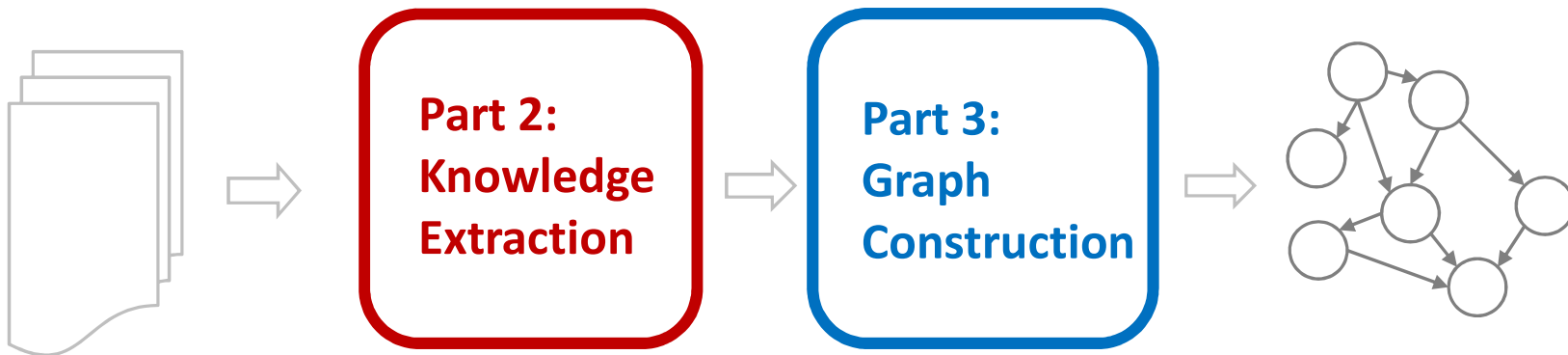


# Tutorial Overview

---

<https://kgtutorial.github.io>

Part 1: Knowledge Graphs



Part 4: Critical Analysis

# Tutorial Outline

---

1. Knowledge Graph Primer

[Jay]



2. Knowledge Extraction Primer

[Jay]



Coffee Break



3. Knowledge Graph Construction

a. Probabilistic Models

[Jay]



b. Embedding Techniques

[Sameer]



4. Critical Overview and Conclusion

[Sameer]



# Critical Overview

---

SUMMARY

SUCCESS STORIES

DATASETS, TASKS, SOFTWARES

EXCITING RESEARCH DIRECTIONS

# Critical Overview

---

SUMMARY

SUCCESS STORIES

DATASETS, TASKS, SOFTWARES

EXCITING RESEARCH DIRECTIONS

# Why do we need Knowledge graphs?

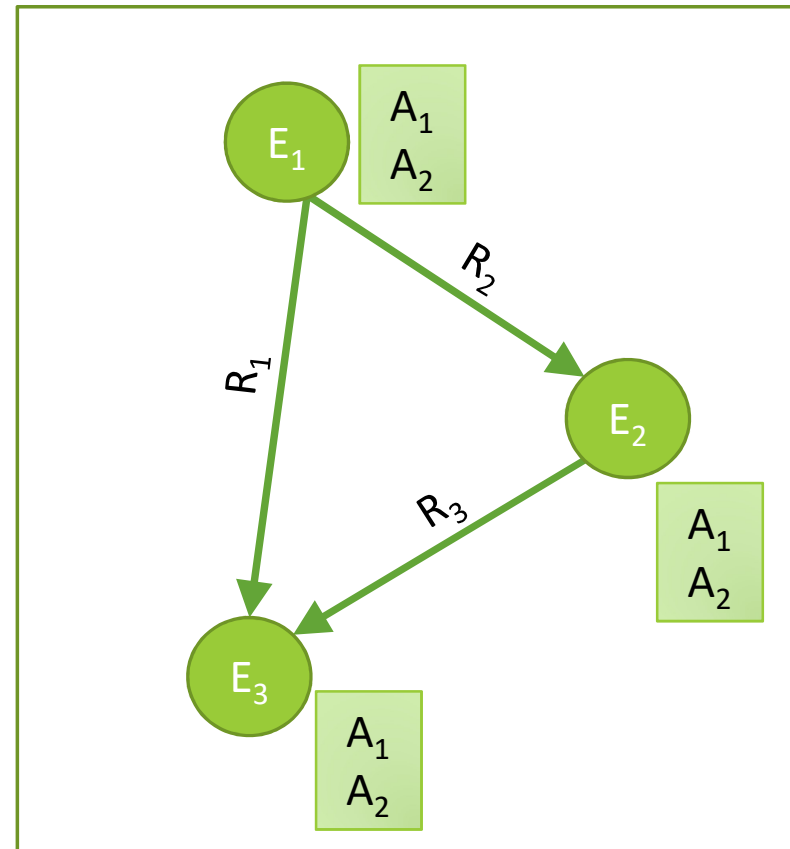
---

- Humans can explore large database in intuitive ways
- AI agents get access to human common sense knowledge

# Knowledge graph construction

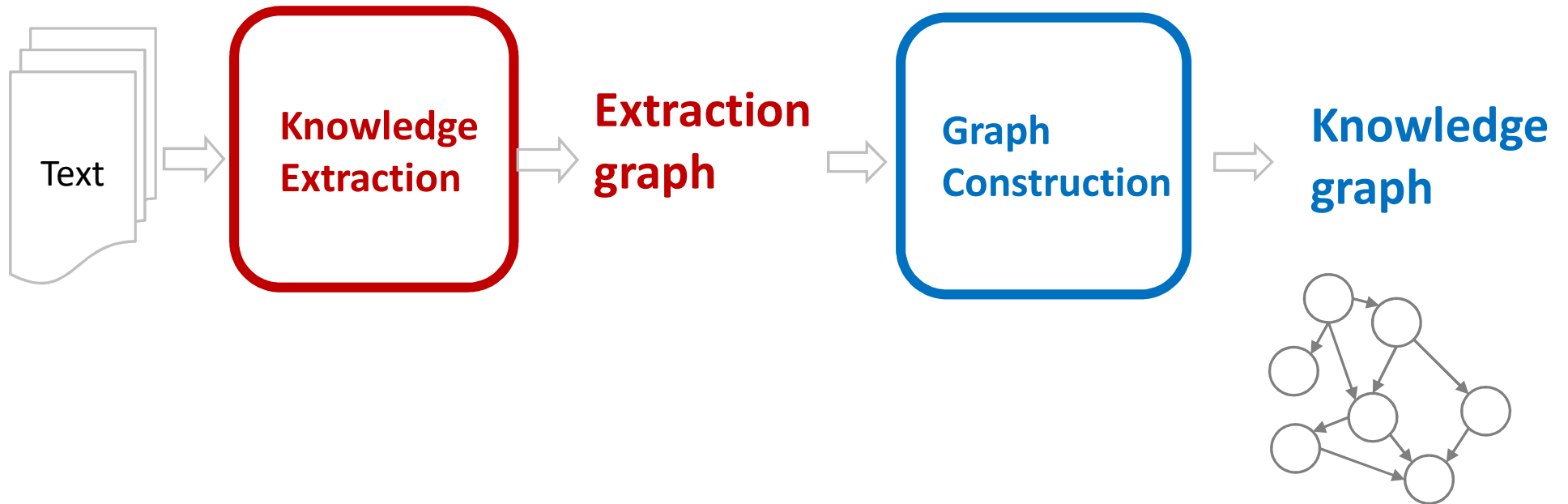
---

- **Who** are the entities (nodes) in the graph?
- **What** are their attributes and types (labels)?
- **How** are they related (edges)?



# Knowledge Graph Construction

---



# Two perspectives

---

	Extraction graph	Knowledge graph
Who are the entities? (nodes)	<ul style="list-style-type: none"><li>• Named Entity Recognition</li><li>• Entity Coreference</li></ul>	<ul style="list-style-type: none"><li>• Entity Linking</li><li>• Entity Resolution</li></ul>
What are their attributes? (labels)	<ul style="list-style-type: none"><li>• Entity Typing</li></ul>	<ul style="list-style-type: none"><li>• Collective classification</li></ul>
How are they related? (edges)	<ul style="list-style-type: none"><li>• Semantic role labeling</li><li>• Relation Extraction</li></ul>	<ul style="list-style-type: none"><li>• Link prediction</li></ul>

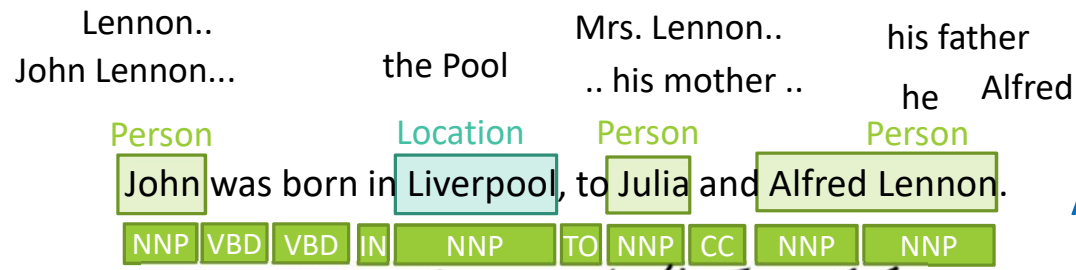


# Knowledge Extraction

John was born in Liverpool, to Julia and Alfred Lennon.

