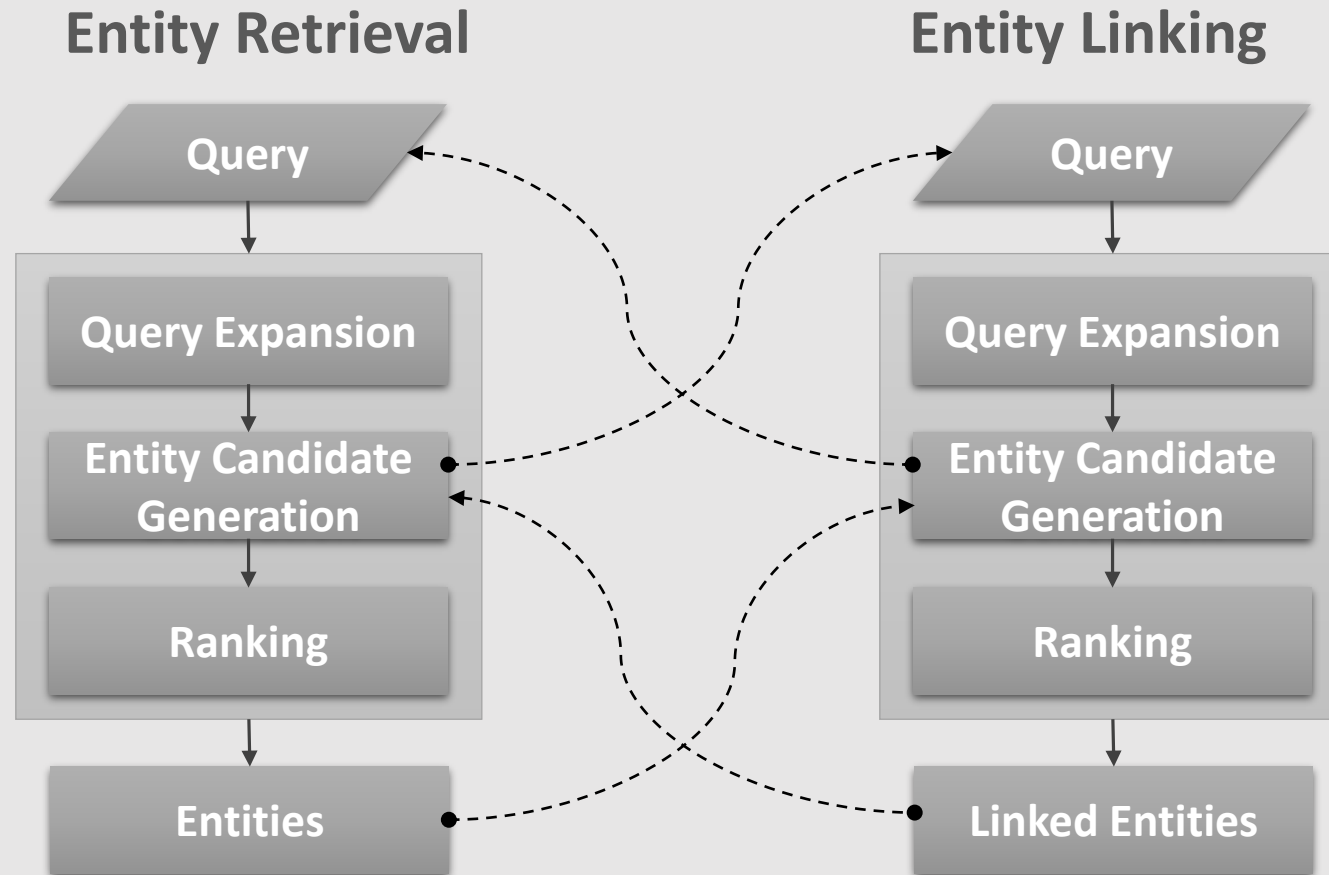
Entity Retrieval, Entity Linking, Question Answering



