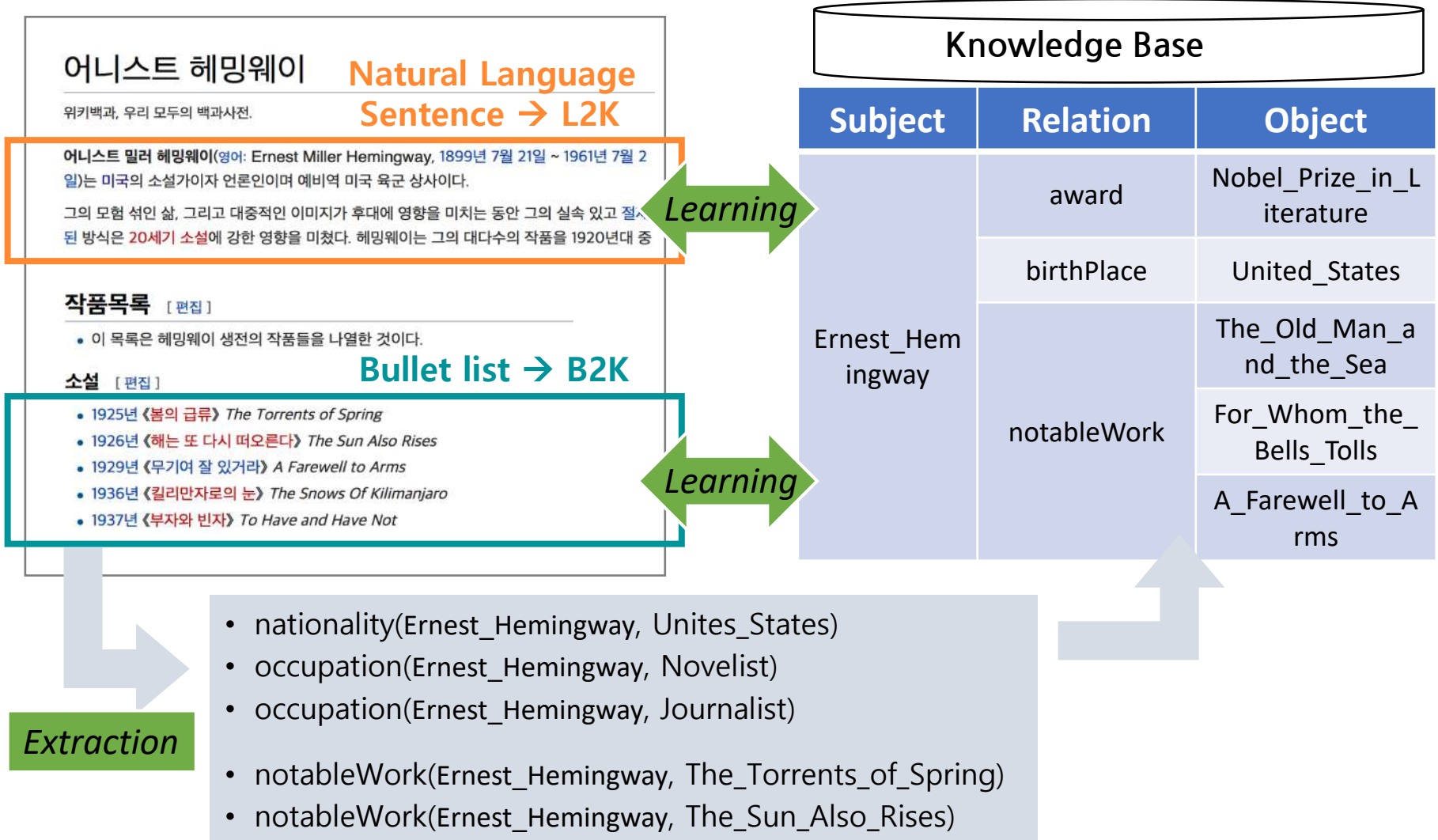


# Enriching DBpedia by Knowledge Base Population and Dark Entity Resolution

Key-Sun Choi  
[kschoi@kaist.ac.kr](mailto:kschoi@kaist.ac.kr)

School of Computing  
Korea Advanced Institute of Science & Technology  
([www.kaist.edu](http://www.kaist.edu))

# Overview of Fact Triplet Extraction



# Goal and Challenges

---

- Goal
  - Constructing an iterative learning platform based on Distant Supervision to improve knowledge learning quality
- Challenges
  - 1) How to build a good quality & quantity initial knowledge base?
  - 2) How to make a reliable knowledge base population system?

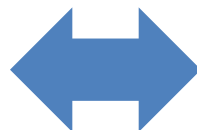
# Machine need to learn from text

---

- Corpus
  - Wikipedia, News, Blog, Twitter, Dialog, ...
- Knowledge Base
  - DBpedia, Freebase, Wikidata, ...