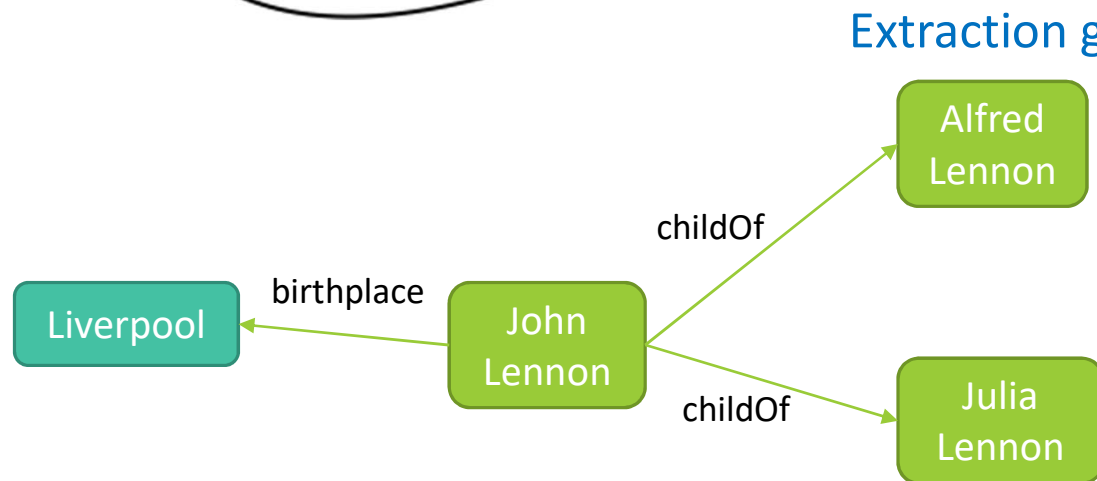
Text

NLP



Annotated text

Information  
Extraction



# Information Extraction

---

## Single extractor

Defining domain

Learning extractors

Scoring candidate facts



Supervised



Semi-supervised



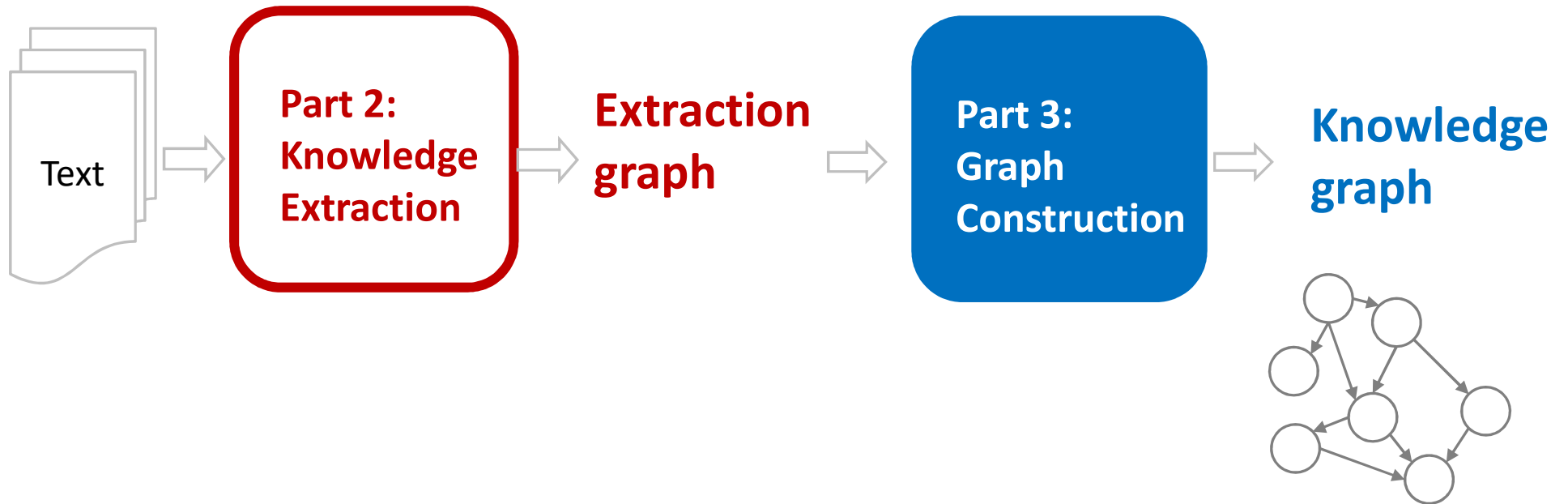
Unsupervised



Fusing multiple extractors

# Knowledge Graph Construction

---

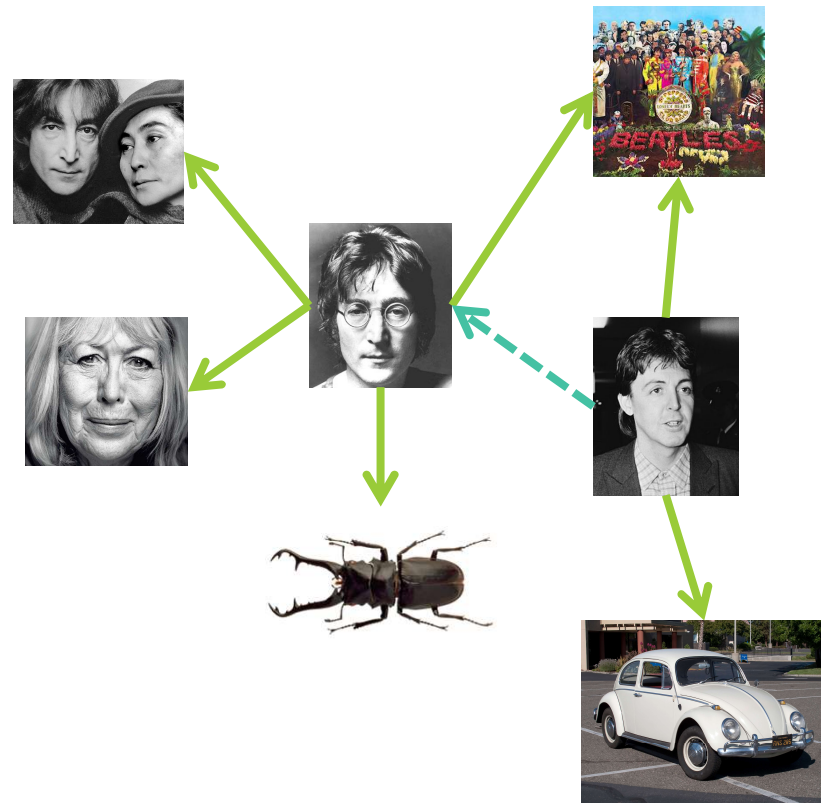


# Issues with Extraction Graph

---

Extracted knowledge could be:

- ambiguous
- incomplete
- inconsistent



# Two approaches for KG construction

---

PROBABILISTIC MODELS

EMBEDDING BASED MODELS

# Two approaches for KG construction

---

PROBABILISTIC MODELS

EMBEDDING BASED MODELS

# Two classes of Probabilistic Models

---

## GRAPHICAL MODEL BASED

- Possible facts in KG are variables
- Logical rules relate facts
- Probability  $\propto$  satisfied rules
- Universal-quantification

## RANDOM WALK BASED

- Possible facts posed as queries
- Random walks of the KG constitute “proofs”
- Probability  $\propto$  path lengths/transitions
- Local grounding

# Illustration of KG Identification

## Uncertain Extractions:

- .5: Lbl(Fab Four, novel)
- .7: Lbl(Fab Four, musician)
- .9: Lbl(Beatles, musician)
- .8: Rel(Beatles, AlbumArtist, Abbey Road)

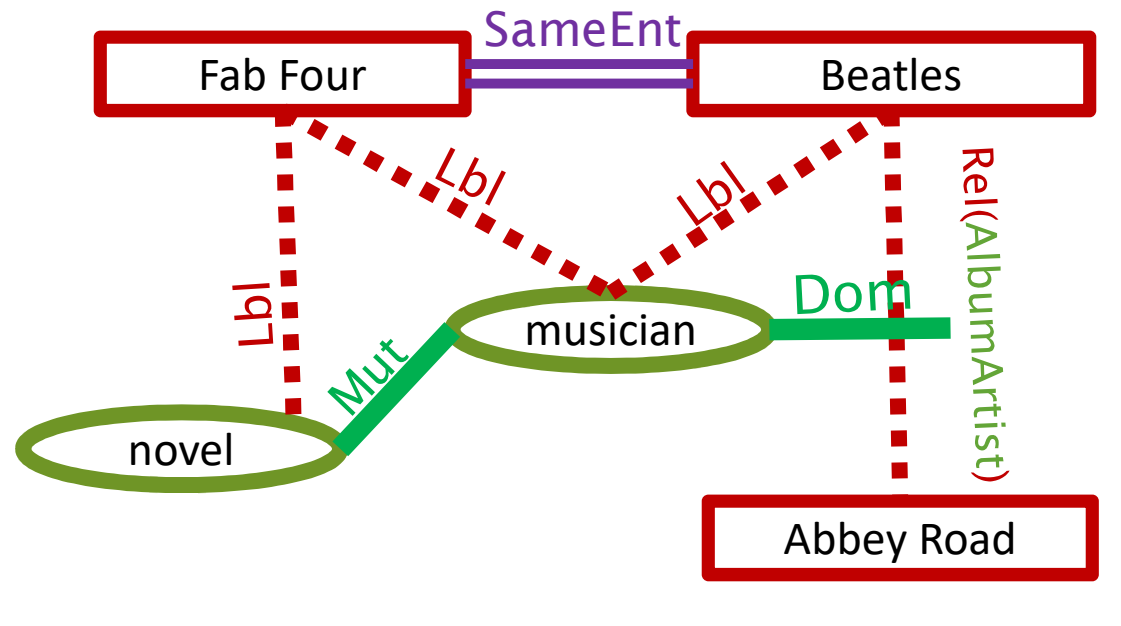
## Ontology:

- Dom(albumArtist, musician)
- Mut(novel, musician)

## Entity Resolution:

- SameEnt(Fab Four, Beatles)

## (Annotated) Extraction Graph



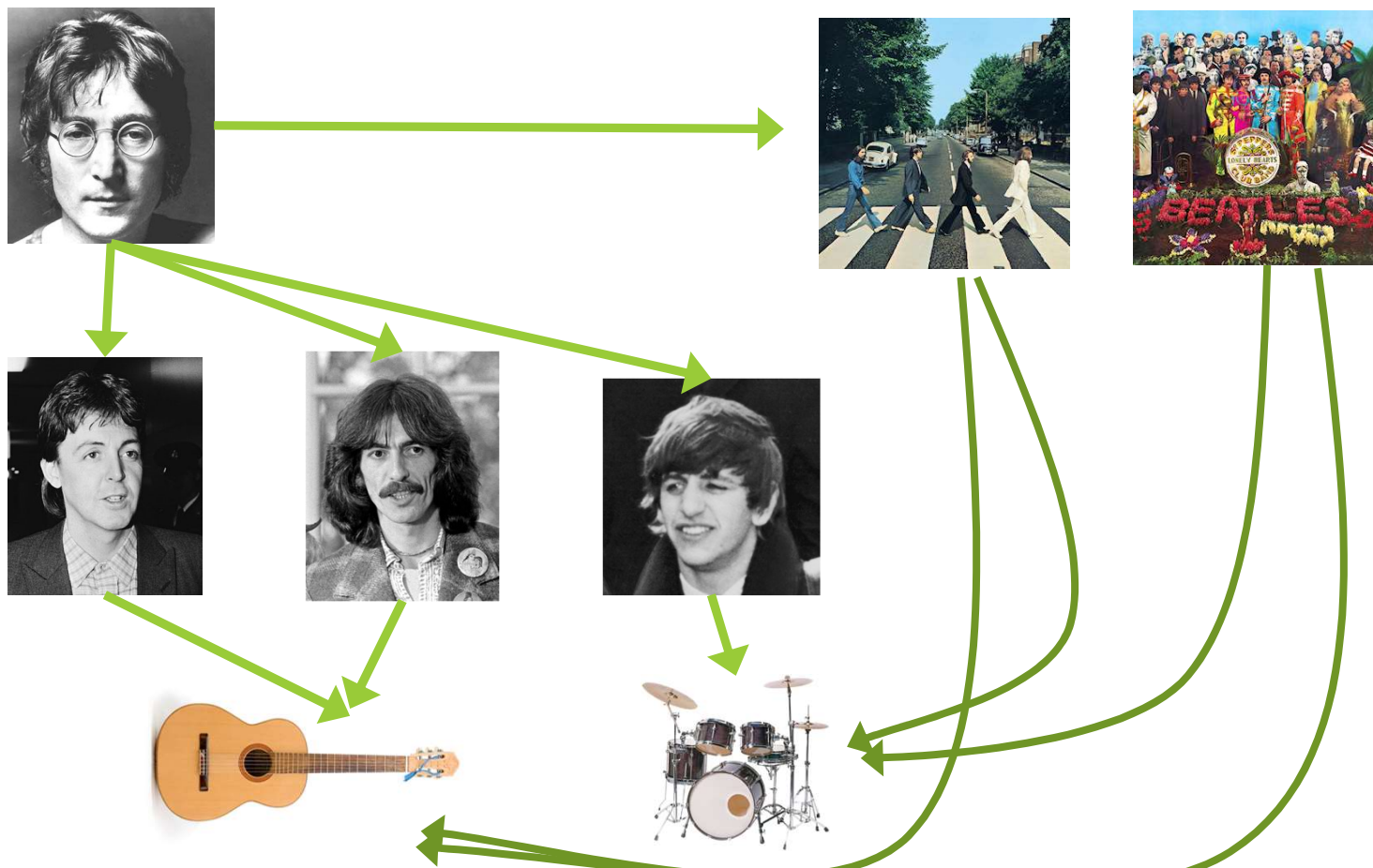
## After Knowledge Graph Identification





# Random Walk Illustration

Query: R(Lennon, PlaysInstrument, ?)