Interaction

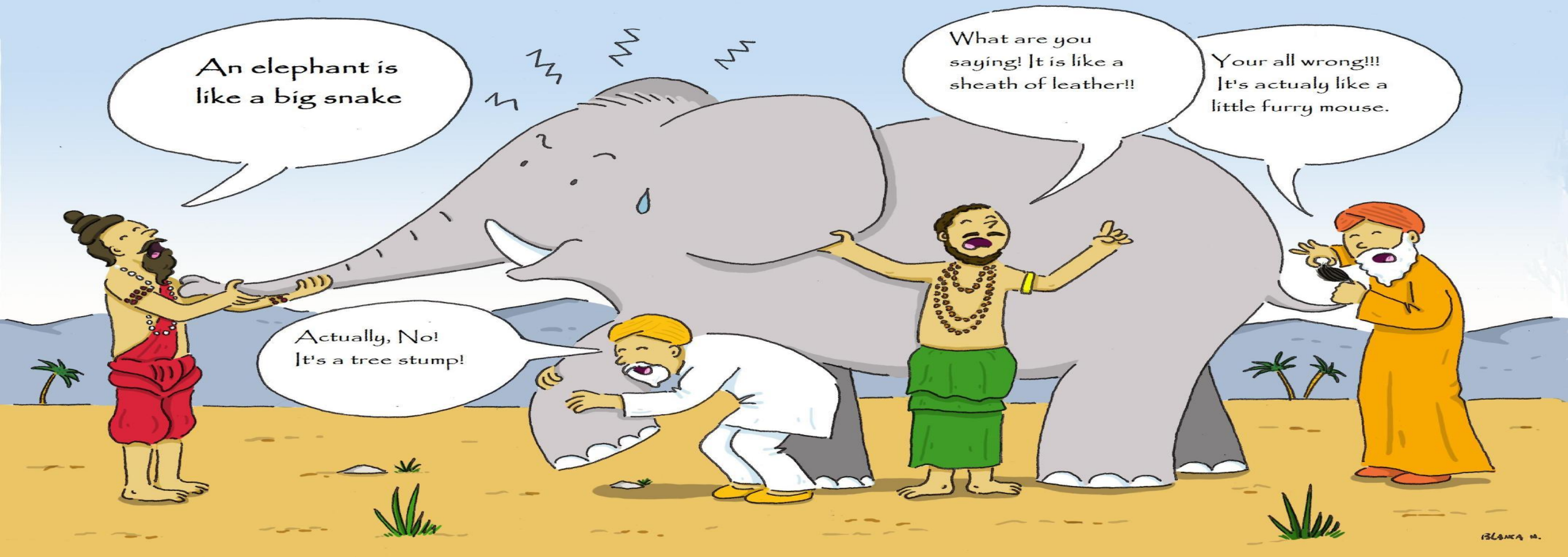


Lastest Architecture

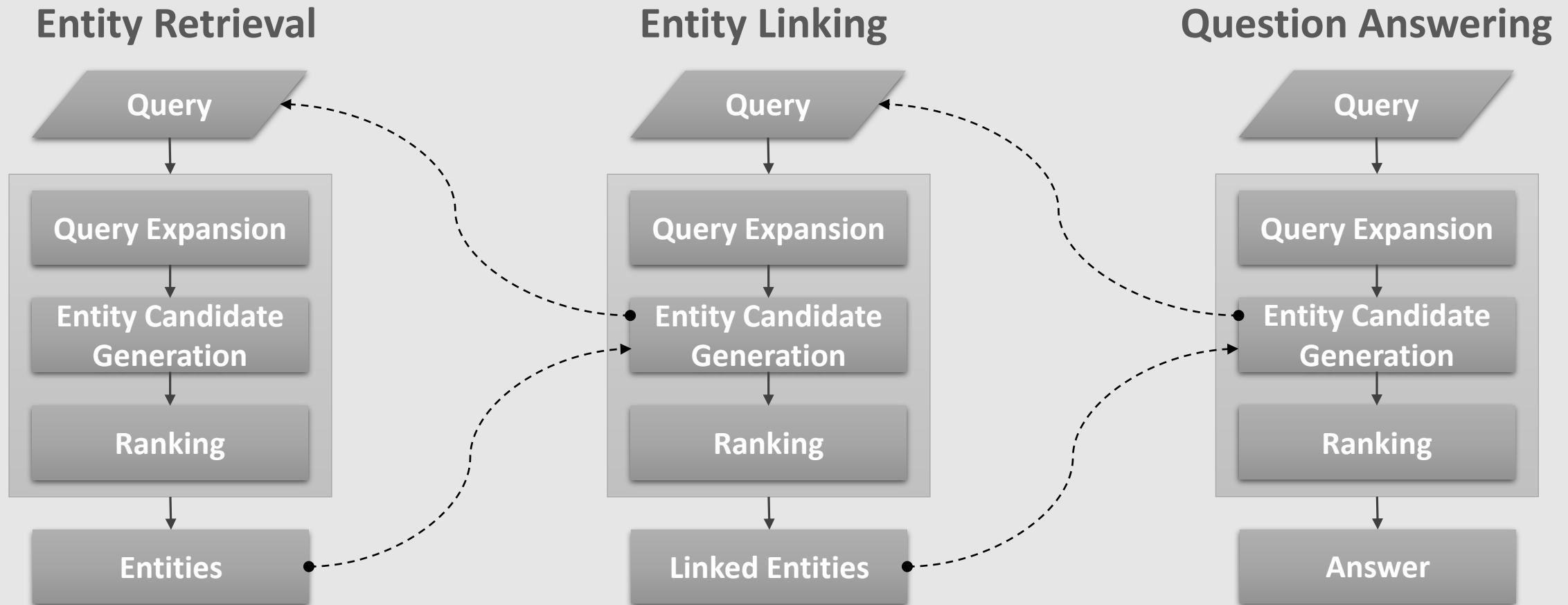


Background

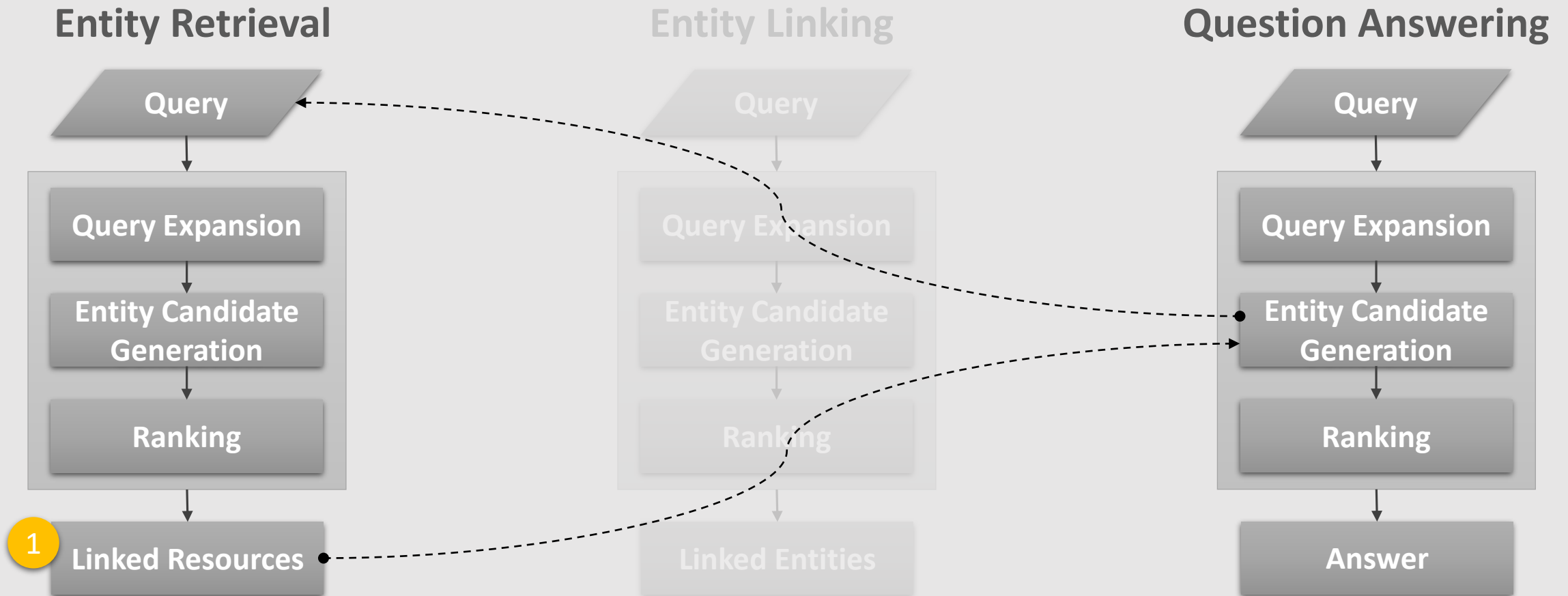
Lastest Architecture



Current Scenario

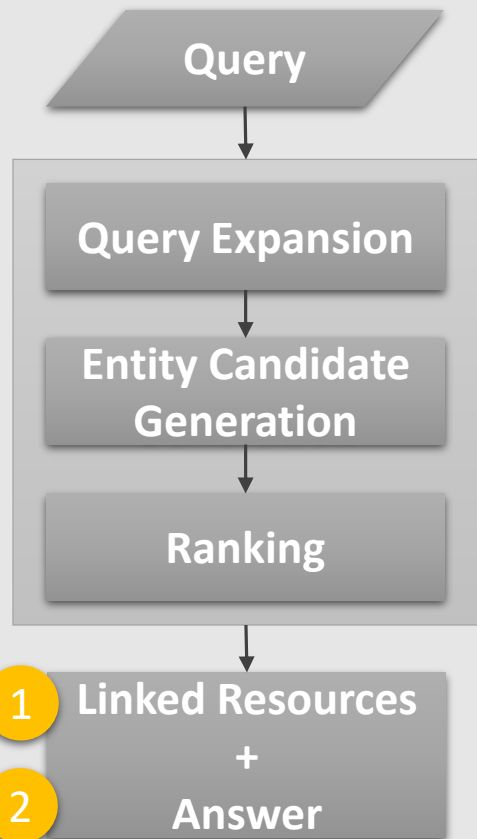


Return linked resources

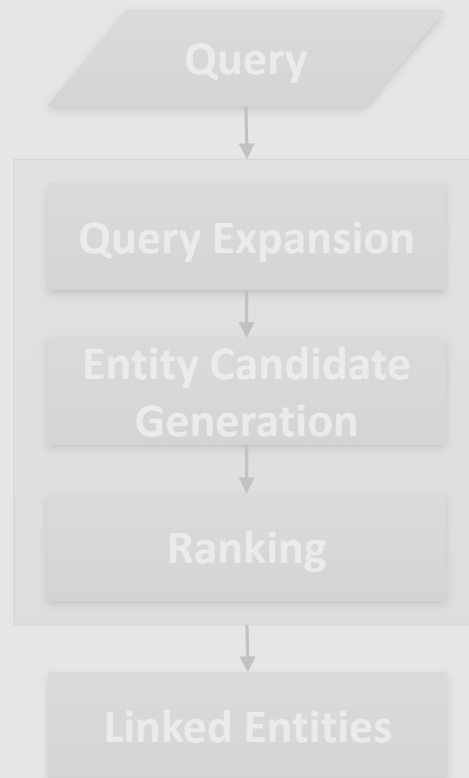