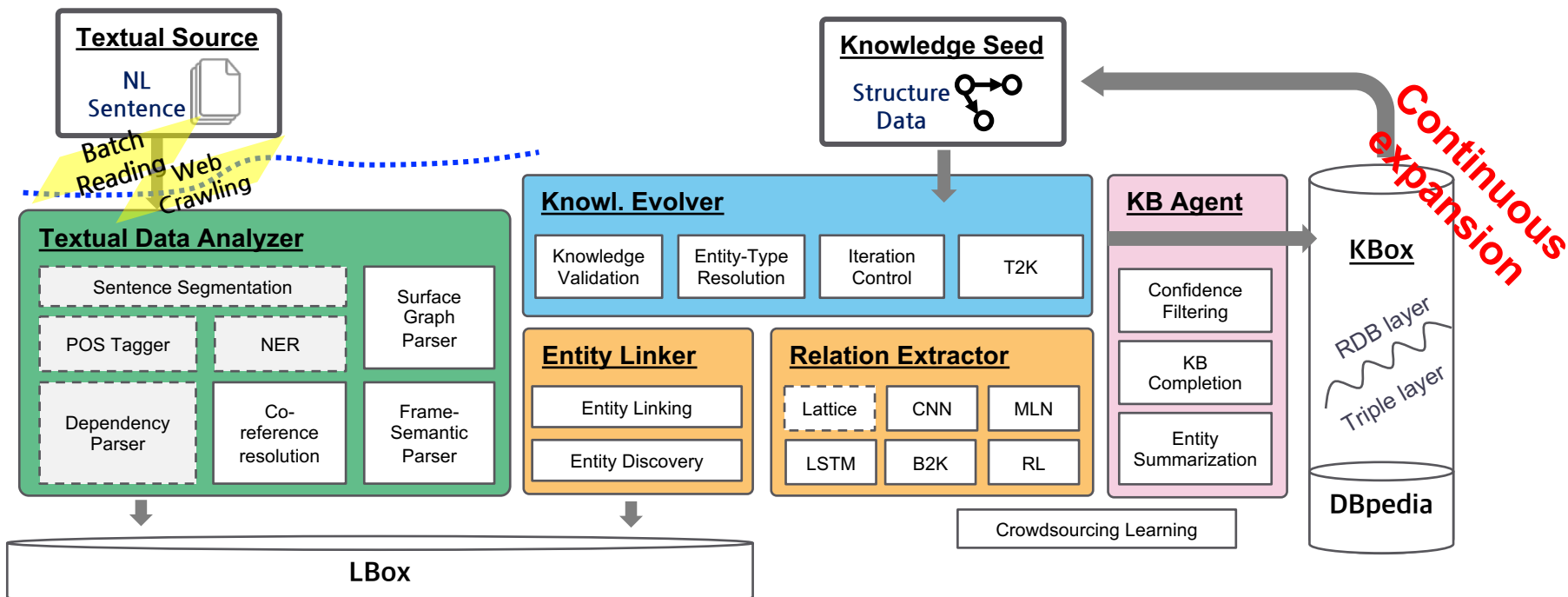


WIKIPEDIA  
The Free Encyclopedia



# KBox : Extended DBpedia

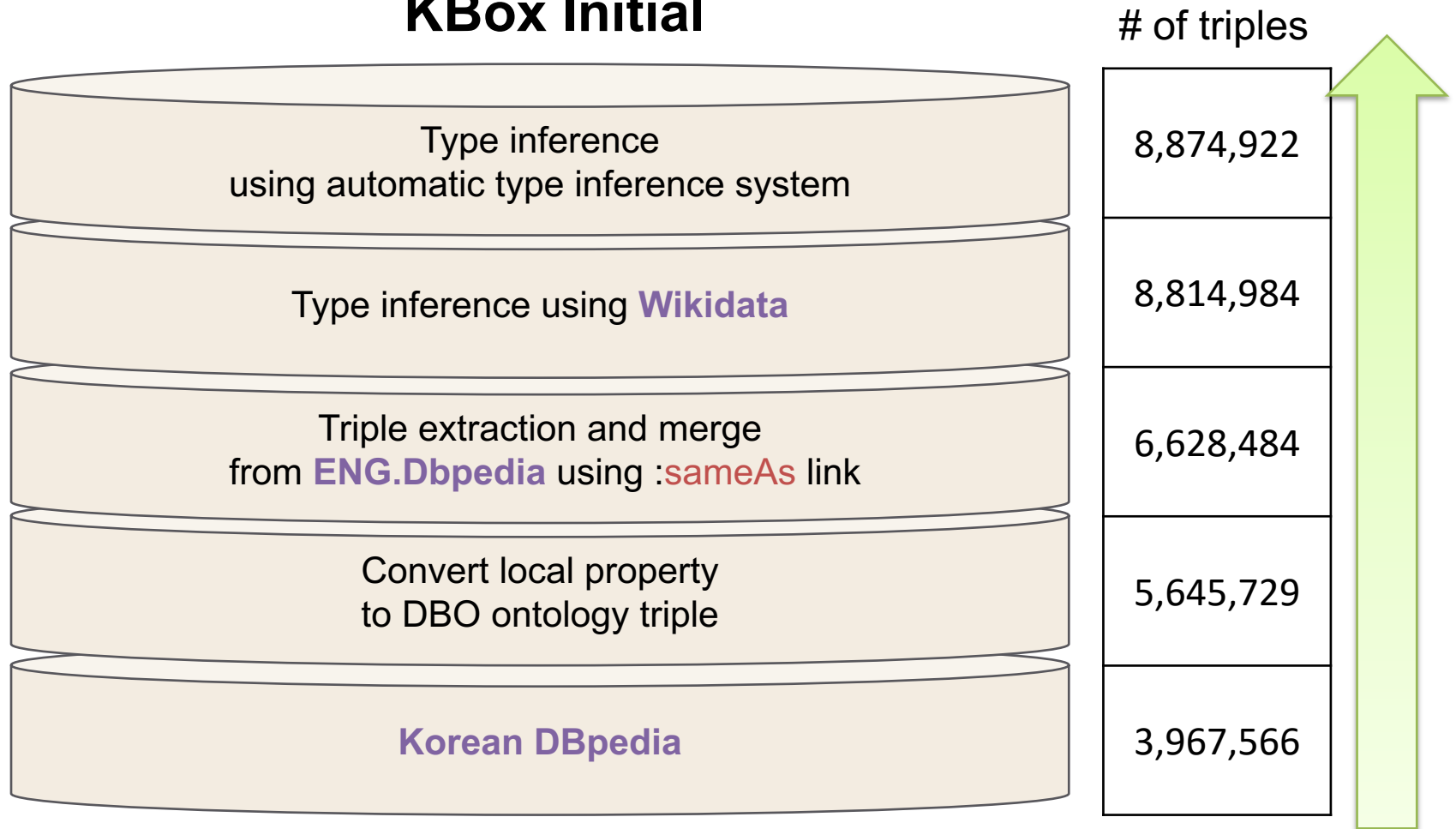
- KBox is a new KB that expands Korean DBpedia.
  - follows the DBpedia ontology schema.
  - continually expands much of the knowledge extracted from text using the Korean KBP system.