# Two approaches for KG construction

---

PROBABILISTIC MODELS

EMBEDDING BASED MODELS

# Why embeddings?

---

## Limitations of probabilistic models

### Limitation to Logical Relations

- Representation restricted by manual design
  - Clustering? Asymmetric implications?
  - Information flows through these relations
- Difficult to generalize to unseen entities/relations

### Computational Complexity of Algorithms

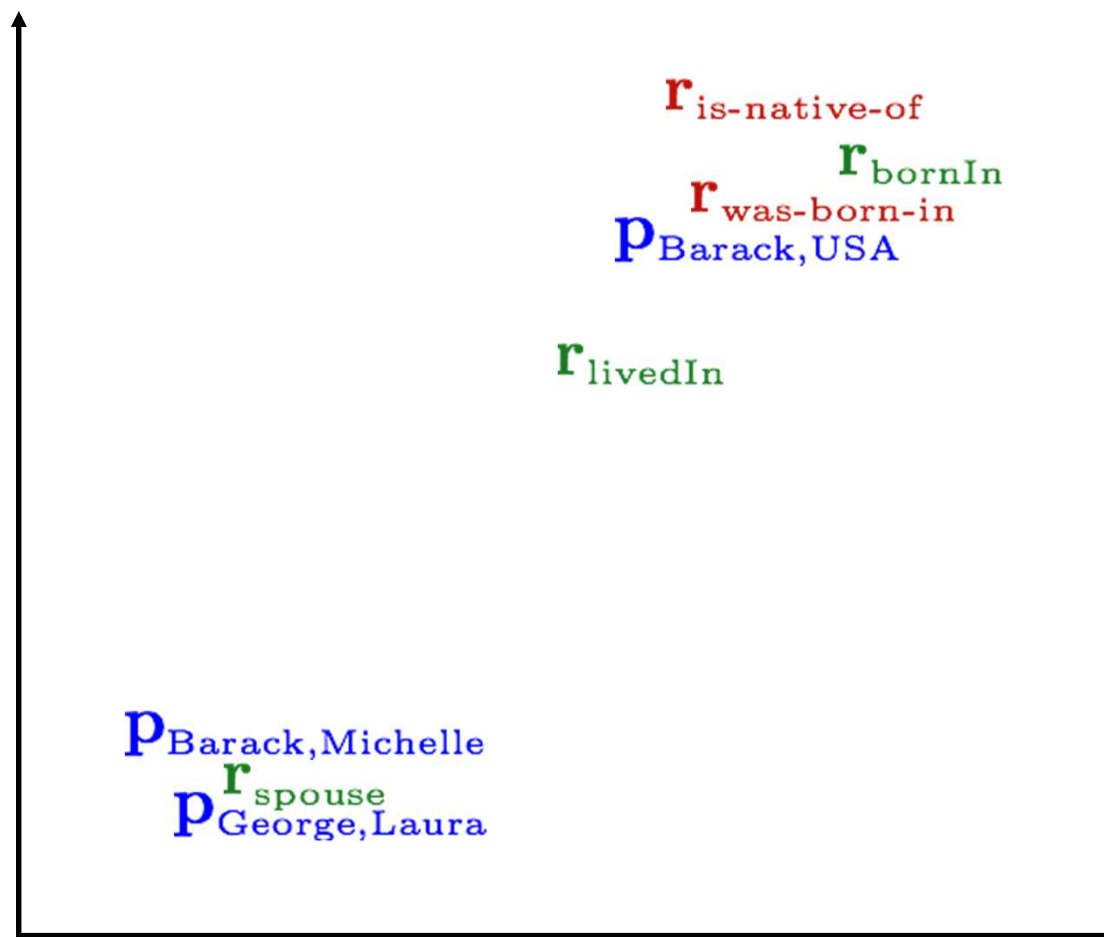
- Learning is NP-Hard, difficult to approximate
- Query-time inference is also NP-Hard
- Not easy to parallelize, or use GPUs
- Scalability is badly affected by representation

## Embedding based models

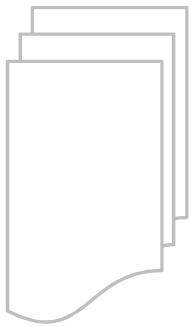
- Can generalize to unseen entities and relations
- Efficient inference at large scale

# Relation Embeddings

---



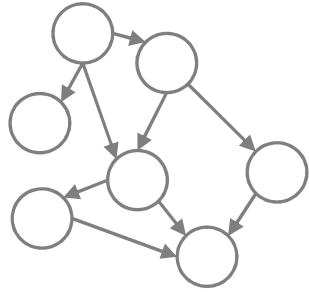
**Part 1: Knowledge Graphs**



**Part 2:  
Knowledge  
Extraction**



**Part 3:  
Graph  
Construction**



# Critical Overview

---

SUMMARY

SUCCESS STORIES

DATASETS, TASKS, SOFTWARES

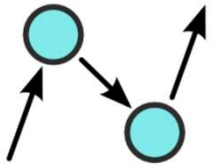
EXCITING RESEARCH DIRECTIONS

# Success stories



## Open Information Extraction

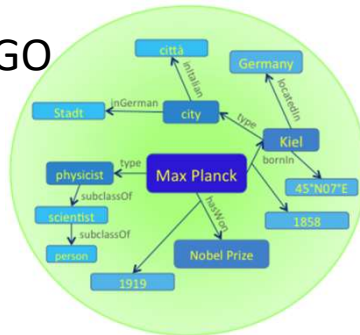
**NELL Knowledge Base Browser**  
CMU Read the Web Project



## ConceptNet

An open, multilingual knowledge graph

YAGO



**db DeepDive v0.8.0**  
Think about features, not algorithms.

# Success story: OpenIE (ReVerb)



Open Information Extraction

[openie.allenai.org](http://openie.allenai.org)



Argument 1:  Relation:   
Argument 2:  All

all location (21) film location (18) statistical region (16) name source (15) travel destination (14) misc.

more types ▾

were bigger than **Jesus** (100)

came to America (95)

appeared on **The Ed Sullivan Show** (88)

broke up in 1970 (56)

Here Comes the Sun (46)

came to America (45)

is for the future (44)

are a great band (42)

perform on **The Ed Sullivan Show** (39)

were **Musical ensemble** (36)

**are a great band** ▶

**Extracted Synonyms:**

were  
is  
was

**Extracted from these sentences:**

are **The Beatles** are **the best band** , hands down but Oasis did make a great cover . (via ClueWeb12)

**The Beatles** are **a great band** . (via ClueWeb12)

**The Beatles** are **the best band** . (via ClueWeb12)

**The Beatles** are **the greatest band** ... Started 1 month ago by georgedcc Yeah , Songs in the Key of Life is a bit much for 1 listen . (via ClueWeb12)

**The Beatles** , arguably , are **the greatest band** , and may or may not have the greatest name . (via ClueWeb12)

The point is , from my view , **The Beatles** are **a good band** , but way behind the greatest artists to ever grace rock . (via ClueWeb12)



# Success story: NELL

## NELL Knowledge Base Browser

CMU Read the Web Project

Search

log in | preferences | help/instructions | feedback

categories

relations

- everypromotedthing
- abstractthing
  - event
    - convention
    - musicfestival
    - protestevent
    - meetingeventtitle
    - conference
      - mlconference
    - weatherphenomenon
    - sportevent
      - sportgame
      - race
      - olympics
      - grandprix
    - crimeorcharge
    - earthquakeevent
    - election
    - bombingevent
    - militaryeventtype
      - militaryconflict
    - productlaunchevent
    - filmfestival
    - roadaccidentevent
    - meetingeventtype
    - eventoutcome
  - mlalgorithm
  - physiologicalcondition
    - disease

### beatles (musicartist)

literal strings: [BEATLES](#), [Beatles](#), [beatles](#)

#### Help NELL Learn!

NELL wants to know if these beliefs are correct.  
If they are or ever were, click thumbs-up. Otherwise, click thumbs-down.

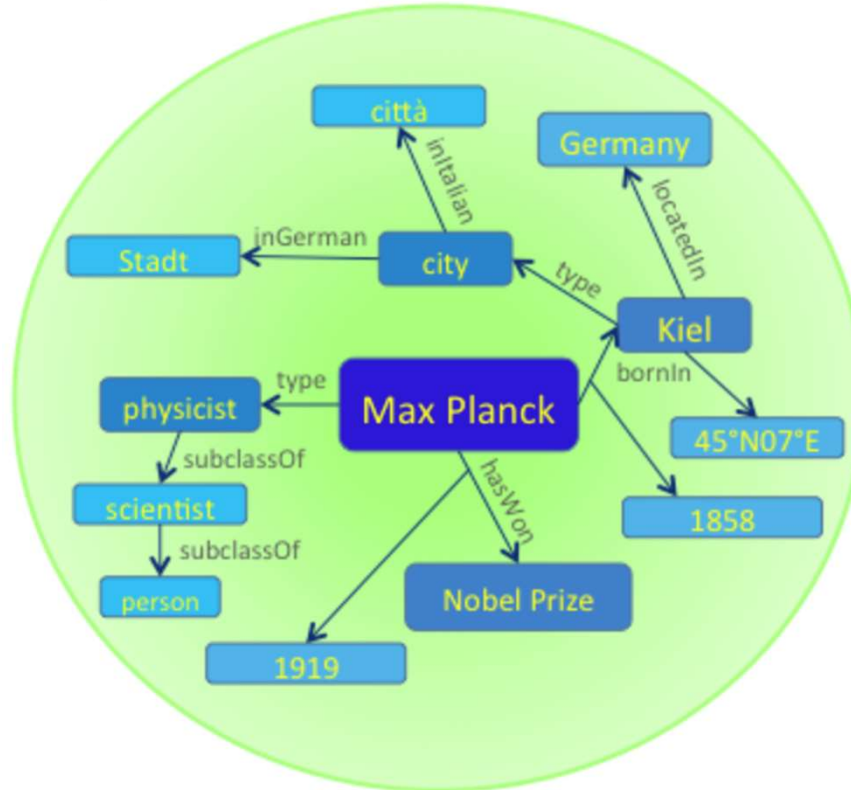
- [beatles](#) is a [musical artist](#)  
- [beatles](#) is a musician in the [genre classic\\_pop](#) (musicgenre)  
- [beatles](#) is a musician in the [genre pop](#) (musicgenre)  
- [beatles](#) is a musician in the [genre rock](#) (musicgenre)  
- [beatles](#) is a musician in the [genre classic\\_rock](#) (musicgenre)  

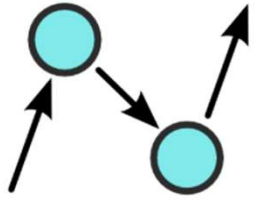
#### categories

- [musicartist](#)(100.0%)
  - MBL @198 (100.0%) on 07-feb-2011 [ Promotion of musicartist:beatles musicartistgenre musicgenre:classic\_rock ]
  - CPL @1021 (80.9%) on 14-oct-2016 [ "numerous other artists including \_ " "traducidas de \_ " " incluidas en \_ " had a guitar player" "early pioneers such as \_ " "controversial photo of \_ " "distressed image of \_ " "D-tracks of \_ " "Beatles Come Together \_ " "ohne die \_ " "opening band for \_ " "American acts like \_ " "classic acts like \_ " "performance footage of \_ " " were the perfect band" \_ " ' record label" "record album by \_ " "les paroles de \_ " " never recorded the song" "such renowned artists as \_ " " did a few songs" "Top artists include \_ " "crazy lives of \_ " "UK artists such as \_ " "Lennon started \_ " " ' musical talent" " ' Birthplace" " ' harmonies" "Tour , starring \_ " " ' last days" " ' fourth album" " ' sixth studio album" " ' original recordings" "They were also pushing \_ " "She Said by \_ " "Other artists featured include \_ " "Post general comments related to \_ " "track also shows \_ " "such major artists as \_ " "time favorite band is \_ " "past masters such as \_ " "pop hooks of \_ " "popular musicians like \_ " "pop icons such as \_ " "music artists like \_ " "music bands like \_ " "pop stars such as \_ " "pop influenced by \_ " " more the Beatles" "had a guitar player" "traducidas de \_ " " incluidas en \_ " had a guitar player" "early pioneers such as \_ " "controversial photo of \_ " "distressed image of \_ " "D-tracks of \_ " "Beatles Come Together \_ " "ohne die \_ " "opening band for \_ " "American acts like \_ " "classic acts like \_ " "performance footage of \_ " " were the perfect band" \_ " ' record label" "record album by \_ " "les paroles de \_ " " never recorded the song" "such renowned artists as \_ " " did a few songs" "Top artists include \_ " "crazy lives of \_ " "UK artists such as \_ " "Lennon started \_ " " ' musical talent" " ' Birthplace" " ' harmonies" "Tour , starring \_ " " ' last days" " ' fourth album" " ' sixth studio album" " ' original recordings" "They were also pushing \_ " "She Said by \_ " "Other artists featured include \_ " "Post general comments related to \_ " "track also shows \_ " "such major artists as \_ " "time favorite band is \_ " "past masters such as \_ " "pop hooks of \_ " "popular musicians like \_ " "pop icons such as \_ " "music artists like \_ " "music bands like \_ " "pop stars such as \_ " "pop influenced by \_ " ]