Return the answer

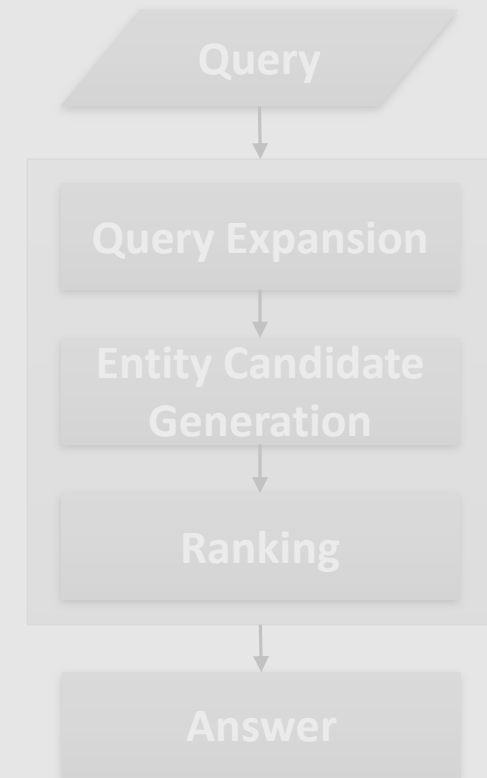
Entity Retrieval



Entity Linking

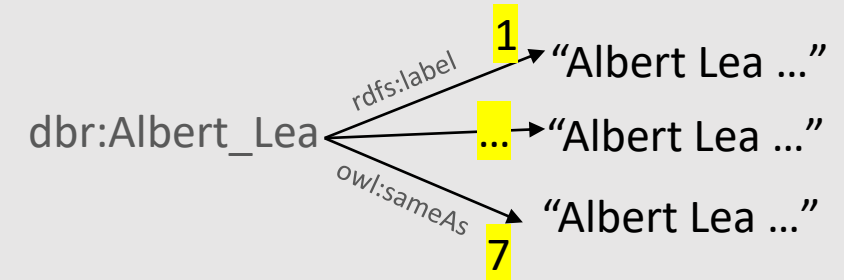
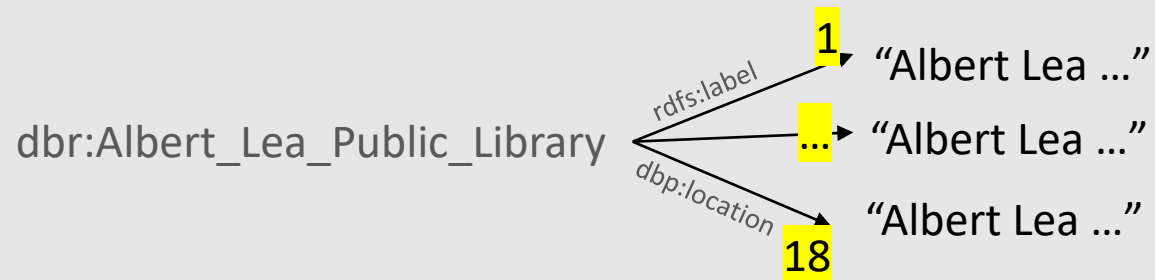


Question Answering



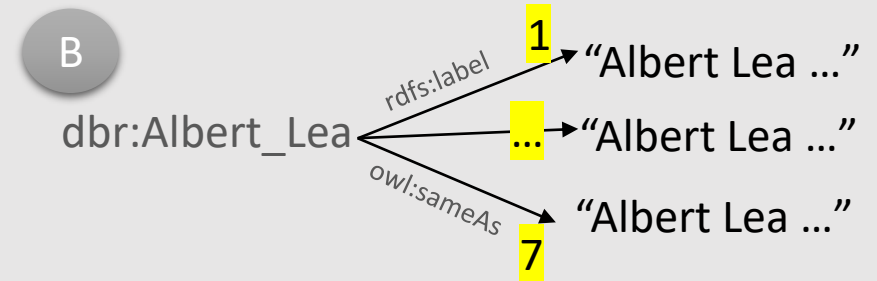
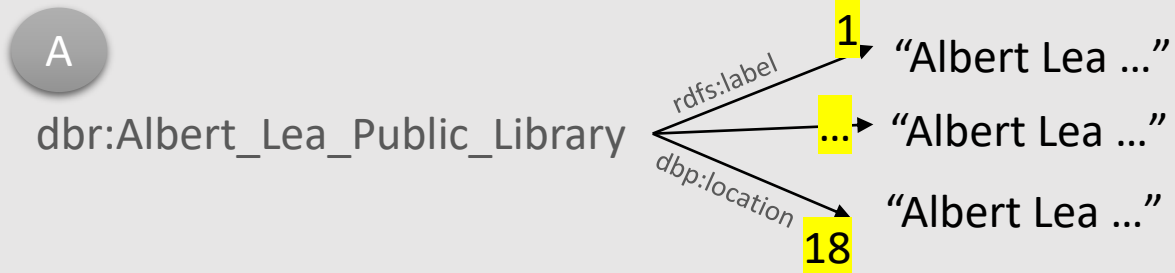
Challenge 1: *Local Term Frequency