# Kbox Initialization, kbox.kaist.ac.kr

## - Korean enriched DBPedia

### KBox Initial



# Kbox Initialization

## - How to build a good quality & quantity initial knowledge base

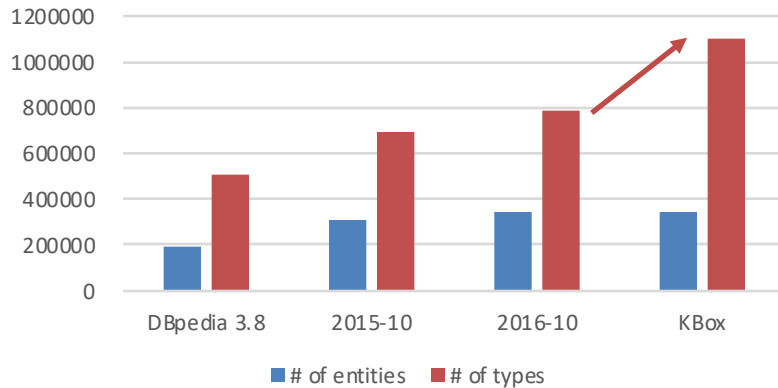
---

- 2-layer Storage
  - RDB (MySQL) : All the information about all the KBox triples, such as triple score, extraction module, source sentence.
  - Triple Store (Stardog, Virtuoso) : Reliable triples
    - 1) The initial triples extracted from the DBpedia and Wikidata
      - 1-1) Convert local property (prop-ko) to DBO property (출생지 → birthPlace)
      - 1-2) Type inference using :sameAs link from EN.DBpedia
      - 1-3) Triple generation using :sameAs link from EN.Dbpedia and Wikidata
    - 2) Automatically extracted triples using the Korean KB Population with a confidence score above 0.9

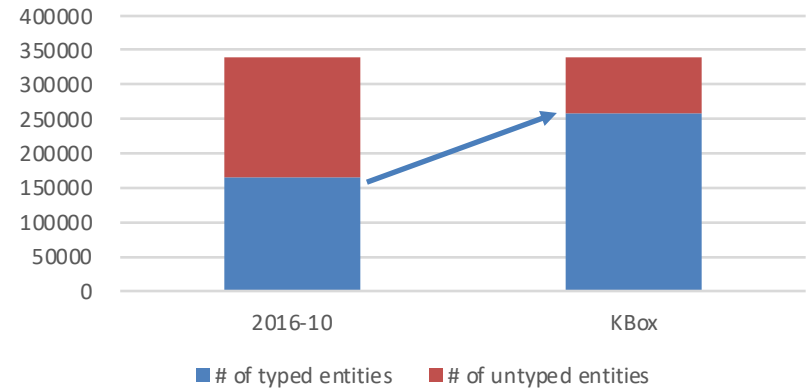
# KBox Statistics

(How to build a good quality & quantity initial knowledge base?)