# Success story: YAGO

- **Input:** Wikipedia infoboxes, WordNet and GeoNames
- **Output:** KG with 350K entity types, 10M entities, 120M facts
- Temporal and spatial information





# ConceptNet

An open, multilingual knowledge graph

[Link](#)

## **en** beatles

An English term in ConceptNet 5.5

### Derived terms

- [en](#) beetle →
- [en](#) beatledom →
- [en](#) beatlemania →
- [en](#) beatlesque →
- [en](#) fourth beetle →

### beatles is a type of...

- [en](#) a British band →
- [en](#) man <sup>(n)</sup> →
- [en](#) band <sup>(n)</sup> →
- [en](#) musician <sup>(n)</sup> →
- [en](#) album <sup>(n)</sup> →

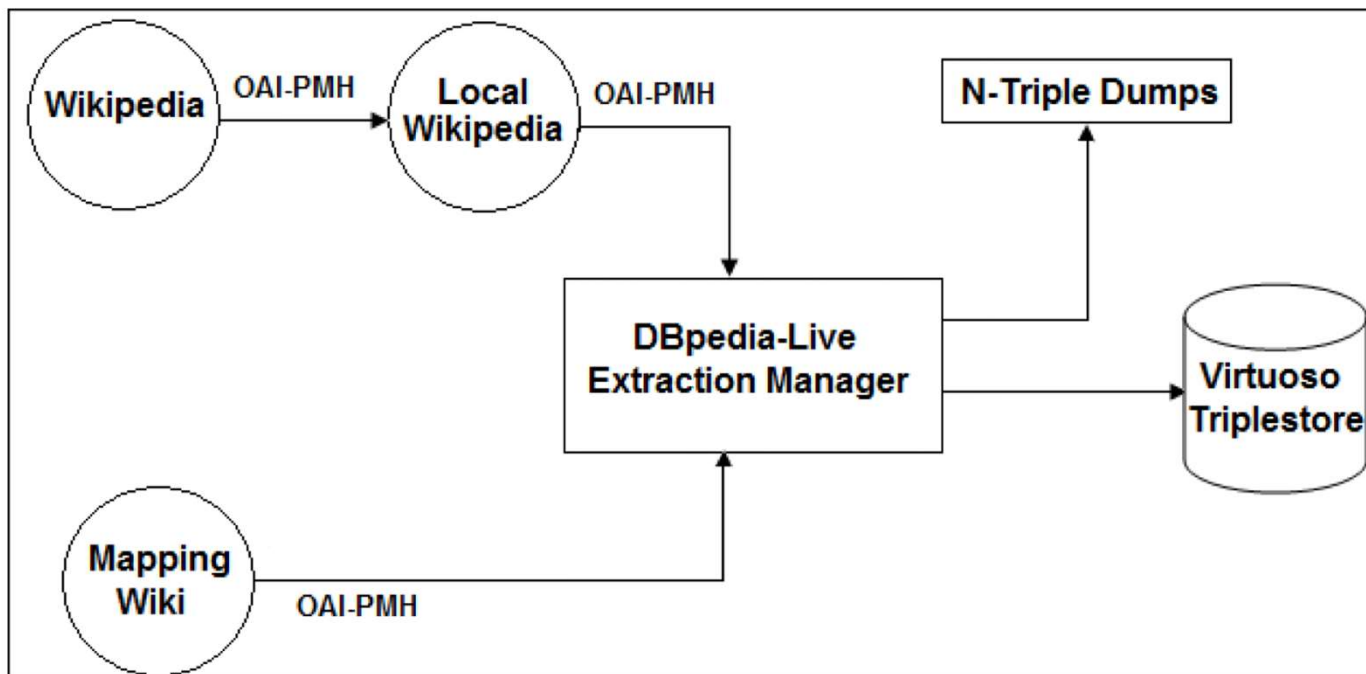
### Links to other sites

- [dbpedia.org](#) The Beatles →
- [sw.opencyc.org](#) Beatle →
- [umbel.org](#) Beatle →
- [wordnet-rdf.princeton.edu](#) 400520405-N →
- [wordnet-rdf.princeton.edu](#) 108386847-n →
- [wikidata.dbpedia.org](#) Q1299 →
- [en.wiktionary.org](#) Beatles →
- [dbpedia.org](#) The Beatles (No. 1) →
- [wikidata.dbpedia.org](#) Q738260 →
- [fr.wiktionary.org](#) Beatles →
- [dbpedia.org](#) The Beatles (The Original Studio Recordings) →
- [wikidata.dbpedia.org](#) Q603122 →

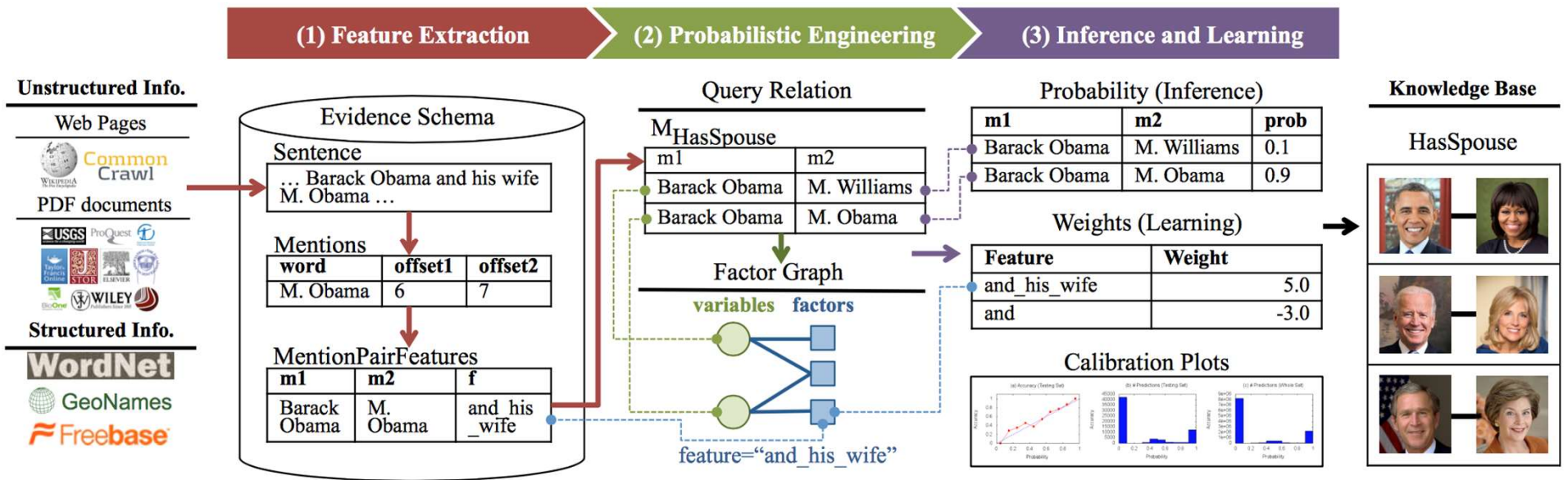
# Success story



- **DBpedia** is automatically extracted structured data from Wikipedia
  - 17M canonical entities
  - 88M type statements
  - 72M infobox statements















# DeepDive



- Best Precision/recall/F1 in KBP-slot filling task 2014 evaluations (31 teams participated)



# IE systems in practice

	Defining domain	Learning extractors	Scoring candidate facts	Fusing extractors
ConceptNet				
NELL				Heuristic rules
Knowledge Vault				Classifier
OpenIE				

# Critical Overview

---

SUMMARY

SUCCESS STORIES

DATASETS, TASKS, SOFTWARES

EXCITING RESEARCH DIRECTIONS

# Datasets

---

- KG as datasets
  - [FB15K-237](#) Knowledge base completion dataset based on Freebase<sup>1</sup>
  - [DBPedia](#) Structured data extracted from Wikipedia
  - [NELL](#) Read the web datasets
  - [AristoKB](#) Tuple knowledge base for Science domain
- Text datasets
  - [Clueweb09](#): 1 billion webpages (sample of Web)
  - [FACC1](#): Freebase Annotations of the Clueweb09 Corpora
  - [Gigaword](#): automatically-generated syntactic and discourse structure
  - [NYTimes](#): The New York Times Annotated Corpus
- Datasets related to Semi-supervised learning for information extraction  
[Link](#): entity typing, concept discovery, aligning glosses to KB, multi-view learning

<sup>1</sup>see Dettmers et al, 2017 for details (<https://arxiv.org/pdf/1707.01476.pdf>)



# Shared tasks

---

- Text Analysis Conference on Knowledge Base Population (TAC KBP)
  - **Slot filling task**
  - **Cold Start KBP Track**
  - **Tri-Lingual Entity Discovery and Linking Track (EDL)**
  - **Event Track**
  - **Validation/Ensembling Track**

# Software: NLP

---

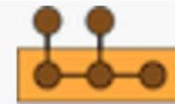
- Stanford CoreNLP: a suite of core NLP tools [\[link\]](#) (Java code)

 **Stanford CoreNLP**

- FIGER: fine-grained entity recognizer assigns over 100 semantic types [link](#) (Java code)

UNIVERSITY *of* WASHINGTON

- FACTORIE: out-of-the-box tools for NLP and information integration [link](#) (Scala code)



**FACTORIE**

- EasySRL: Semantic role labeling [link](#) (Java code)

UNIVERSITY *of* WASHINGTON

# Software: Extracting and Reasoning

---

- **Open IE**

(University of Washington)

Open IE 4.2 [link](#) (Scala code)

Stanford Open IE [link](#) (Java code)



- **Interactive Knowledge Extraction (IKE)**  
(Allen Institute for Artificial Intelligence)

[link](#) (Scala code)



- **PSL: Probabilistic soft logic**  
[link](#) (Java code)



- **ProPPR: Programming with Personalized PageRank**  
[link](#) (Java code)

**Carnegie Mellon University**

# Critical Overview

---

SUMMARY

SUCCESS STORIES

DATASETS, TASKS, SOFTWARES

EXCITING RESEARCH DIRECTIONS

# Exciting Active Research

---

- INTERESTING APPLICATIONS OF KG
- MULTI-MODAL INFORMATION EXTRACTION
- KNOWLEDGE AS SUPERVISION
- COMMON KNOWLEDGE

# Exciting Active Research

---

- INTERESTING APPLICATIONS OF KG
- MULTI-MODAL INFORMATION EXTRACTION
- KNOWLEDGE AS SUPERVISION
- COMMON KNOWLEDGE

# Interesting application of Knowledge Graphs

---

## The Literome Project [[link](#)]

- Automatic system for extracting genomic knowledge from PubMed articles
- Web-accessible knowledge base

**Search for directed genic interactions:**

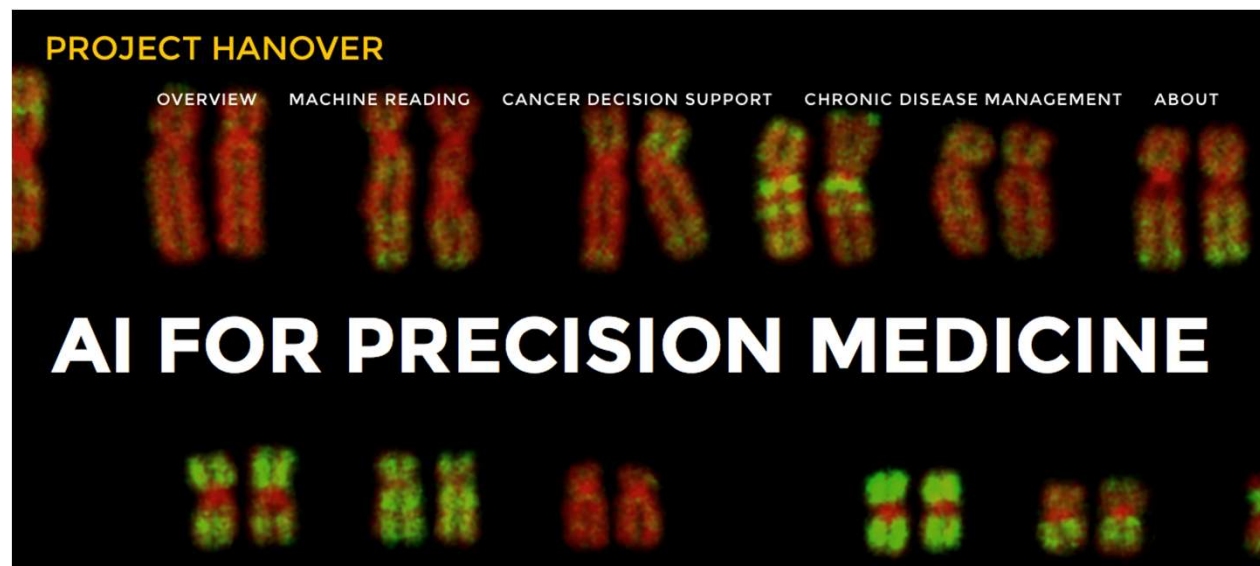
[help](#)