Q: “Albert Lea”



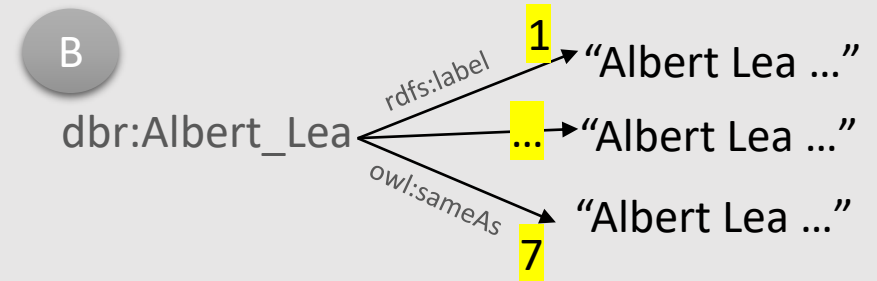
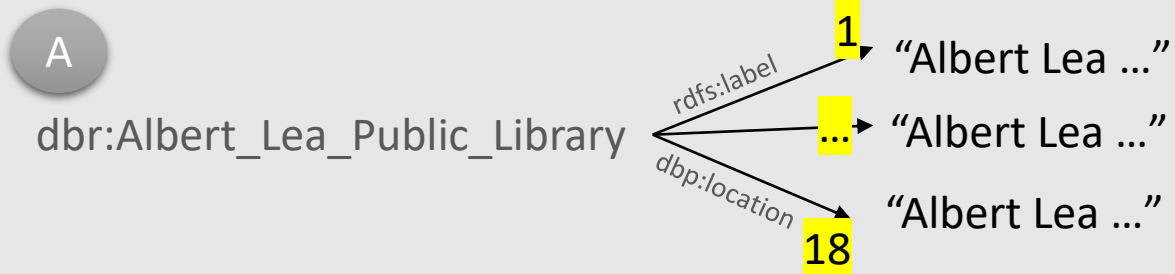
Challenge 2: Cohesion

Q: "Albert Lea"



Challenge 2: Cohesion

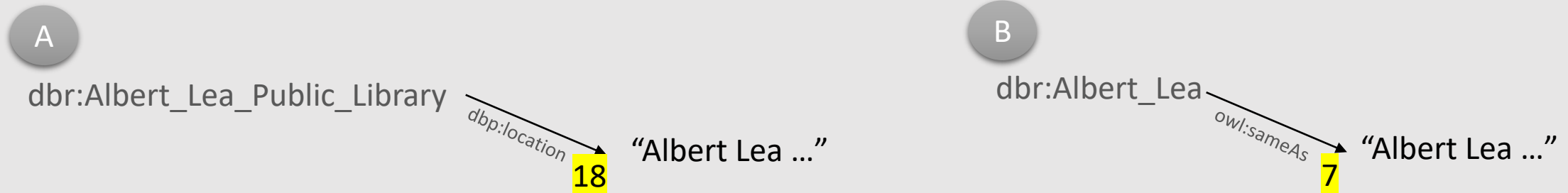
Q: “Albert Lea”



Preposition 1 [1][2]: maximize the number of tokens and reduce the number of mapped entities

Challenge 2: Cohesion

Q: “Albert Lea”

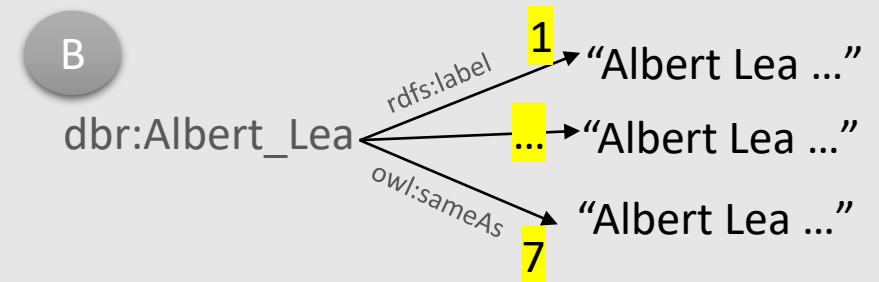
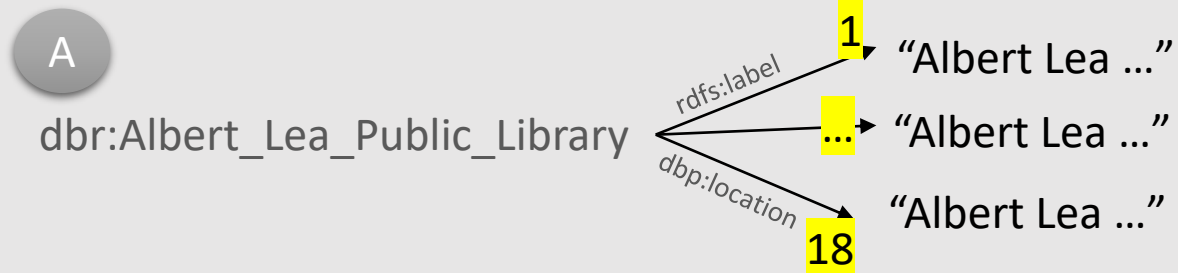


Proposition 1 [1][2]: maximize the number of tokens and reduce the number of mapped entities

Does **NOT** satisfy: Either A and B comply with the proposition 1

Challenge 2: Cohesion

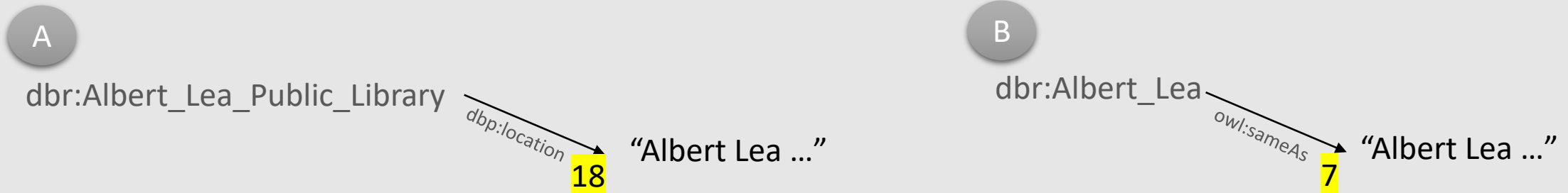
Q: “Albert Lea”