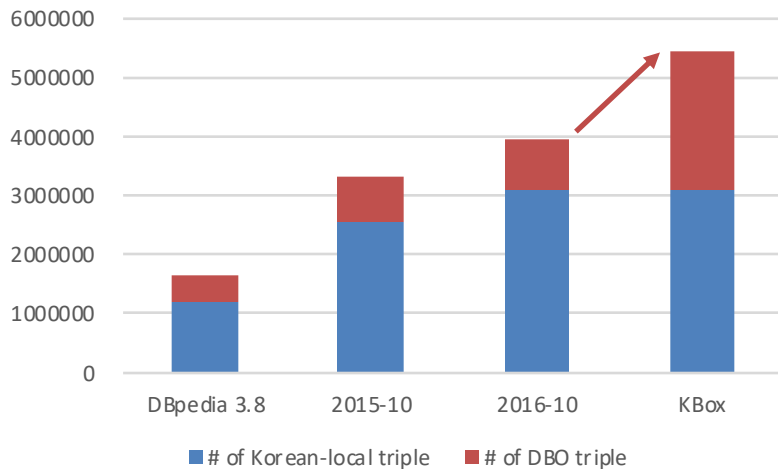
Enrichment in Entity and Types



Enrichment in Typed entities



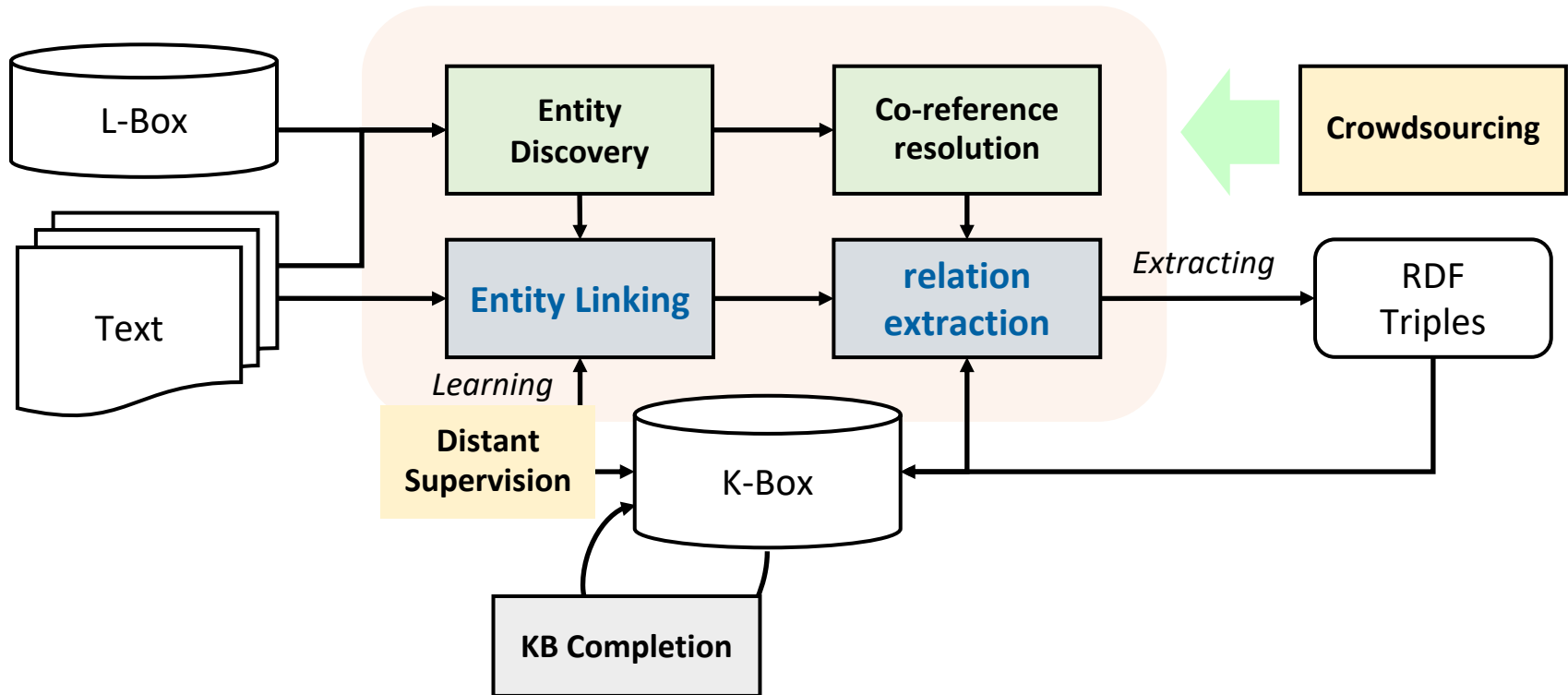
Enrichment in Triples



Version	# of Korean local triple	# of DBO triple	# of relational triples
DBpedia 3.8	1,219,355	432,600	1,651,649
2015-10	2,569,432	738,599	3,308,031
2016-10	3,077,297	890,269	3,967,383
<b>KBox</b>	<b>3,077,297</b>	<b>2,350,292</b>	<b>5,426,857</b>



# Essential Tasks for KB Population



# Example of Knowledge Learning / Extracting

## Extracted Knowledge

• Agent Person Artist

- 마사\_겔혼 (Martha\_Gellhorn) ★
  - spouse **어니스트\_헤밍웨이** (Ernest\_Hemingway)
  - is spouse of **어니스트\_헤밍웨이** (Ernest\_Hemingway)
- **어니스트\_헤밍웨이** (Ernest\_Hemingway) ★
  - award **노벨\_문학상** (Nobel\_Prize\_in\_Literature)
  - birthPlace **미국** (United\_States)
  - birthPlace **일리노이\_주** (Illinois)
  - deathPlace **미국** (United\_States)
  - deathPlace **아이다호\_주** (Idaho)
  - notableWork **노인과\_바다** (The\_Old\_Man\_and\_the\_Sea)
  - notableWork **누구를\_위하여\_좋은\_올리나** (For\_Whom\_the\_Bell\_Tolls)
  - notableWork **무기여\_잘\_있거라** (A\_Farewell\_to\_Arms)
  - occupation **소설가** (Novelist)

**어니스트 헤밍웨이**는 미국의 **소설가**이자 **저널리스트**이다. **1854년** **노벨 문학상**을 수상하였다. **헤밍웨이**는 **1899년 7월 21일** **일리노이주**에서 태어났다. **헤밍웨이**는 폴린 **파이퍼**와 이혼한 뒤 **마사 겔혼**과 재혼하였다. **19** **아이다호 주**에서 **업** 총으로 **62세**의 나이에 자살했다.

Type: Artist  
String: 헤밍웨이  
Kor\_Entity: 어니스트 헤밍웨이  
Eng\_Entity: Ernest\_Hemingway

- **1926년** 《**해는 또다시 떠오른다**》 **The Sun** Also Rises.
- **1929년** 《**무기여 잘 있거라**》 A Farewell to Arms.
- **1940년** 《**누구를 위하여 좋은 올리나**》**For Whom the** Bell Tolls.
- **1950년** 《**강 건너 숲속으로**》 Across the River and Into the Trees.
- **1952년** 《**노인과 바다**》 The Old Man and the Sea.

Relation Extraction ← Entity Linking

# Problem and solutions in KB Population

How to make a reliable knowledge base population system?

---

## Problems:

Noise data on distant supervision

Hard to understand a long and complex sentence

Difficulty in interpreting ontology relation for all collected sentences

Incompleteness of knowledge base

## Solutions:

Enhanced noise filtering using crowdsourcing and reinforcement learning → **Increase precision**

Diversify relation extraction model, Co-reference resolution and Zero anaphora detection