**Search for genotype-phenotype associations:**

[help](#)

# Interesting application of Knowledge Graphs

---



Microsoft  
Research

## **Chronic disease management:**

develop AI technology for predictive and preventive personalized medicine to reduce the national healthcare expenditure on chronic diseases (90% of total cost)

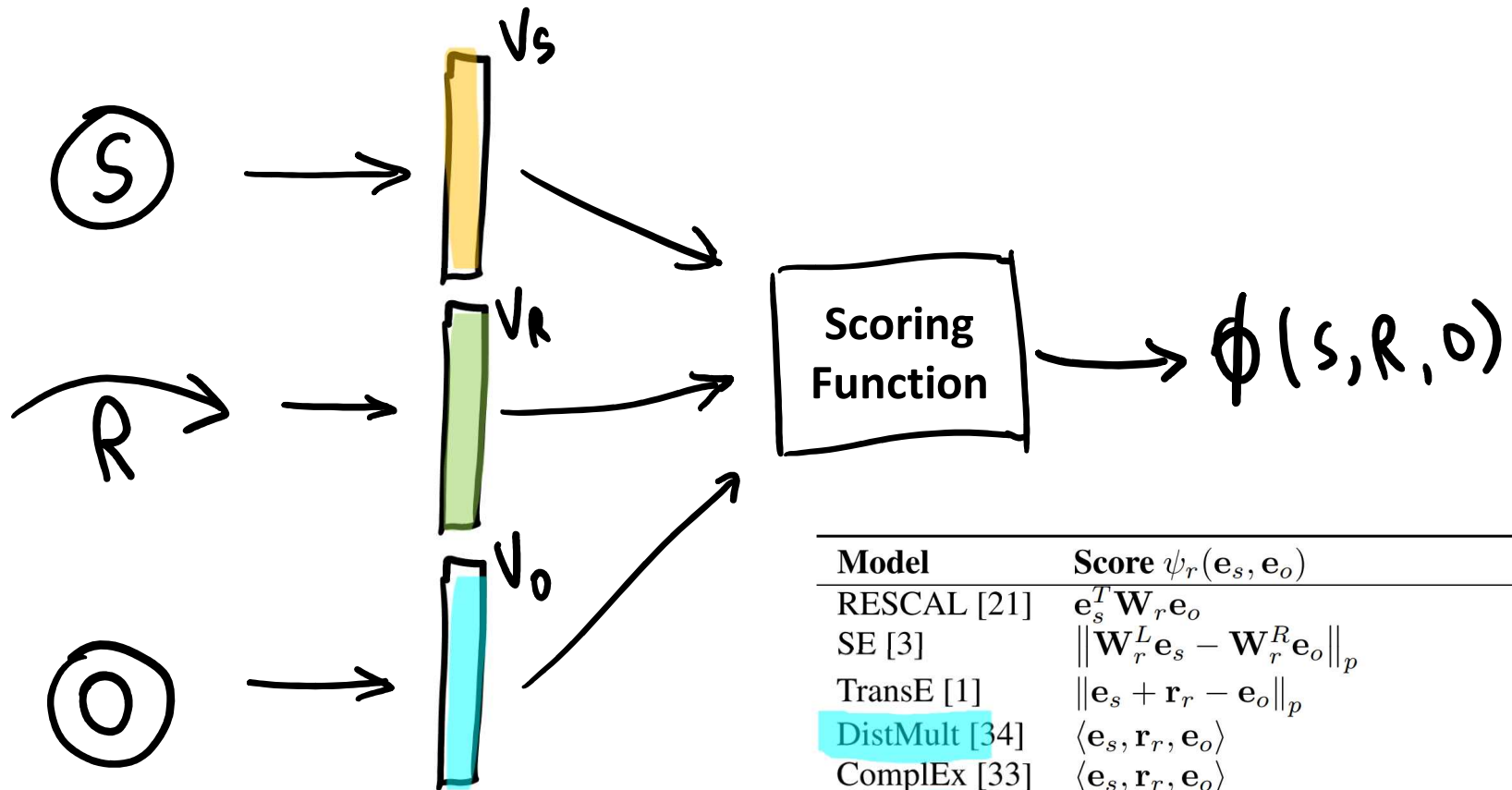


# Exciting Active Research

---

- INTERESTING APPLICATIONS OF KG
- MULTI-MODAL INFORMATION EXTRACTION
- KNOWLEDGE AS SUPERVISION
- COMMON KNOWLEDGE

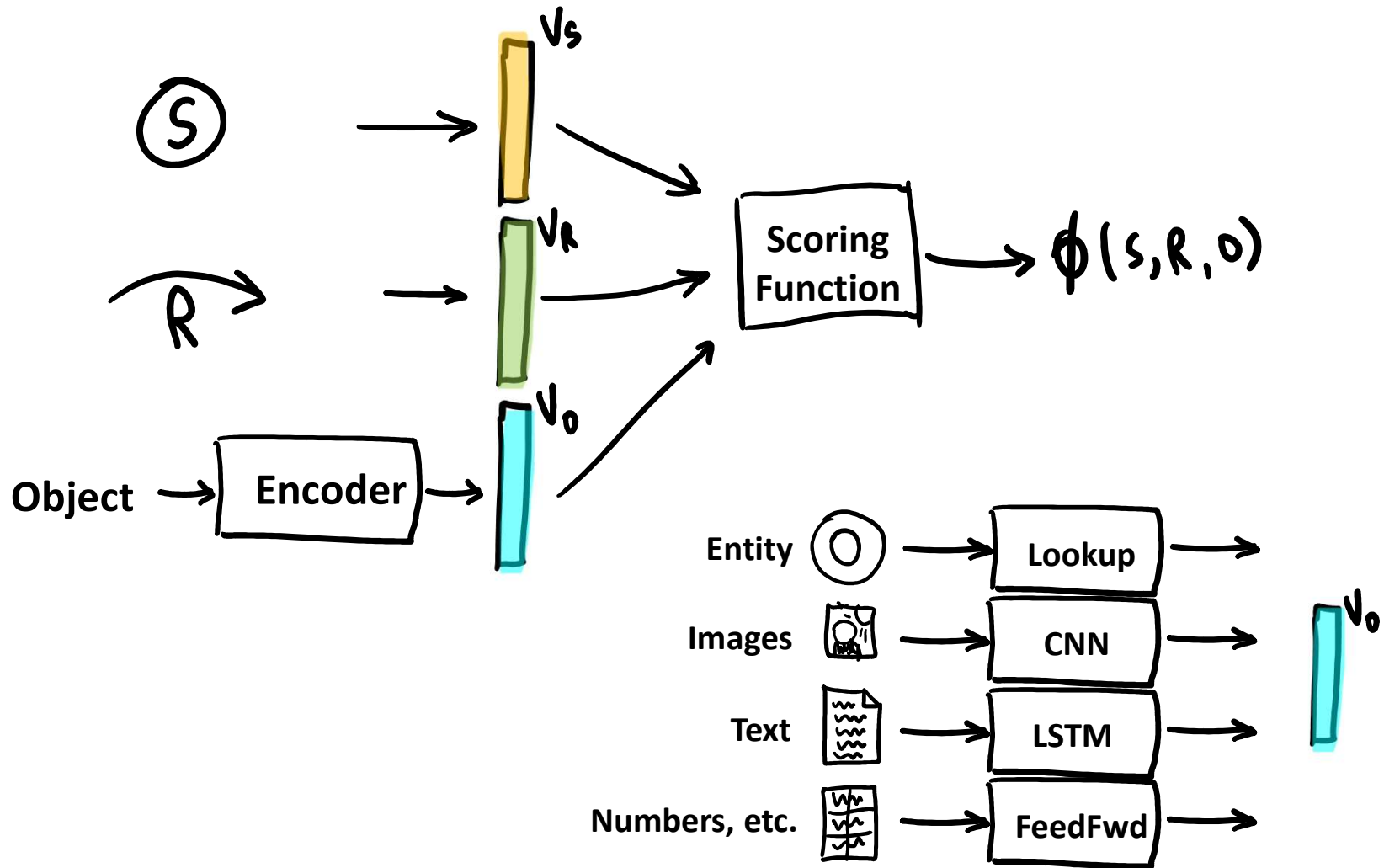
# Knowledge Base Completion



Model	Score $\psi_r(\mathbf{e}_s, \mathbf{e}_o)$
RESCAL [21]	$\mathbf{e}_s^T \mathbf{W}_r \mathbf{e}_o$
SE [3]	$\ \mathbf{W}_r^L \mathbf{e}_s - \mathbf{W}_r^R \mathbf{e}_o\ _p$
TransE [1]	$\ \mathbf{e}_s + \mathbf{r}_r - \mathbf{e}_o\ _p$
DistMult [34]	$\langle \mathbf{e}_s, \mathbf{r}_r, \mathbf{e}_o \rangle$
Complex [33]	$\langle \mathbf{e}_s, \mathbf{r}_r, \mathbf{e}_o \rangle$
ConvE	$f(\text{vec}(f([\overline{\mathbf{e}}_s; \overline{\mathbf{r}}_r] * \omega))) \mathbf{W} \mathbf{e}_o$

Table from Dettmers, et al. (2017)

# Multimodal KB Embeddings

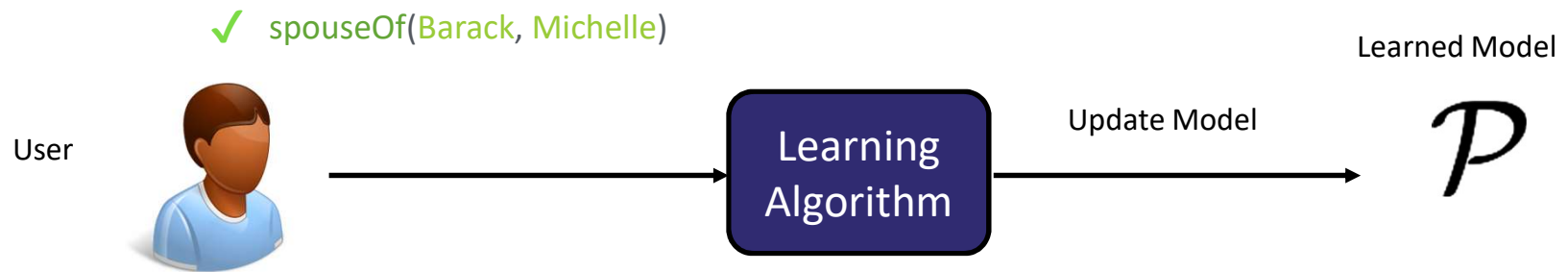


# Exciting Active Research

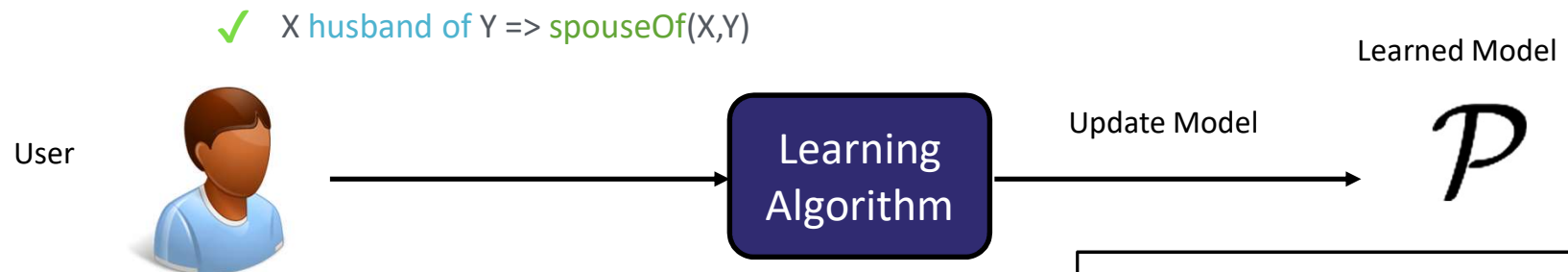
---

- INTERESTING APPLICATIONS OF KG
- MULTI-MODAL INFORMATION EXTRACTION
- KNOWLEDGE AS SUPERVISION
- COMMON KNOWLEDGE

# Knowledge as Supervision



**Problem 1:** Each annotation takes time  
**Problem 2:** Each annotation is a drop in the ocean



Many different options

- Generalized Expectation
- Posterior Regularization
- Labeling functions in SNORKEL

