



# Knowledge Graph

- Machine Handlable Knowledge



# Outline

---

1. Knowledge Graph
2. Knowledge Extraction
3. Knowledge Graph Construction
4. Overview and Conclusion



# What is a knowledge graph?

---



# What is a knowledge graph?

---

- Knowledge in graph form!



# What is a knowledge graph?

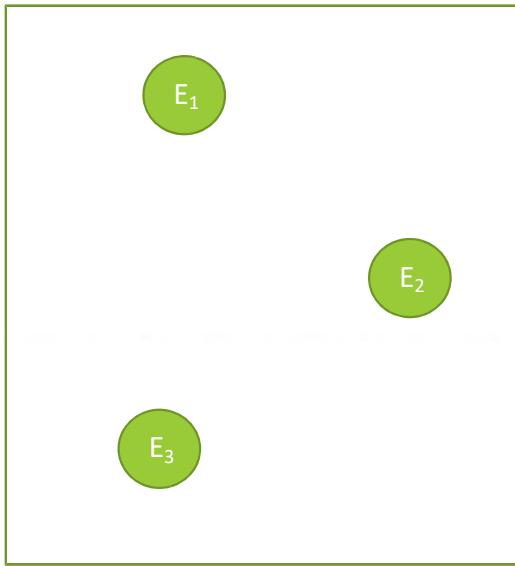
---

- Knowledge in graph form!
- Captures entities, attributes, and relationships

# What is a knowledge graph?

- Knowledge in graph form!

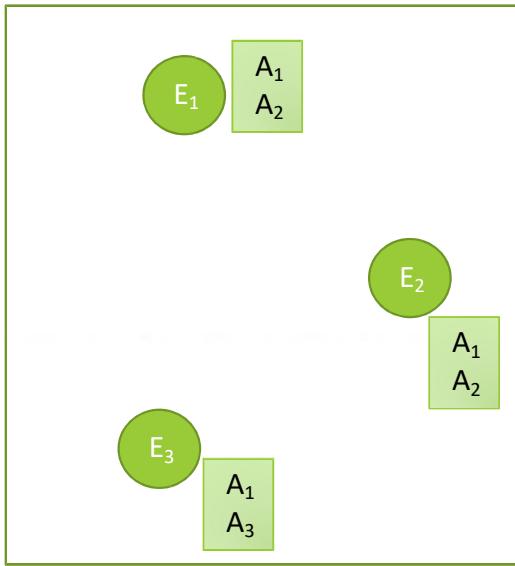
- Captures entities, attributes, and relationships
- Nodes are entities



# What is a knowledge graph?

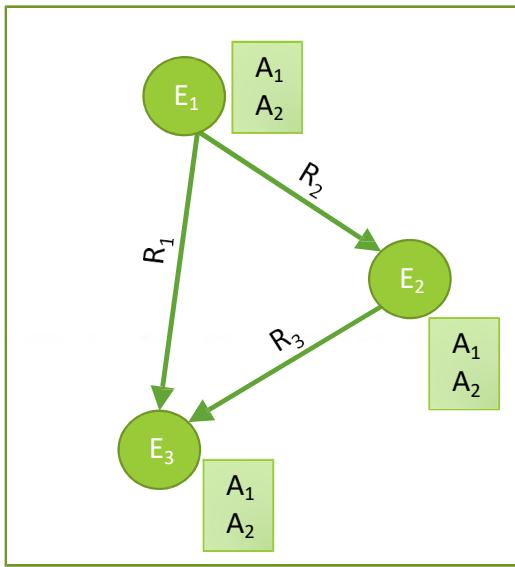
- Knowledge in graph form!