Proposition 2 [3][4]: use the context

Challenge 2: Cohesion

Q: “Albert Lea”

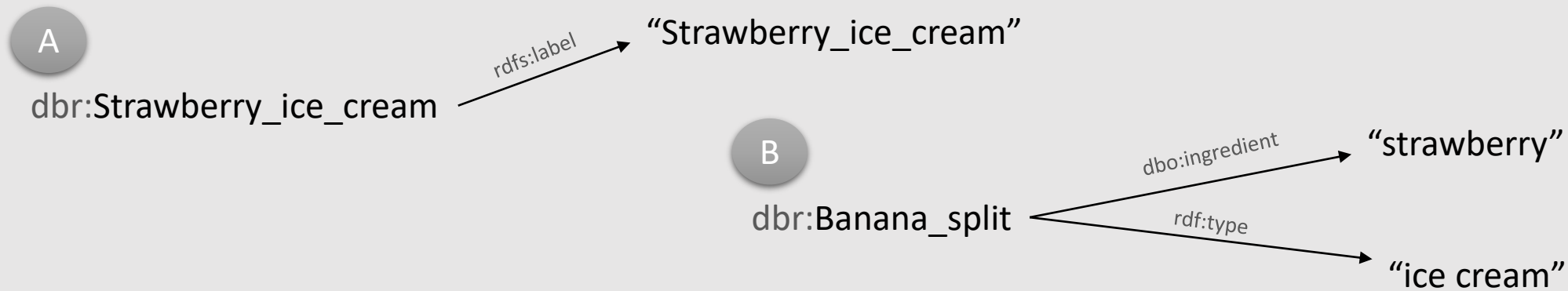


Proposition 2 [3][4]: use the context

Does **NOT** satisfy: Either A and B comply with preposition 2

Challenge 2: Cohesion

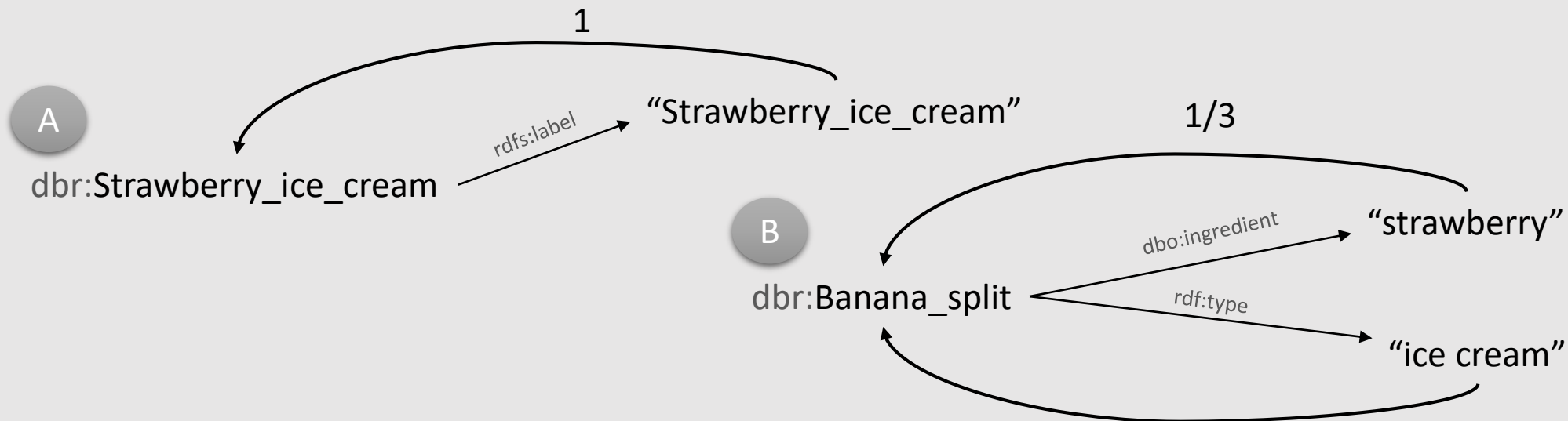
Q: “strawberry ice cream”



Proposition 3 [5][6]: Max Similarity (Levenshtein and Jaccard)

Challenge 2: Cohesion

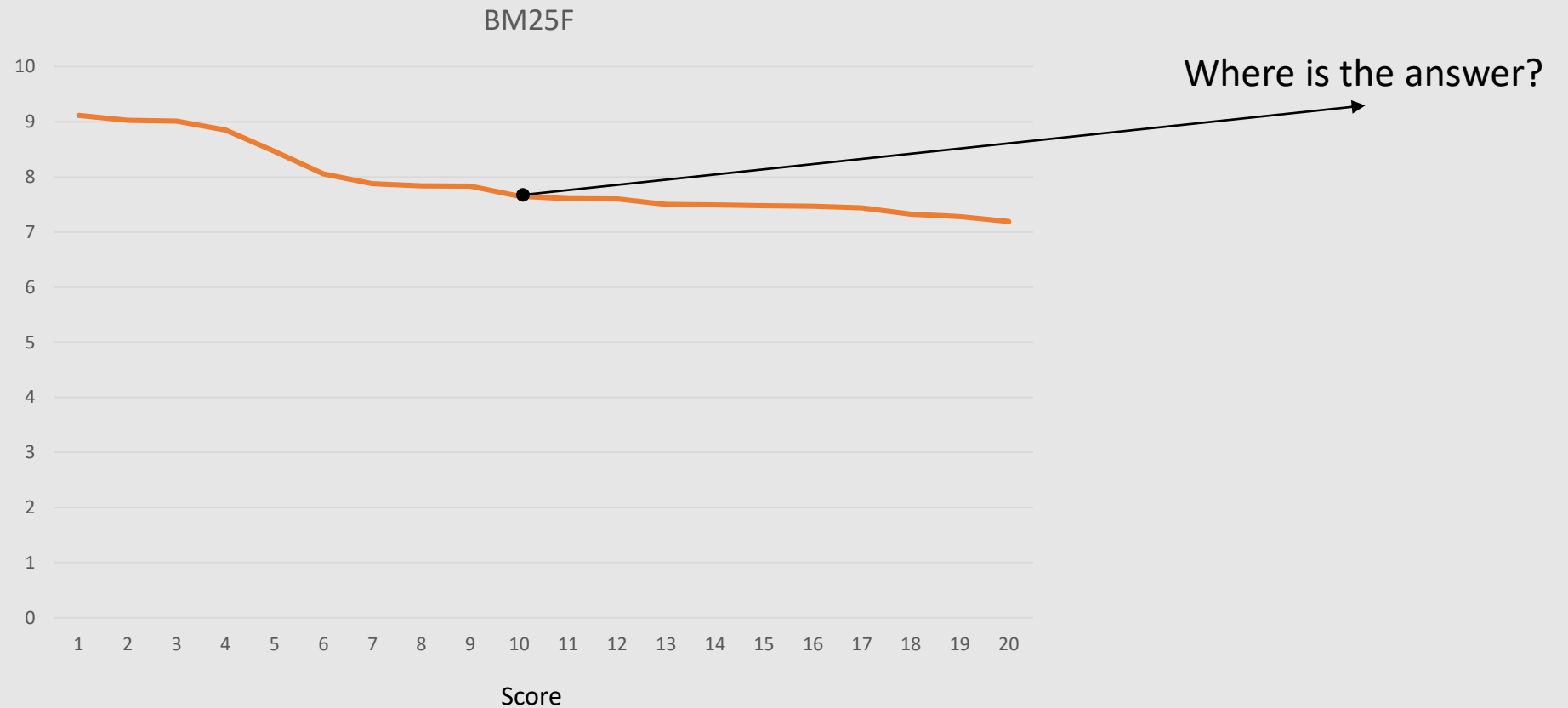
Q: "strawberry ice cream"