Full-text graph generation → **Increase recall**

Dark entity, Iterative Learning, Knowledge Completion and Summarization

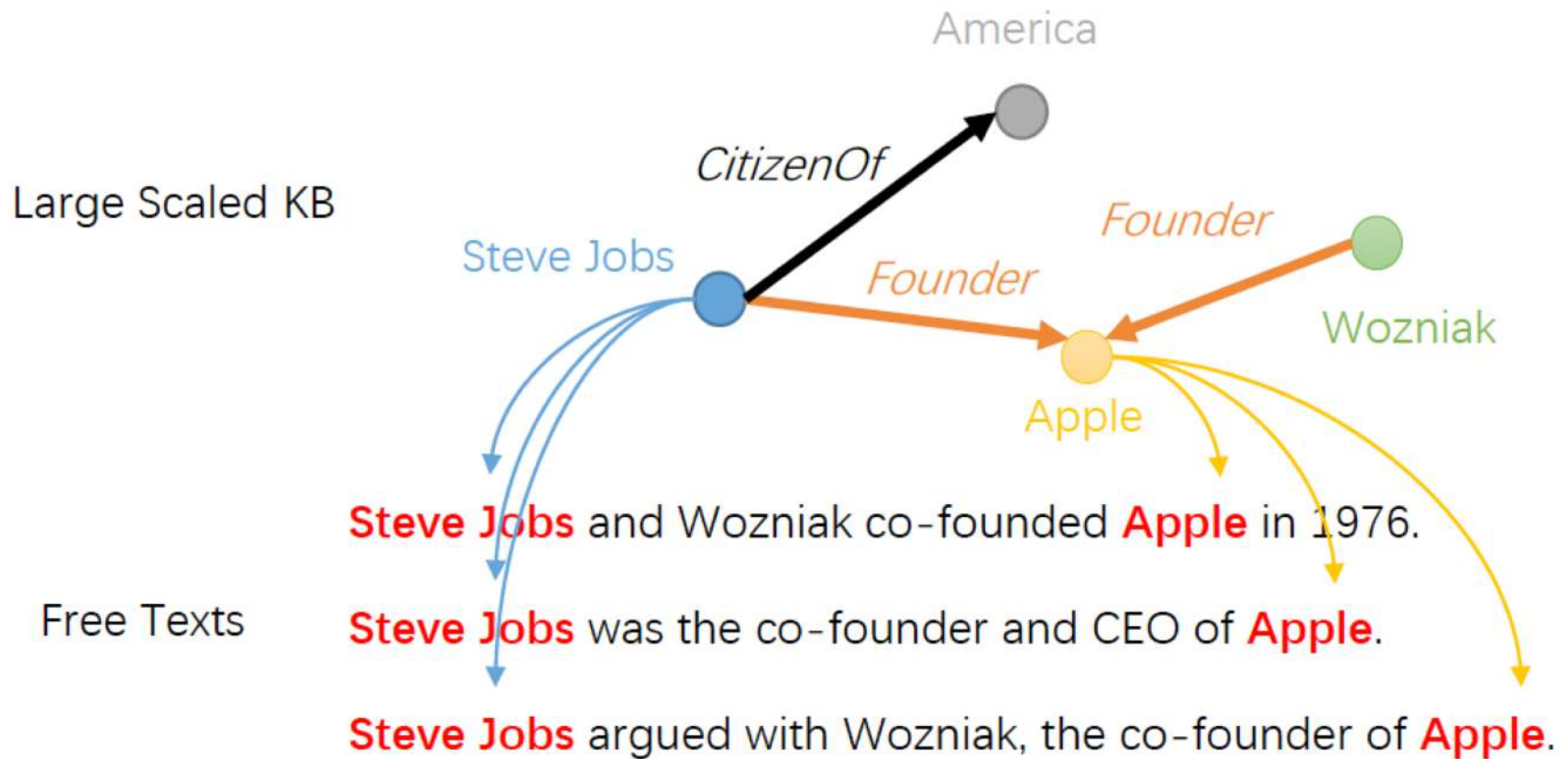
# Distant Supervision (Relation Extraction)

---

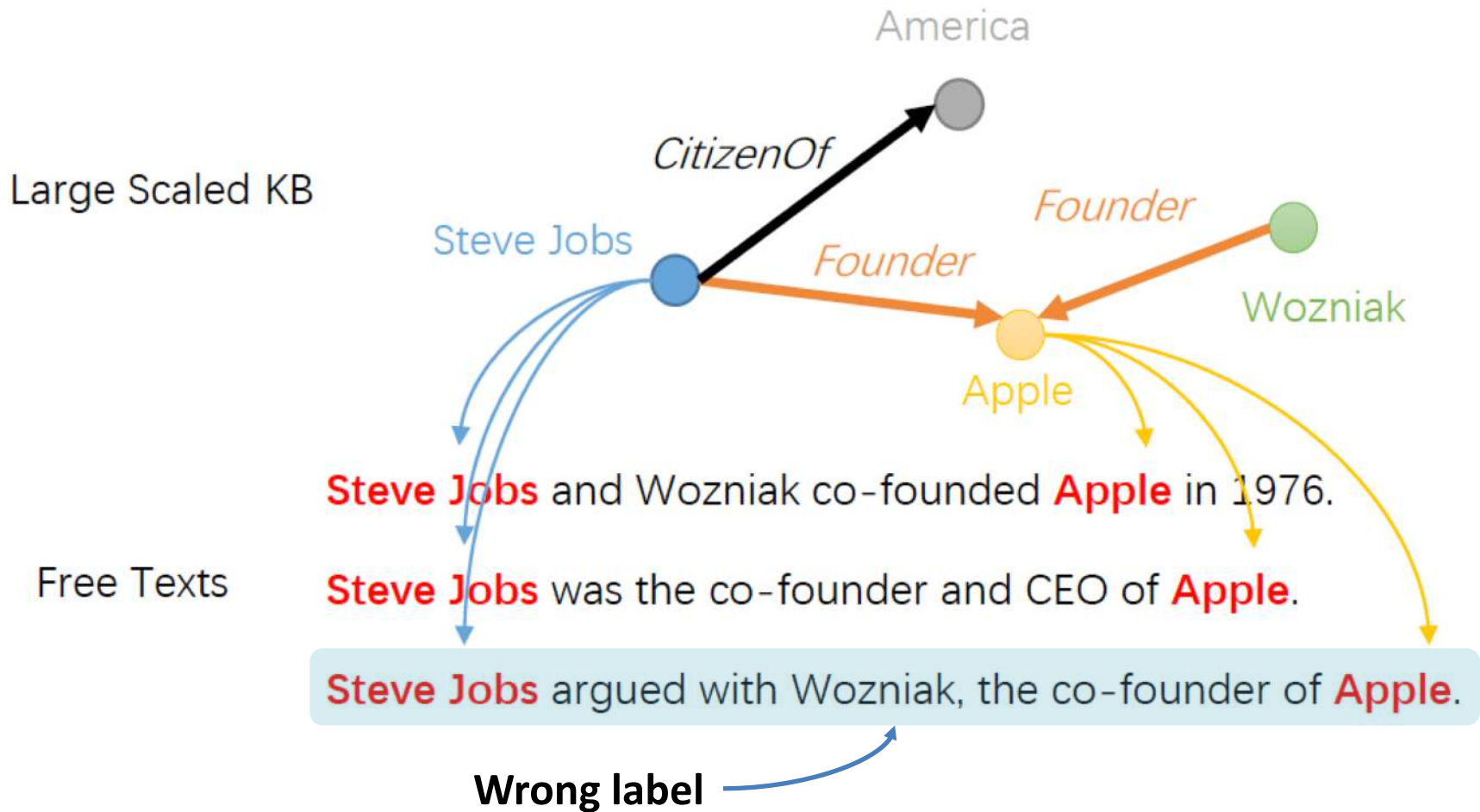
- Distant Supervision (Mintz et al. 2009)
  - For a triple fact  $r(e1, e2)$  in a KB, all sentences that mention both entities  $e1$  and  $e2$  are aligned with relation  $r$ .