- Captures entities, attributes, and relationships
- Nodes are entities
- Nodes are labeled with attributes (e.g., types)



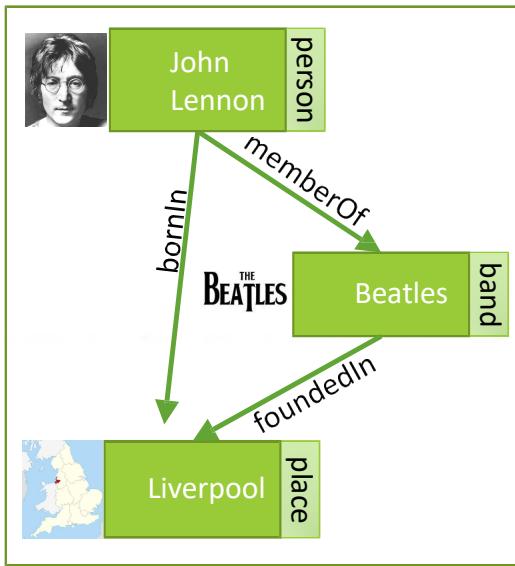
# What is a knowledge graph?

- Knowledge in graph form!
- Captures entities, attributes, and relationships
- Nodes are entities
- Nodes are labeled with attributes (e.g., types)
- Typed edges between two nodes capture a relationship between entities



# Example knowledge graph

- Knowledge in graph form!
- Captures entities, attributes, and relationships
- Nodes are entities
- Nodes are labeled with attributes (e.g., types)
- Typed edges between two nodes capture a relationship between entities





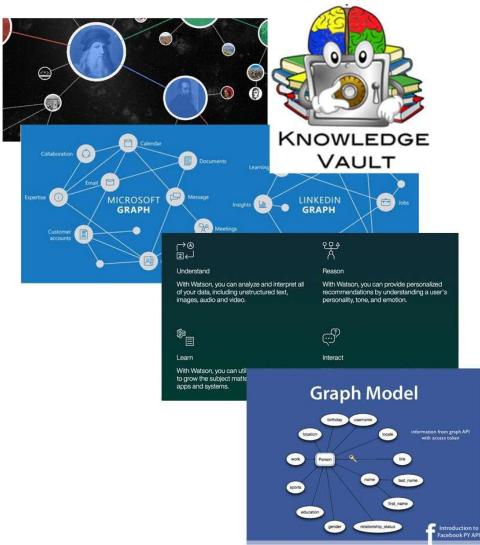
# Why knowledge graphs?

---

- Humans:
  - Combat information overload
  - Explore via intuitive structure
  - Tool for supporting knowledge-driven tasks
  
- Als:
  - Key ingredient for many AI tasks
  - Bridge from data to human semantics
  - Use decades of work on graph analysis

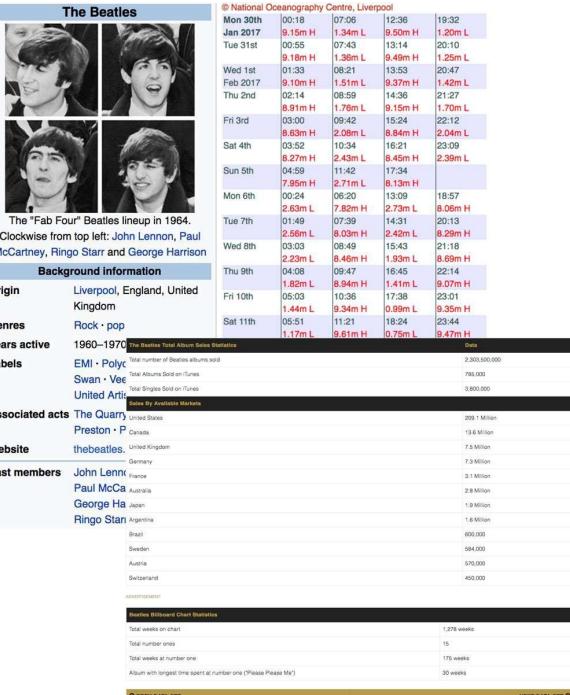
# Knowledge Graphs & Industry

- Google Knowledge Graph
  - Google Knowledge Vault
- Amazon Product Graph
- Facebook Graph API
- IBM Watson
- Microsoft Satori
  - Project Hanover/Literome
- LinkedIn Knowledge Graph
- Yandex Object Answer