Proposition 3 [5][6]: Max Similarity (Levenshtein and Jaccard)

Does **NOT** satisfy: Either A and B comply with preposition 3

Challenge 3: Scoring Function



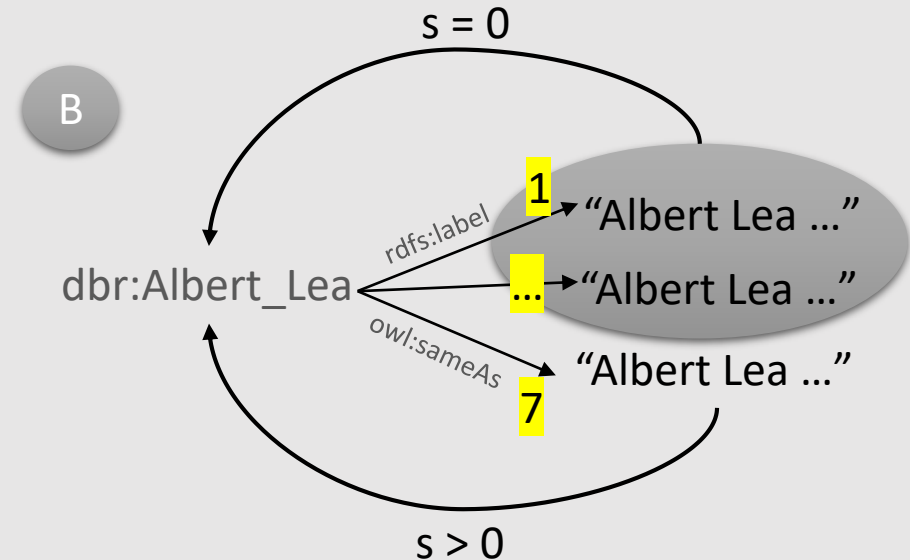
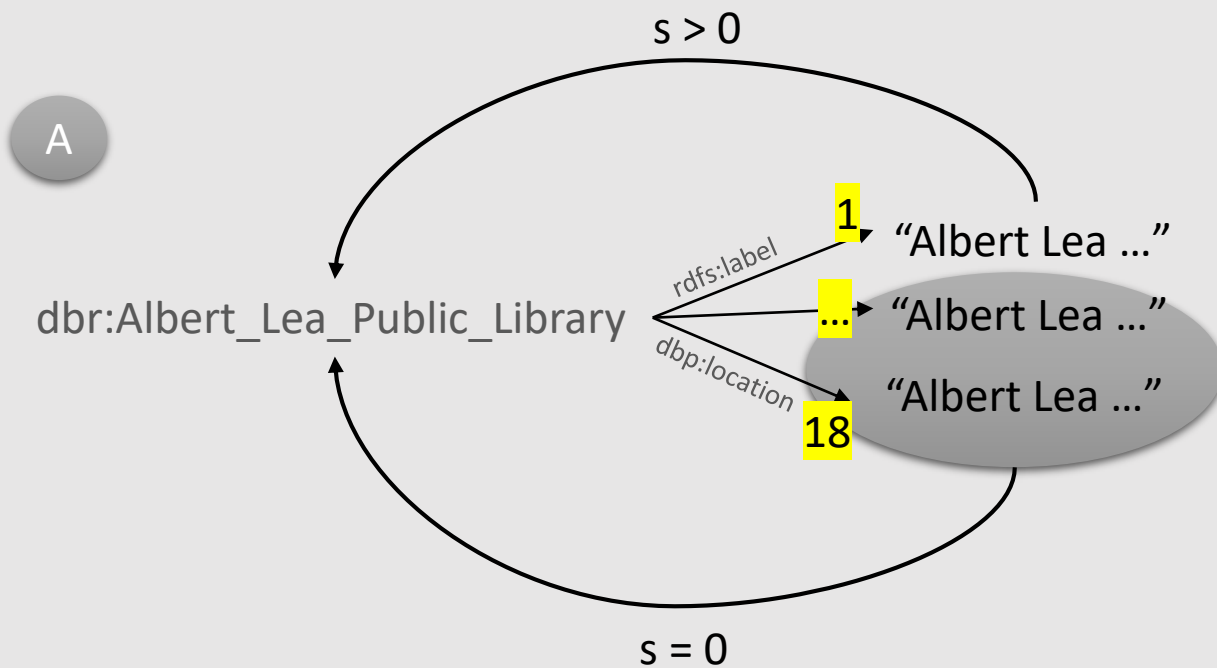
Approach

Can we do better?