# Exciting Active Research

---

- INTERESTING APPLICATIONS OF KG
- MULTI-MODAL INFORMATION EXTRACTION
- KNOWLEDGE AS SUPERVISION
- COMMON KNOWLEDGE

# Aristo Science QA challenge

---

- Science questions dataset

~5K 4-way multiple choice questions

Frogs lay eggs that develop into tadpoles and then into adult frogs. This sequence of changes is an example of how living things \_\_\_\_\_

(A) go through a life cycle

(B) form a food web

(C) act as a source of food

(D) affect other parts of the ecosystem

Science knowledge

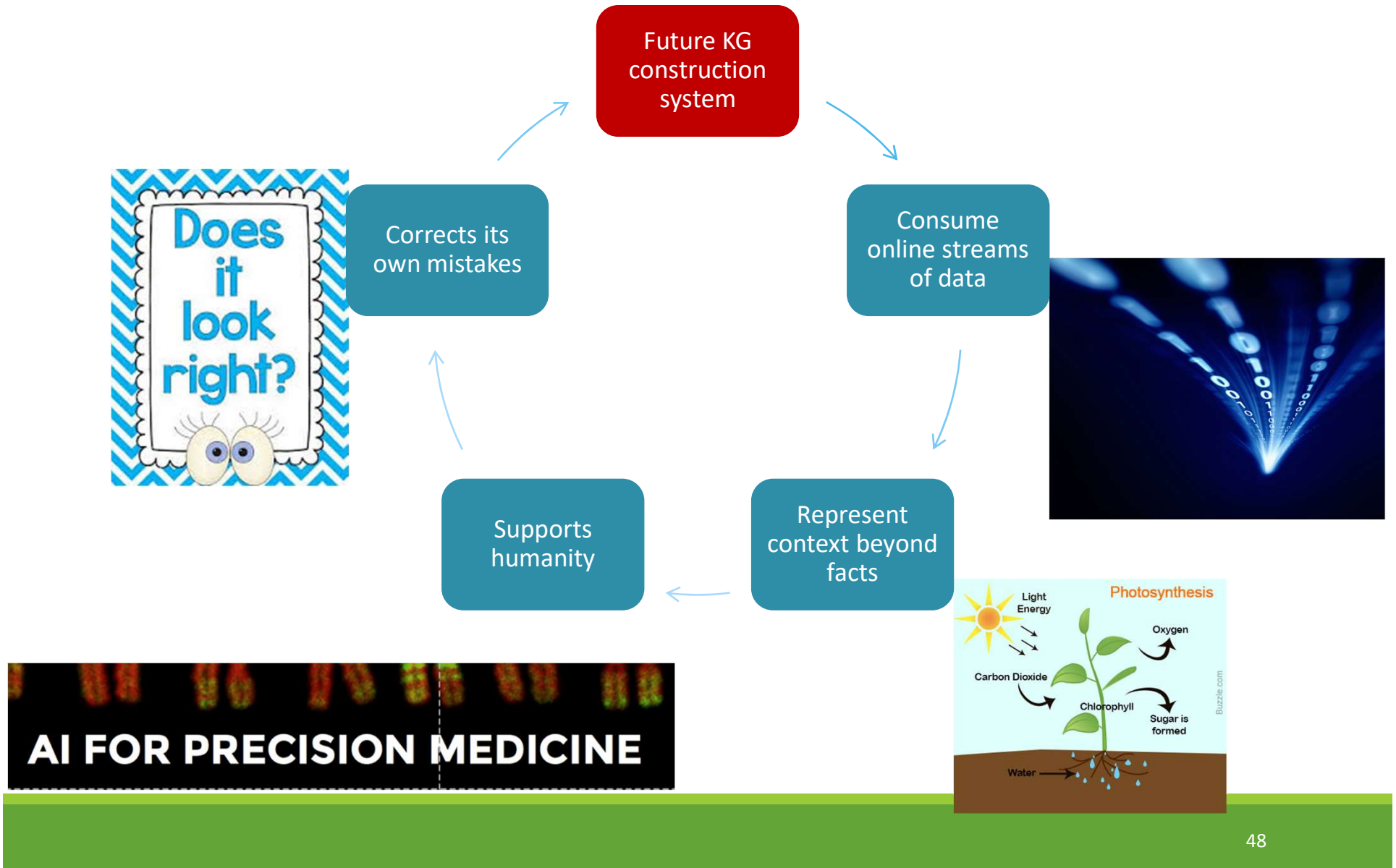
frog's life cycle,  
metamorphosis



Common sense  
knowledge

frog is an animal,  
animals have life cycle

# Future.....





# Thank You

---



**Jay Pujara**  
[jaypujara.org](http://jaypujara.org)  
[jay@cs.umd.edu](mailto:jay@cs.umd.edu)  
@jay\_mlr



**Sameer Singh**  
[sameersingh.org](http://sameersingh.org)  
[sameer@uci.edu](mailto:sameer@uci.edu)  
@sameer\_

# Two perspectives

---

	Extraction graph	Knowledge graph
Who are the entities? (nodes)		
What are their attributes? (labels)		
How are they related? (edges)		

# Natural Language Processing

Document

Within-doc Coreference...

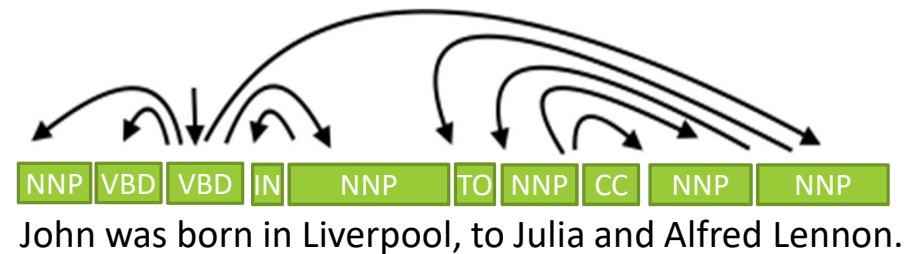
Lennon..                      Mrs. Lennon..                      his father  
John Lennon...                      the Pool                      .. his mother ..                      he    Alfred

Person                      Location                      Person                      Person

John was born in Liverpool, to Julia and Alfred Lennon.

Sentence

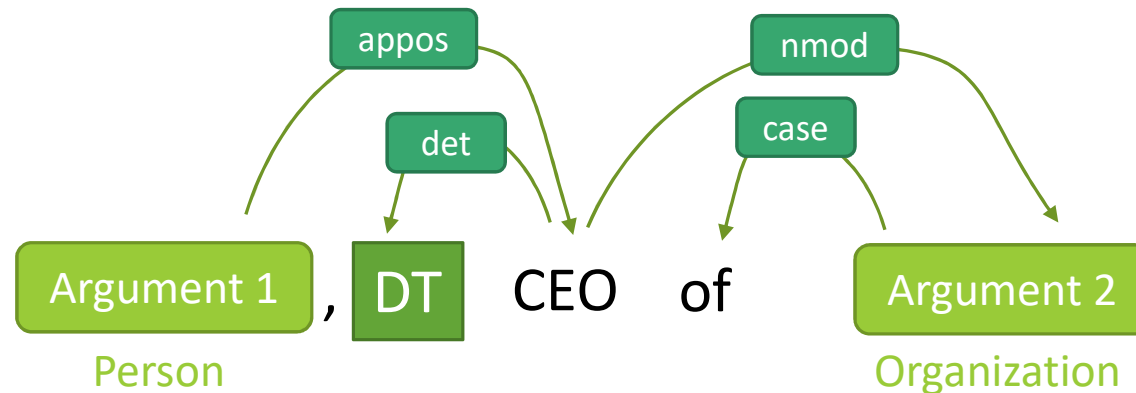
Dependency Parsing,  
Part of speech tagging,  
Named entity recognition...



# NLP annotations → features for IE

---

Combine tokens, dependency paths, and entity types to define rules.



Bill Gates, the CEO of Microsoft, said ...

Mr. Jobs, the brilliant and charming CEO of Apple Inc., said ...

... announced by Steve Jobs, the CEO of Apple.

... announced by Bill Gates, the director and CEO of Microsoft.

... mused Bill, a former CEO of Microsoft.

*and many other possible instantiations...*

# Success story: OpenIE

---

- **Key contributions:**

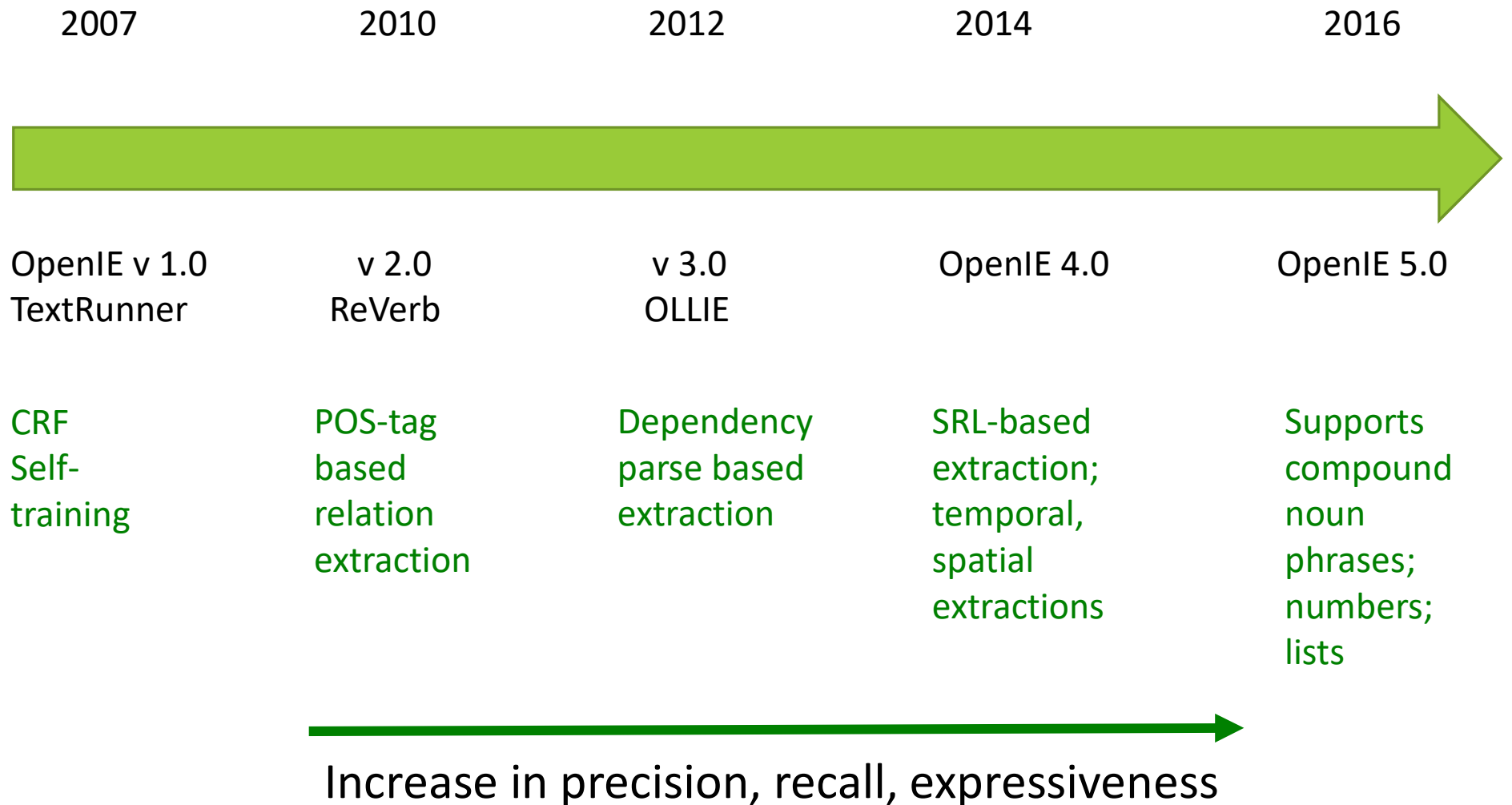
- No need for human defined relation schemas
- First ever successful open-source open domain IE system