# Where do knowledge graphs come from?

- Structured Text
  - Wikipedia Infoboxes, tables, databases, social nets





# Where do knowledge graphs come from?

- Structured Text
  - Wikipedia Infoboxes, tables, databases, social nets
- Unstructured Text
  - WWW, news, social media, reference articles

## Beatles last live performance

Published: Thursday, January 26th 2017, 5:24 am PST

Updated: Monday, January 30th 2017, 4:06 am PST

Written by Jim Eftink, Producer [CONNECT](#)



(KFVS) - How about a little Beatles history.

It was on this date in 1969, the band performed their last live public performance.

**Allan Williams, First Manager of the Beatles, Dies at 86**

(Source: Stock image) By ALLAN KOZINN DEC. 31, 2016

OU [View Profile](#) [Follow](#) [Report](#) [174 of 10,000 users](#)

The Harrison family is proud to announce the release of George Harrison – [The Vinyl Collection](#) box set featuring all of George Harrison's solo studio albums on vinyl.

**IWO GEORGE HARRISON - THE VINYL COLLECTION**  
Released on 24th February, 2017, the vinyl box set includes all twelve of George Harrison's solo studio albums, plus his original release track listing and artwork. Also included in the box set are George's classic album *Love in Japan* (2001).

See More



[George Harrison - The Vinyl Collection - Released February 24th 2017](#)

George Harrison - The Vinyl Collection, available to purchase now with an exclusive limited edition...

Comments

908 shares

[Write a comment...](#)

**Jeffrey Smith** What I would really be interested in is an "All Things Must Pass" box set with just the solo tracks without Phil Spector's Wall of Sound. It better won't be released in the US, so I would buy it in a heartbeat.

[Like](#) [Reply](#) [102d](#) · January 17 at 10:50am

[More](#)

[Dave Standring](#)

I have just seen the new box set from the great music industry rockers at EMI. They are really rubbing their hands with glee since more whitewashing various methods to make people buy their already bought and paid for record collection.

[Like](#) [Reply](#) [14d](#) · January 17 at 10:51am · Edited

[40 Replies](#)

[View more comments](#)

2 of 166

OU [View Profile](#) [Follow](#) [Report](#) [January 17 at 8:50am · v4](#)

"Of very few individual songs can it be said, "This changed the course of popular music." "A Day in the Life" did. It did. Right here.

**The Beatles - A Day In The Life**

A Day in The Life The Beatles' Video Collection is Out Now. Get your copy here!

[More](#)

[View more posts](#)

31

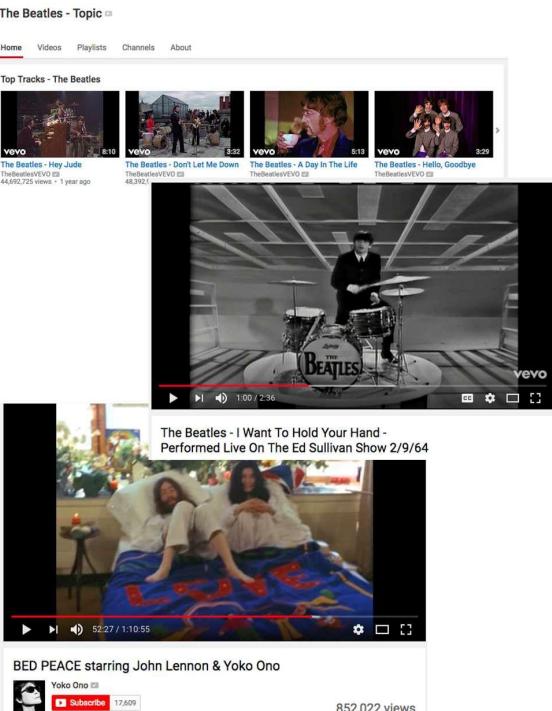
# Where do knowledge graphs come from?