Sentences	Relation
1. <b>Steve Jobs</b> and Wozniak co-founded <b>Apple</b> in 1976.	<i>Founder</i>
2. <b>Michael Jordan</b> is an American retired professional <b>basketball player</b> .	<i>Career</i>
3. <b>Washington D.C.</b> is the capital of <b>United states</b> .	<i>CapitalOf</i>
.....	.....

# Automatic labeling between KB and sentence



# Problem : Noise data



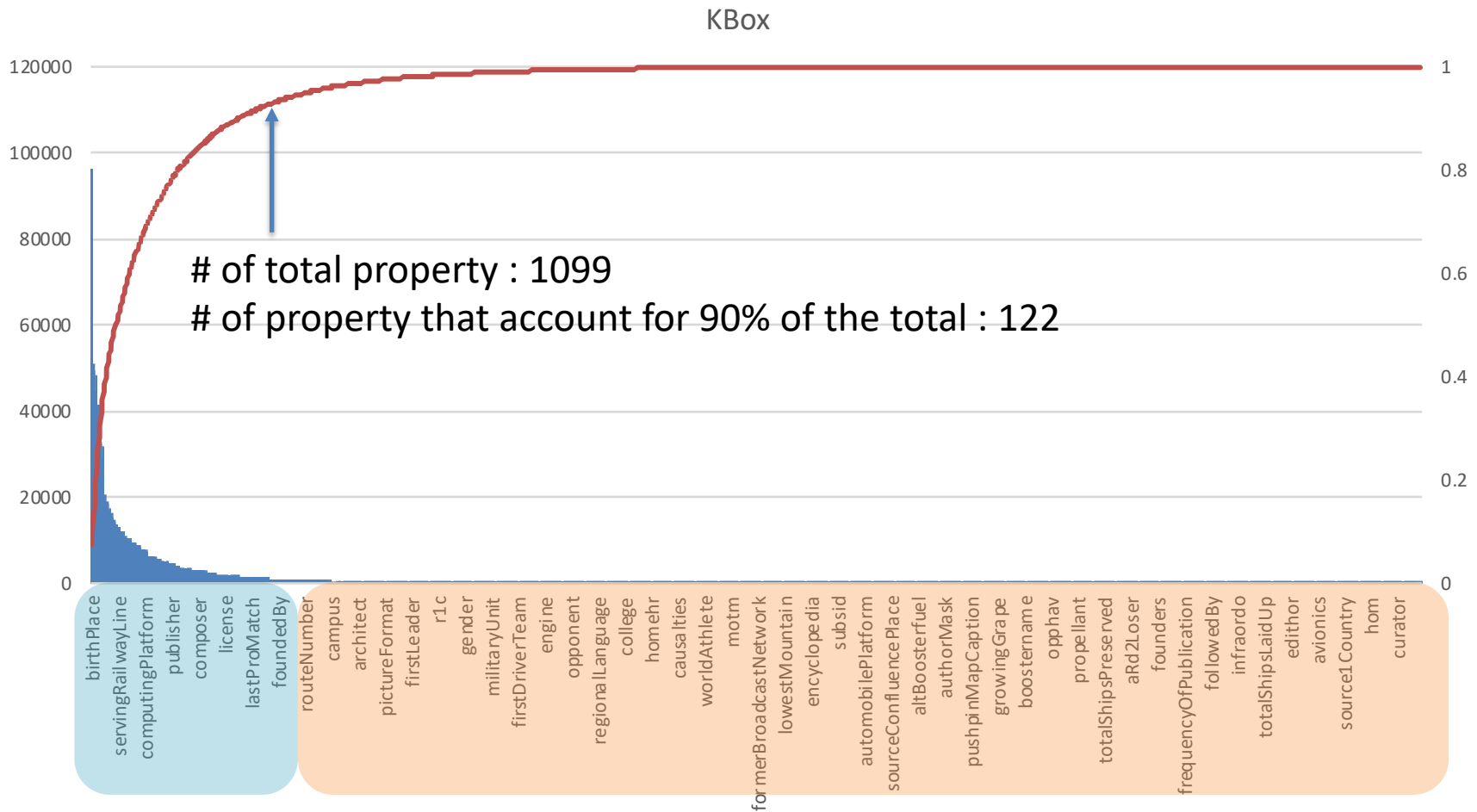
# Problem : Noise data

---

- Average error rate reported in this paper<sup>1</sup>
  - 74.1% (DS data from Wikipedia – YAGO)
  - 31.0% (DS data from NYT News – Freebase)
- **Intent** : Up to 20% performance improvement by learning models with noise-free data (in our experiments)

1. Ru, C., Tang, J., Li, S., Xie, S., & Wang, T. (2018). Using semantic similarity to reduce wrong labels in distant supervision for relation extraction. Information Processing & Management.