- **ReVerb**

- Input = [Clueweb09 corpus](#) (1B web pages)
- Output = 15M high-precision extractions

# Open IE Systems

---



# Success story: NELL

---

- **Key technical contributions:**
  - “Never ending learning” paradigm
  - “Coupled bootstrap learning” to reduce semantic drift
  
- Input: Clueweb09 corpus (1B web pages)
- Ontology: ~2K predicates  
O(100K) constraints between predicates
- Output: 50 million candidate facts  
3 million high-confidence facts

# Success story: YAGO

---

- **Key contributions:**

- **Rich Ontology:** Linking Wikipedia categories to WordNet
- **High Quality:** High precision extractions (~95%)



# Success story: ConceptNet

---

- Commonsense knowledge base
- **Key contributions:**
  - **Freely available resource:** covers wide range of common sense concepts and relations organized in a easy-to-use semantic network
  - **NLP toolkit based on this resource:** supports analogy, text summarization, context dependent inferences
- ConceptNet4 was manually built using inputs from thousands of people
  - 28 million facts expressed in natural language
  - spanning 304 different languages

# DeepDive

---



- Machine learning based extraction system
- Key contributions:
  - **scalable, high-performance inference and learning engine**
  - **Developers contribute features (rules) not algorithms**
  - **Combines data from variety of sources (webpages, pdf, figures, tables)**

# Future.....

---



# Aristo ScienceKB

---

- AI2's TupleKB dataset: [link](#)
- **Open problems**
  - Best KR for Science domain
  - Domain targeted KB completion
  - Measuring recall w.r.t. end task

# (1) Future research directions:

## Going beyond facts

---

- Most of the existing KGs are designed to represent and extract binary relations → good enough for search engines
- Applications like QA demand in depth knowledge about higher level structures like activities, events, processes

## (2) Future research directions: Online KG Construction

---

- One shot KG construction → Online KG construction
  - Consume online stream of data
  - Temporal scoping of facts
  - Discovering new concepts automatically
  - Self-correcting systems

## (2) Future research directions: Online KG Construction

---

- **Continuously learning and self-correcting systems**
  - *[Selecting Actions for Resource-bounded Information Extraction using Reinforcement Learning, Kanani and McCallum, WSDM 2012]*
    - Presented a reinforcement learning framework for budget constrained information extraction
  - *[Never-Ending Learning, Mitchell et al. AAI 2015]*
    - Tom Mitchell says “Self reflection and an explicit agenda of learning subgoals” is an important direction of future research for continuously learning systems.

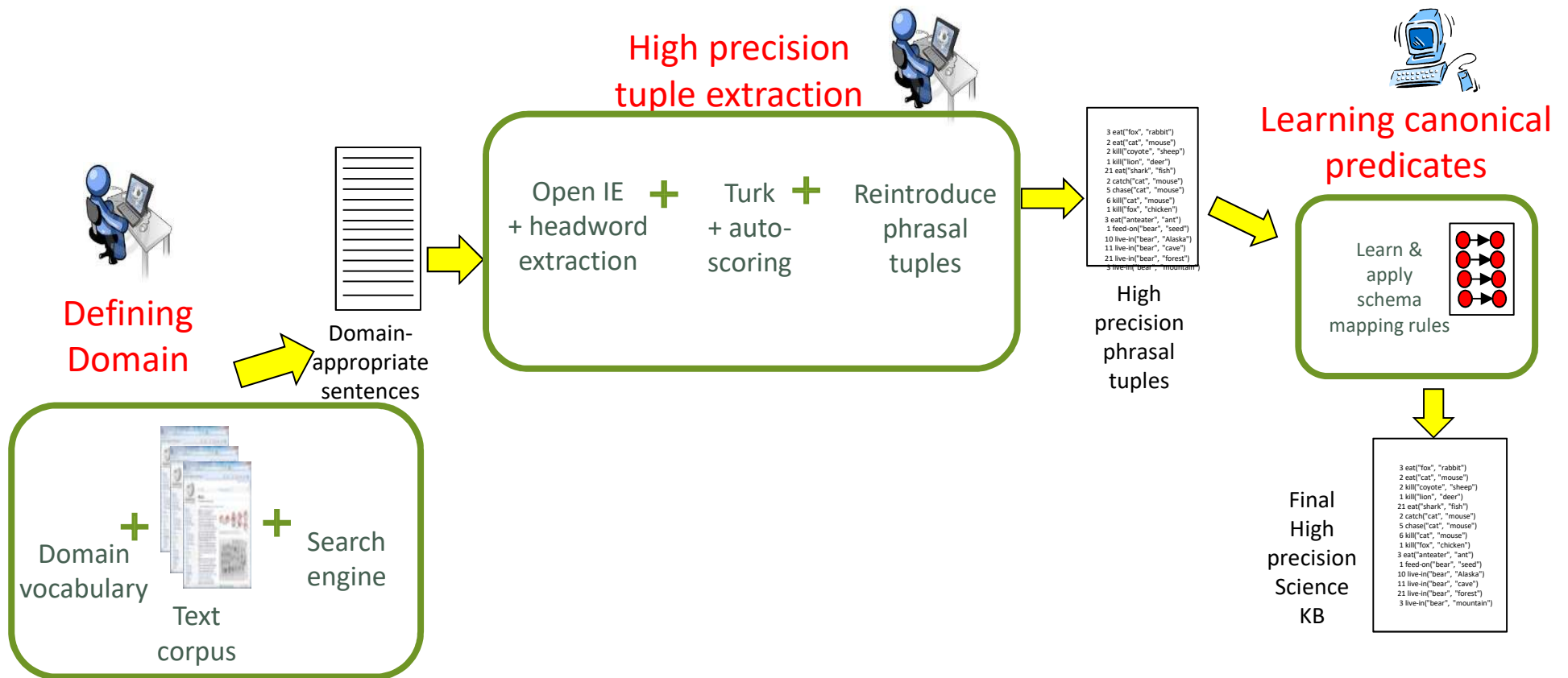


## Existing knowledge graphs

- Too named entity centric (no domain relevance)
- Too noisy (not directly usable by inference systems)

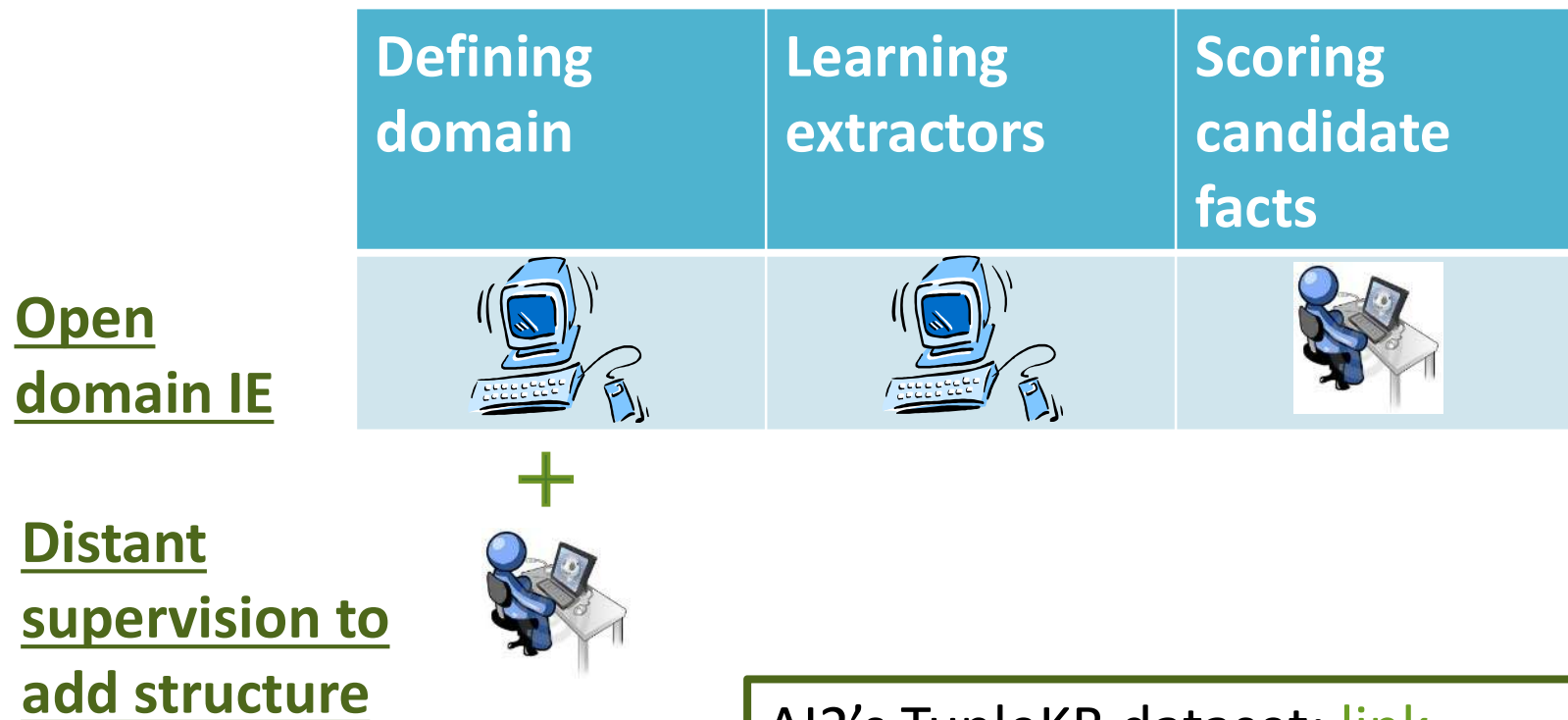


# AI2's ScienceKB



\*\*Upcoming article on "High Precision Knowledge Extraction for Science domain"

## Hybrid Approach: Adding structure to Open domain IE



AI2's TupleKB dataset: [link](#)

- > 300K common-sense and science facts
- > 80% precision

# Future research directions:

## Going beyond facts

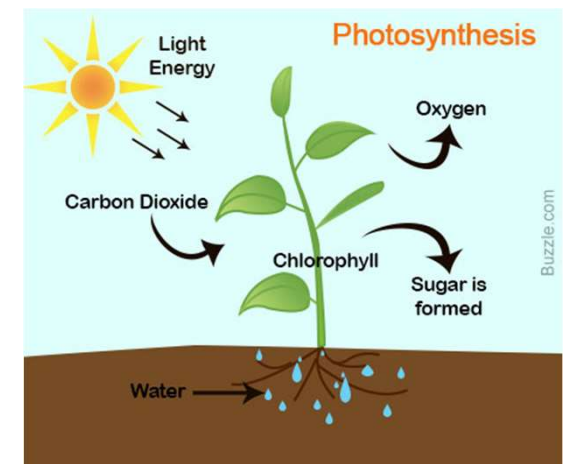
---

- **Fact:** Individual knowledge tuples  
(plant, take in, CO2)

subject	plant
predicate	Take in
object	CO2
time	daytime

- **Event frame:**  
more context how, when, where?