Conditional Spread Activation

Q: "Albert Lea"

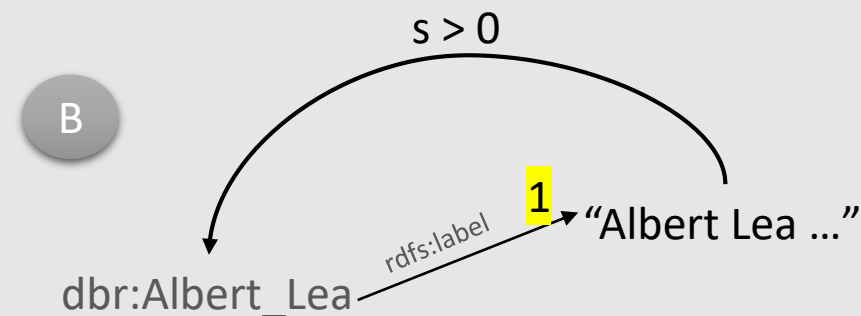
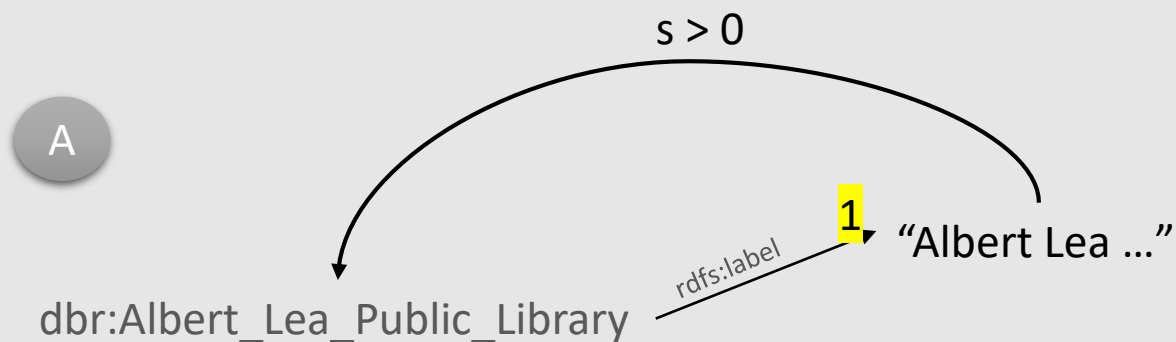
~~Challenge 1: Local Term Frequency~~



Field Weight

Q: "Albert Lea"

~~Challenge 1: Local Term Frequency~~
Challenge 2: Cohesion



$$\text{weight}(p_{\text{type}}) > \text{weight}(p_{\text{label}}) > \text{weight}(p_{\text{others}})$$

Field Weight

Q: "Albert Lea"

~~Challenge 1: Local Term Frequency~~
Challenge 2: Cohesion

A



B

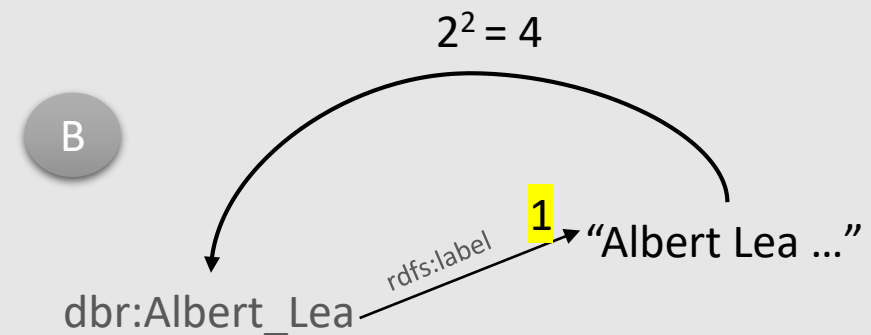
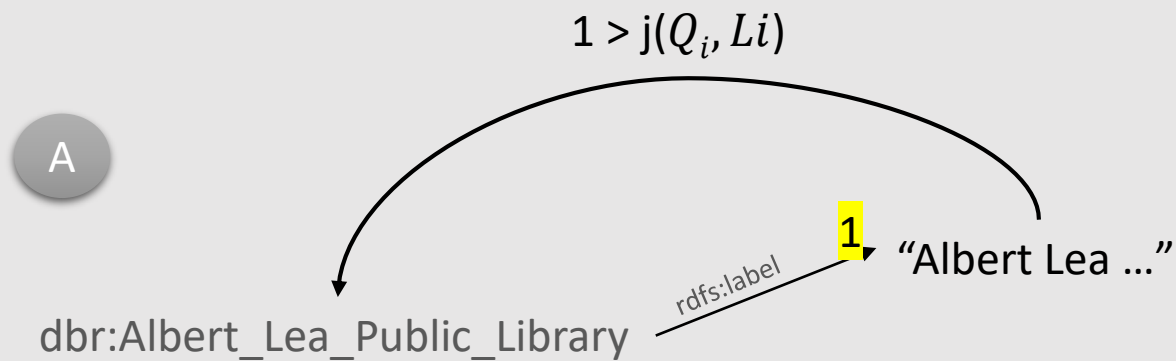


$$S(Q_i, L_i) = \sum Q_i L_i \text{ if } \sum Q_i L_i = 1$$

Field Weight

Q: "Albert Lea"

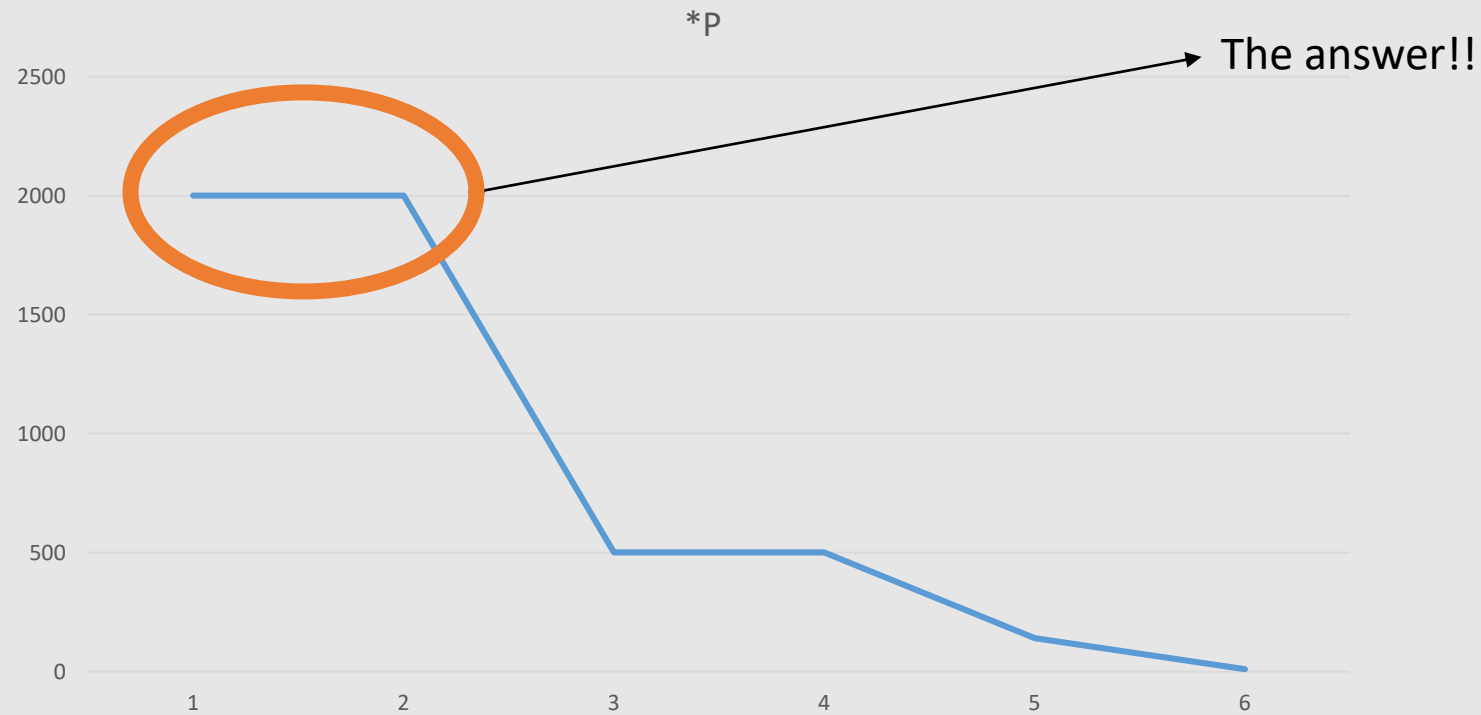
~~Challenge 1: Local Term Frequency~~
~~Challenge 2: Cohesion~~