# Sentence collecting problem and Imbalanced-labeled data



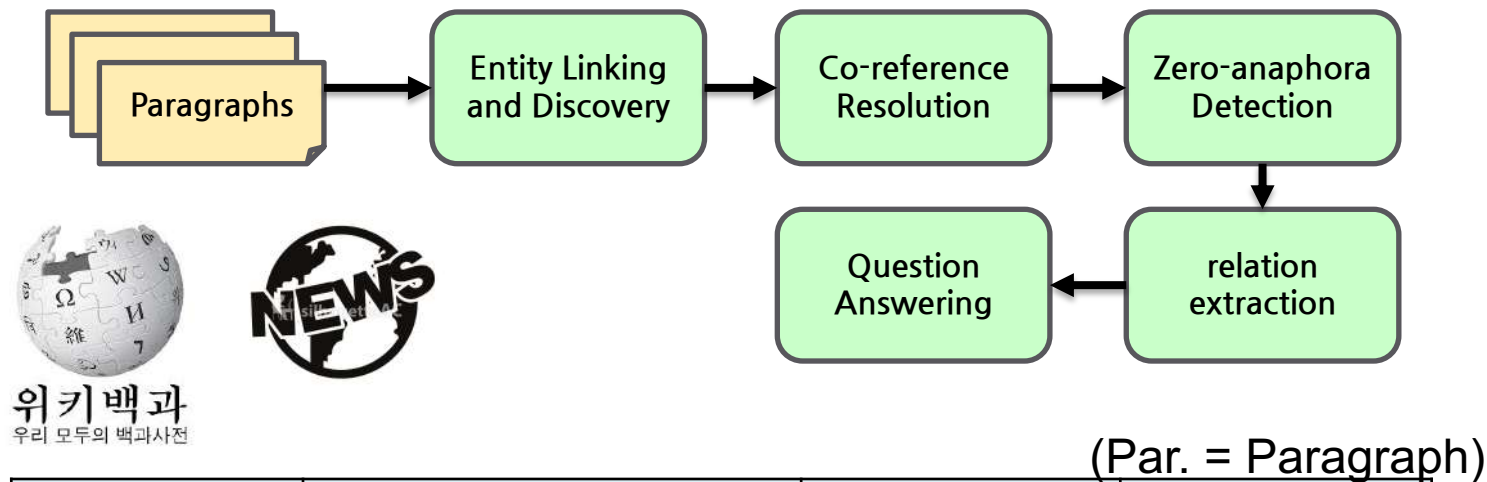
Deep-learning Approach

Rule-based Approach



# Solution 1 : Crowdsourcing

- Goal : Improving the performance of knowledge learning models with human-annotated dataset on 5 tasks



Corpus	Target Tasks	Now	Goal
Wikipedia	Task 1-5	500 Par.	40K Par.
News	Task 4 relation extraction	5,482 Par.	60K Par.

# What makes increase Recall

---

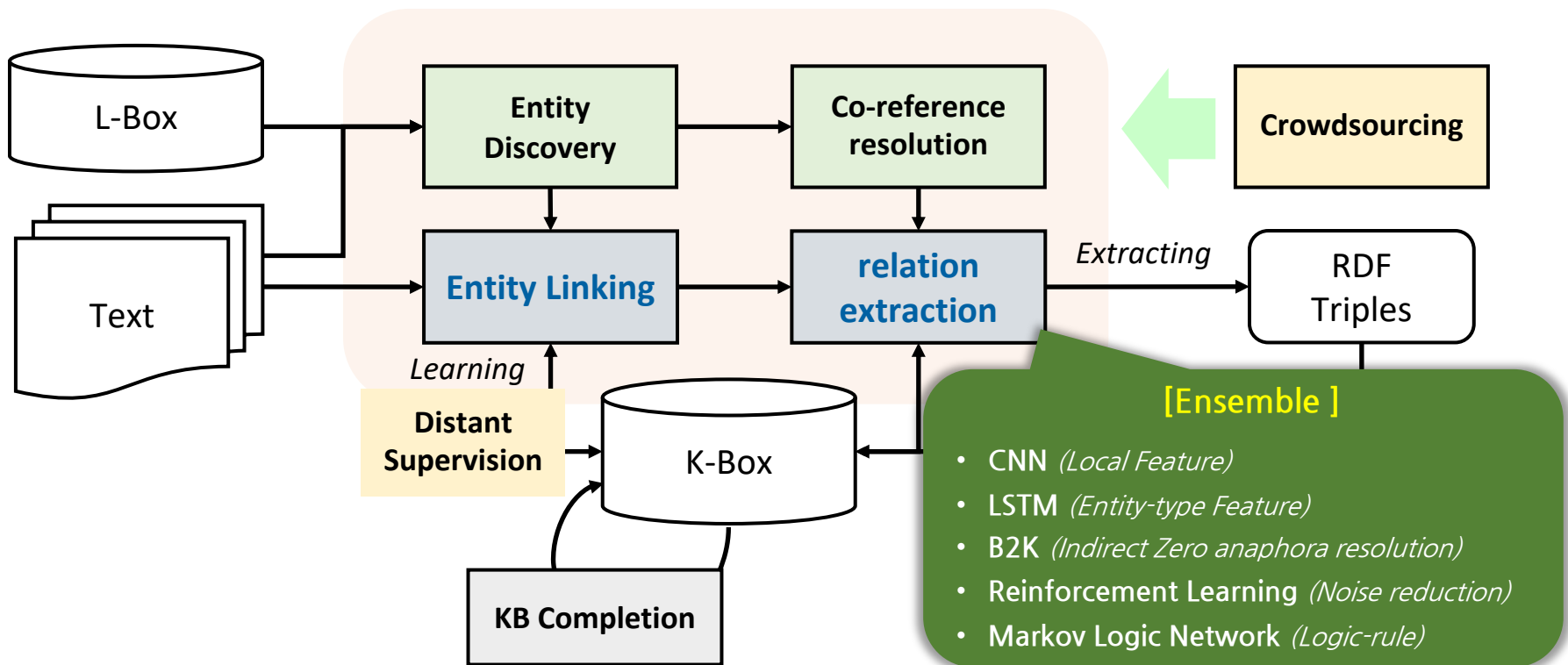
- Entity Discovery
  - Finding new dark entities that can be registered to the KB
- Co-reference resolution
  - **Intent** : Up to 21.6% improved knowledge extraction coverage using Entity Discovery and Co-reference resolution

## Crowdsourced Gold Set Analysis Result

# of sentence	Sentence with no or only one entity. (Exception from knowledge extraction target)	Sentences that are included in knowledge extraction targets with Entity Discover / Coreference Resolution	Coverage
1,435	445	96	21.6% (can be improved)

# Solution 2 : Ensemble

- Even a state-of-the-art relation extraction model shows low performance (F1-score, 40-50%).