- **Processes:**  
representing larger structures, sequence of events  
e.g. Photosynthesis



# (3) Exciting active research: Ambitious Project


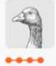





ALLEN INSTITUTE  
for ARTIFICIAL INTELLIGENCE

## The Allen AI Science Challenge

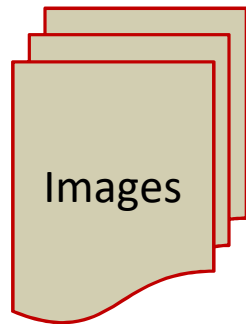
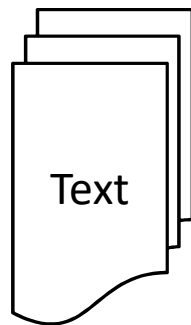
Is your model smarter than an 8th grader?

\$80,000 · a year ago

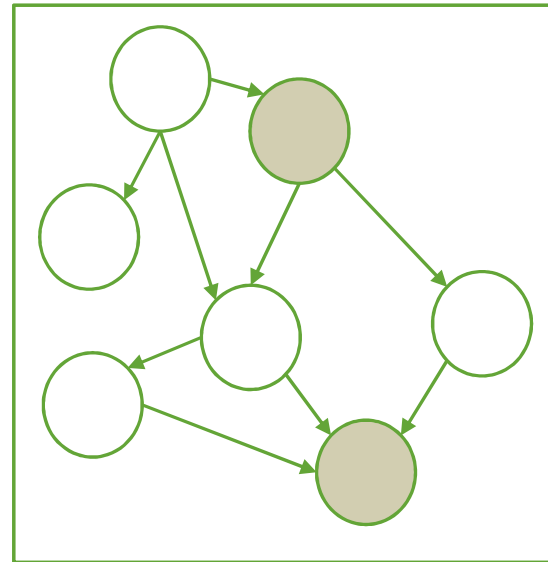
#	Δ1w	Team Name * in the money	Kernel	Team Members	Score ?	Entries	Last
1	—	* Alejandro Mosquera			0.59375	2	1y
2	new	* Cardal			0.59000	2	1y
3	new	* poweredByTalkwalker		   +4	0.59000	4	1y

## (2) Exciting active research: Multi-modal information extraction

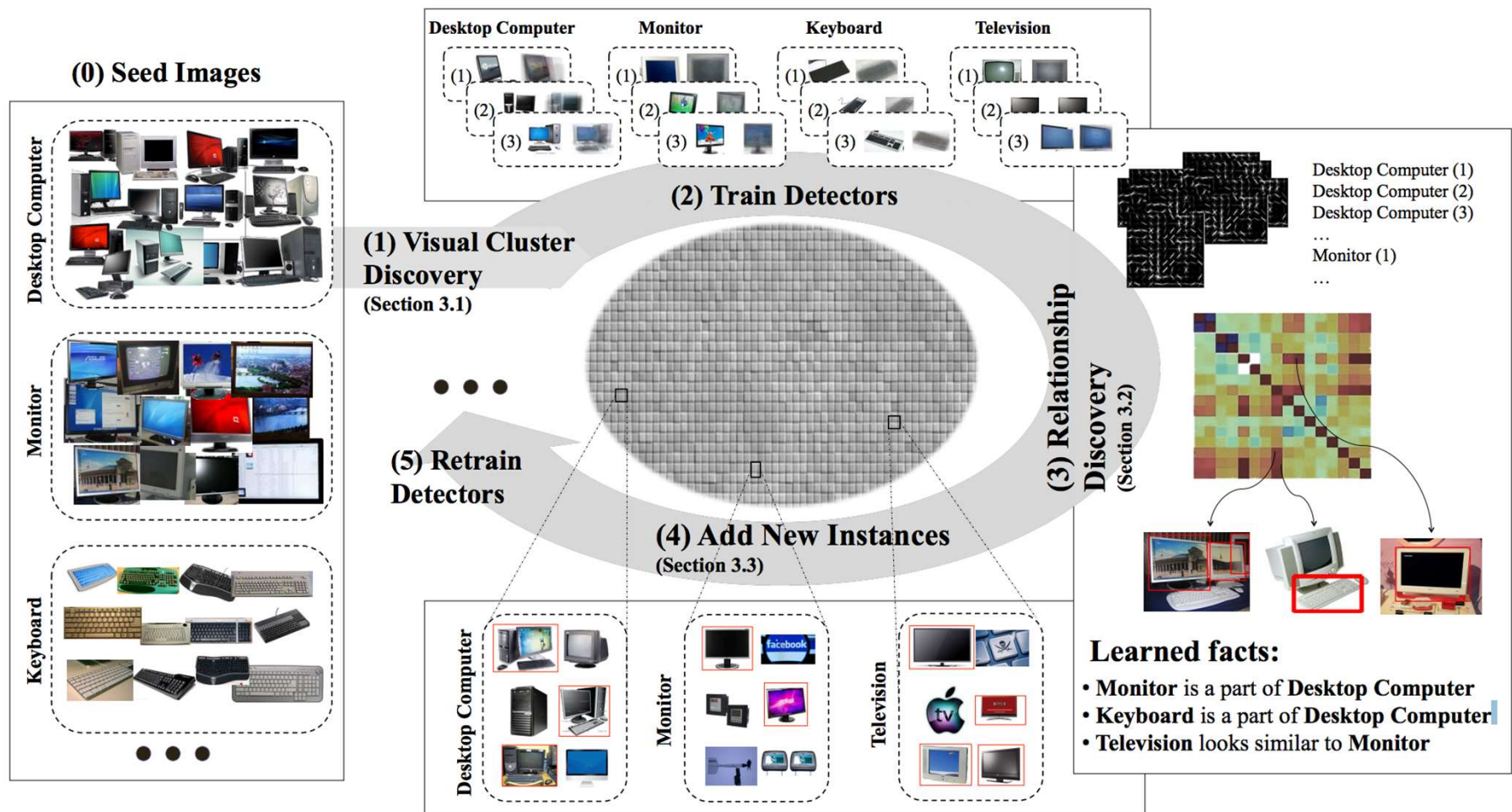
---



### Multi-modal Knowledge Graph



# NEIL: Extracting Visual Knowledge from Web Data



# NEIL: Extracting Visual Knowledge from Web Data

---



## Learned facts:

- **Monitor** is a part of **Desktop Computer**
- **Keyboard** is a part of **Desktop Computer**
- **Television** looks similar to **Monitor**



# WebChild: Text + Images

WEBCHILD Commonsense Browser

e.g. car,bicycle OR car OR a:fix bicycle



## Guess the concept

Domain ▲

Comparable ▲

Physical Part ▲

Activity ▲

Property ▲

Location ▲

Ask me!

mouse



*a hand-operated electronic device that controls the coordinates of a cursor on your computer screen as you move it around on a pad; on the bottom of the device is a ball that rolls on the surface of the pad; 'a mouse takes much more room than a trackball'*

keyboard



*device consisting of a set of keys on a piano or organ or typewriter or typesetting machine or computer or the like*



# Knowledge Base Completion

---



Entity Prediction

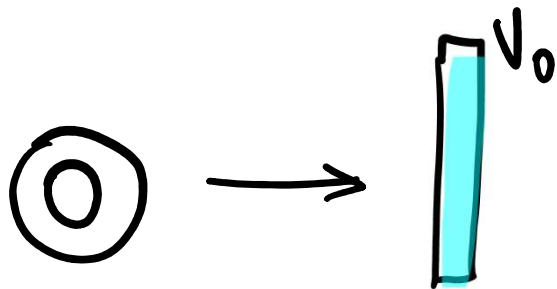


Link Prediction



# Restrictions in the Model

---



Each object has a vector representation:

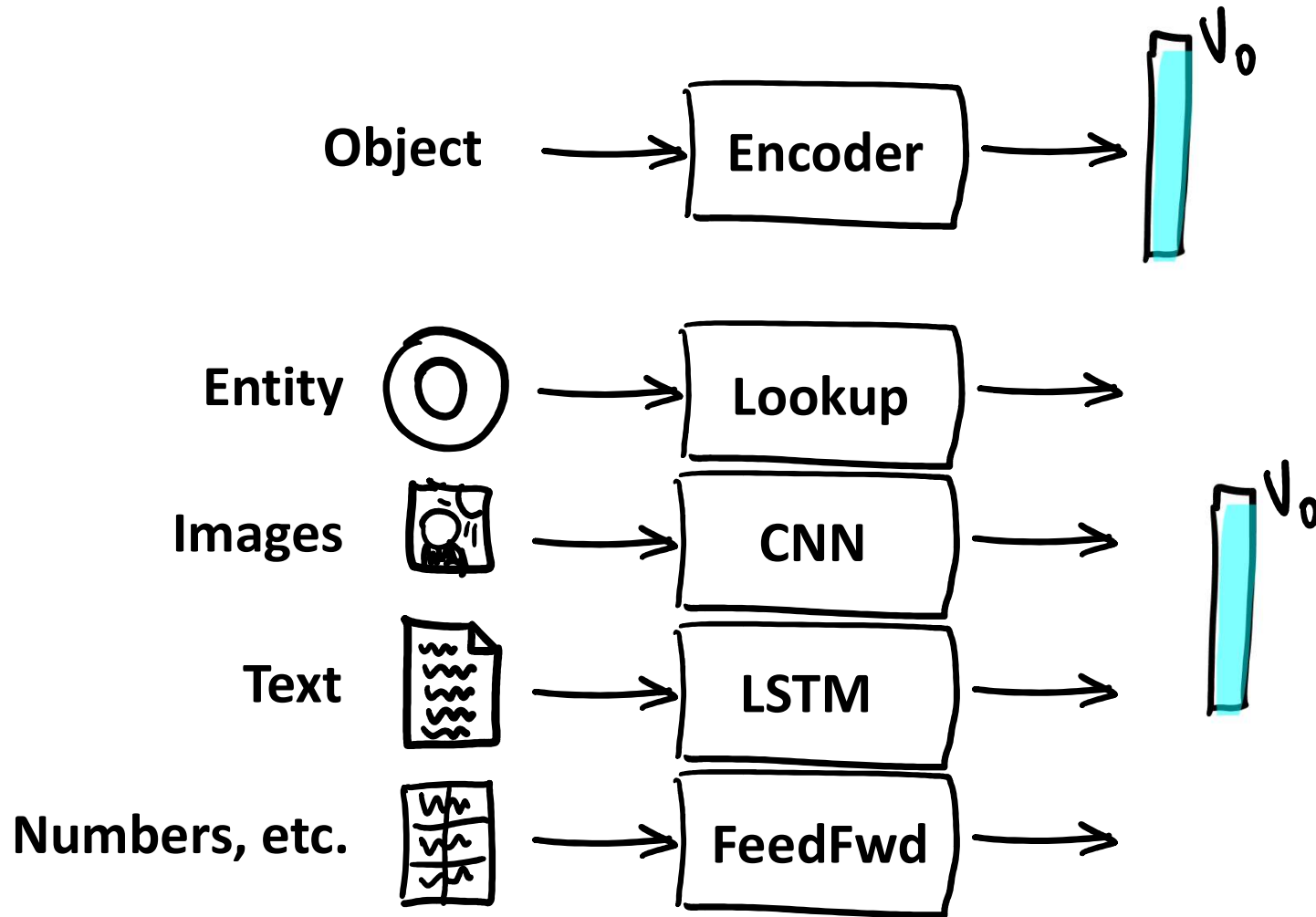
- Limits number of objects
- Large number of parameters
- Is not compositional (doesn't generalize)

What about other kinds of objects?

- Dates and Numbers: should generalize
- Text: Names and Descriptions
- Images: Portraits, Posters, etc.

# Multimodal KB Embeddings

---



# Augmenting Existing Datasets

---

MovieLens-100k-plus	
Relations	13
Users	943
Movies	1682
Posters	1651
Ratings	100,000

YAGO3-10-plus	
Relations	37 → 45
Entities	123,182
Structure Triples	1,079,040
Numbers (Years)	1651
Descriptions	107,326
Images	61,246