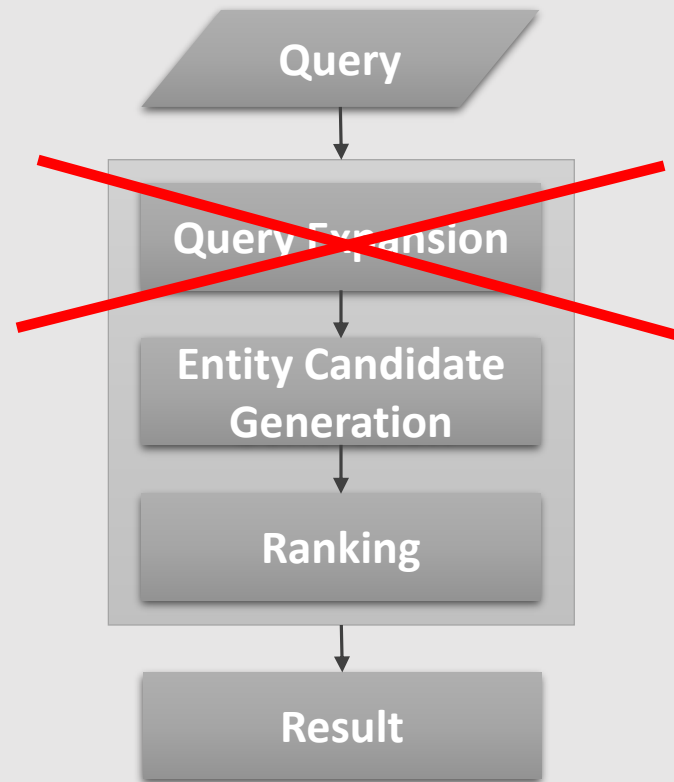
$$S(Q_i, L_i) = \sum Q_i^{L_i} \text{ if } \sum Q_i L_i = 1$$

Field Weight

- ~~Challenge 1: Local Term Frequency~~
- ~~Challenge 2: Cohesion~~
- ~~Challenge 3: Scoring Function~~



Pipeline



Benchmark

Benchmark	#Queries
QALD-2 (DBpedia-Entity)	140
QALD-4	50

DBpedia-Entity

Approach	MAP
CACAO	0.2417
BM25-CA	0.1939
SDM+EL	0.1887
FSDM+EL	0.1719
BM25F-CA	0.1671
LTR	0.1629
LM+EL	0.1534
SDM	0.1533
LM	0.1424
FSDM	0.1403
MLM-CA	0.1273
PRMS	0.1103
BM25	0.1092
MLM-all	0.1062

Table 2: Mean Average Precision (**MAP**) achieved by different Entity Retrieval models on QALD-2 *DBpedia-Entity* benchmark data set.

Approach	P@10
CACAO	0.3057
BM25-CA	0.2527
LM+EL	0.2362
SDM+EL	0.2249
LM	0.2144
FSDM+EL	0.2113
BM25F-CA	0.2053
FSDM	0.2000
SDM	0.1883
PRMS	0.1871
MLM-CA	0.1844
MLM-all	0.1843
LTR	0.1732
BM25	0.0986

Table 3: Precision 10 (**P@10**) achieved by different Entity Retrieval models on QALD-2 *DBpedia-Entity* benchmark data set.

QALD-4

Approach	P	R	F_1
CACAO+F	0.19	0.19	0.19
CACAO	0.11	0.11	0.11
CACAO _{P65}	0.09	0.09	0.09
Levenshtein _b	0.04	0.05	0.04
BM25F [2]	0.03	0.03	0.03
Jaccard _b	0.01	0.04	0.01
Levenshtein _a	0.00	0.00	0.00
Jaccard _a	0.00	0.00	0.00

Table 4: *Precision*, *Recall* and F_1 -*measure* achieved by different ER approaches on QALD-4 benchmark data set.

Approach	P	R	F_1
CACAO _{P65}	1	1	1
CACAO	0.90	0.90	0.90
MAG [24]	0.80	0.80	0.80
DBpedia Spotlight [7]	0.70	0.70	0.70
Levenshtein _b	0.60	0.60	0.60
Jaccard _b	0.60	0.60	0.60
AGDISTIS [36]	0.30	0.30	0.30
BM25F [2]	0.30	0.30	0.30
Levenshtein _a	0.00	0.00	0.00
Jaccard _a	0.00	0.00	0.00

Table 5: *Precision*, *Recall* and F_1 -*measure* achieved by different EL approaches on QALD 4 benchmark data set.

Conclusion & Future Works

- We have shown a **promising** algorithm that outperform state-of-the-art ER engines on standard benchmarks by addressing:
 - Local Term Frequency;
 - Cohesion, and;
 - Score Function problems
- It is still too soon to trace conclusions (law of small numbers) achieves that
- Improve the algorithm runtime so it can be used at scale

Acknowledgements



LiberAI



@theLiberAI

@HTWKLeipzig

@edgardmarx

@akswgroup

Acknowledgements

Thank you!!!

- * 100+ GitHub stars
- * 40+ forks

@theLiberAI @HTWKLeipzig @edgardmarx @akswgroup