# Full-text Graph

---

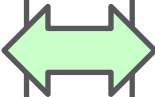
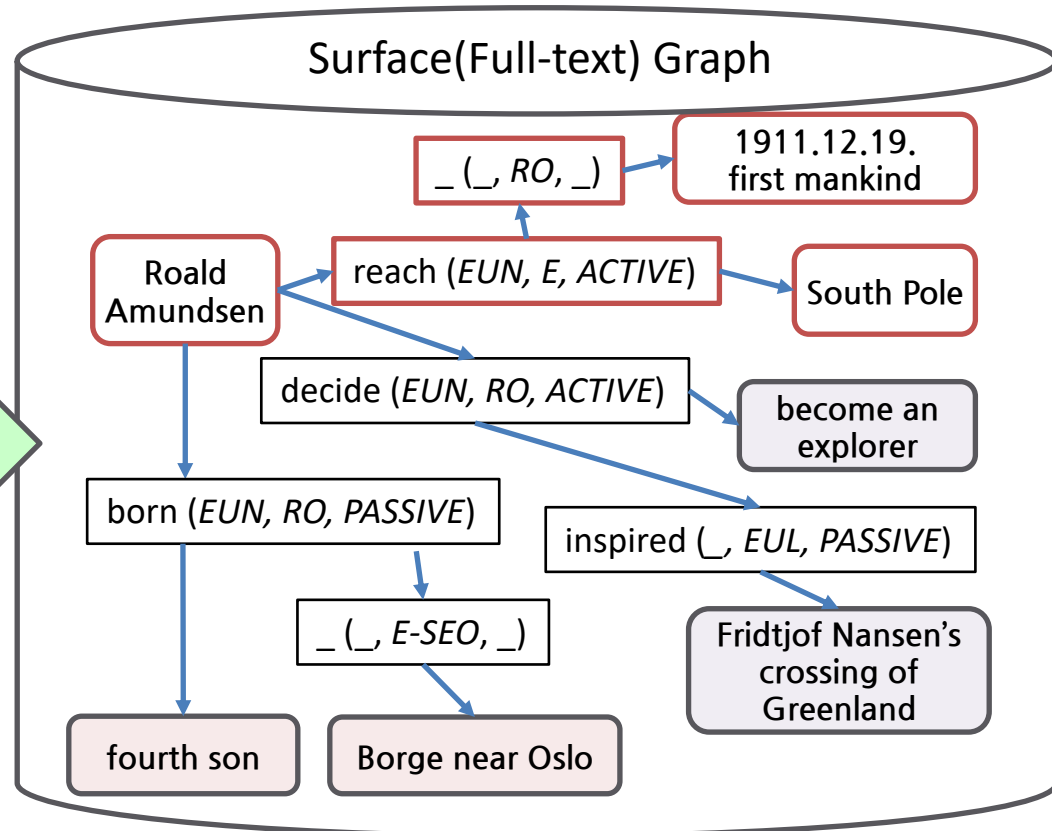
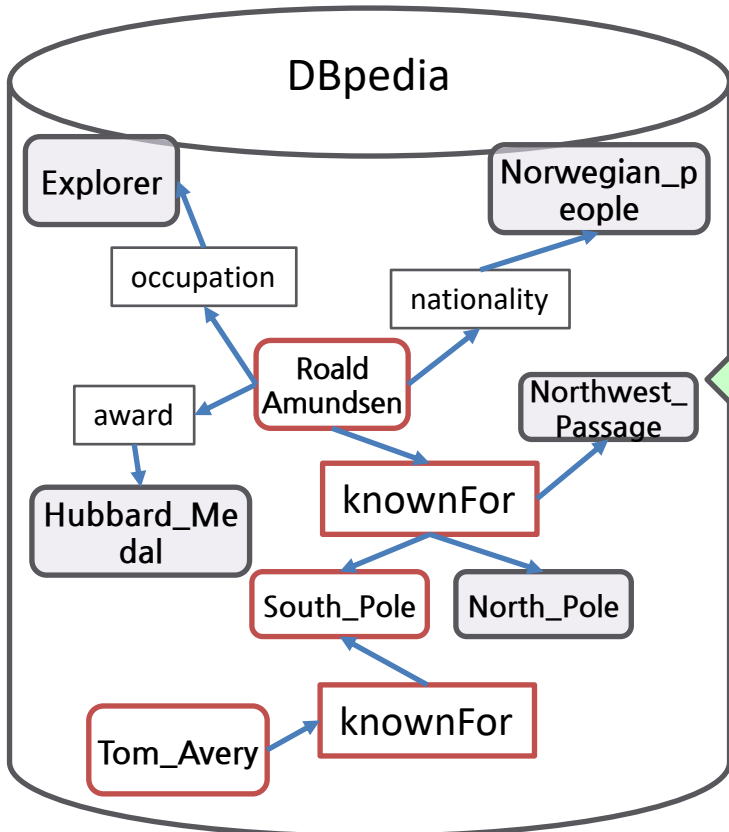
- Purpose
  - KBP/QA coverage expansion by combining ontological KB and full-text graph
- Lack of expression
  - There is a lot of knowledge that can not be expressed in an ontological relation.
- ISO/TC37/SC4 Proposal: Surface Knowledge graph with Linguistic Content

# Full-text Graph: Example of Application

Question : Who was the first person to reach the South Pole? (남극점에 최초로 도달한 사람은 누구인가?)

Answer : Roald Amundsen

*Hard to get an exact answer.*



# Application: [www.OKBQA.org](http://www.OKBQA.org) (Open Knowledge Base and Question Answering)