- Structured Text
  - Wikipedia Infoboxes, tables, databases, social nets
- Unstructured Text
  - WWW, news, social media, reference articles
- Images



# Where do knowledge graphs come from?

- Structured Text
  - Wikipedia Infoboxes, tables, databases, social nets
- Unstructured Text
  - WWW, news, social media, reference articles
- Images
- Video
  - YouTube, video feeds



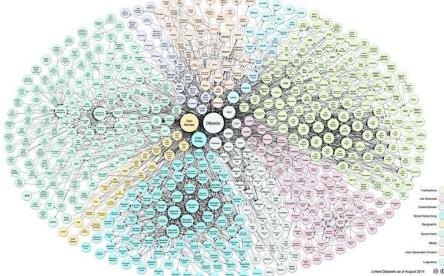


# Knowledge Representation

- Decades of research into knowledge representation
- Most knowledge graph implementations use RDF triples
  - <rdf:subject, rdf:predicate, rdf:object> : r(s,p,o)
  - Temporal scoping, reification, and skolemization...
- ABox (assertions) versus TBox (terminology)
- Common ontological primitives
  - rdfs:domain, rdfs:range, rdf:type, rdfs:subClassOf, rdfs:subPropertyOf, ...
  - owl:inverseOf, owl:TransitiveProperty, owl:FunctionalProperty, ...

# Semantic Web

- Standards for defining and exchanging knowledge
  - RDF, RDFa, JSON-LD, schema.org
  - RDFS, OWL, SKOS, FOAF
- Annotated data provide critical resource for automation
- Major weakness: annotate everything?





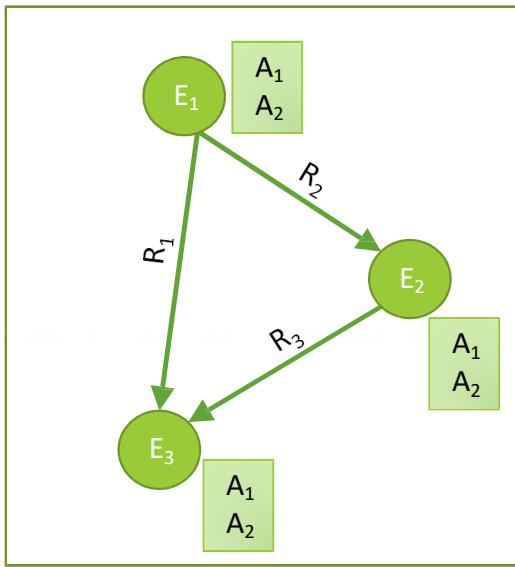
# Information Extraction from Text

---

- Answer to the knowledge acquisition bottleneck
- Many challenges:
  - chunking
  - polysemy/word sense disambiguation
  - entity coreference
  - relational extraction

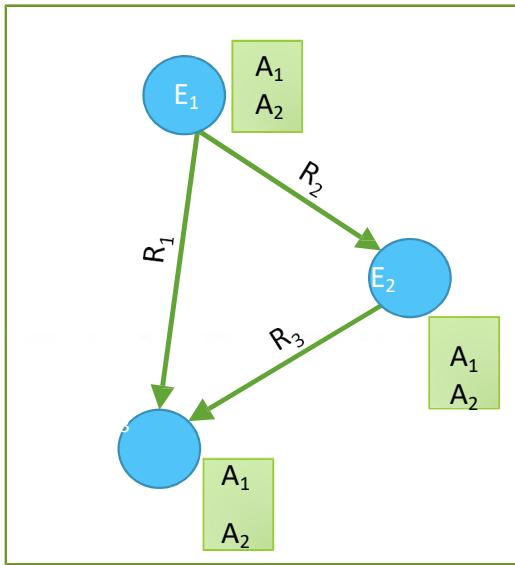
# What is a knowledge graph?

- Knowledge in graph form!
- Captures entities, attributes, and relationships
- Nodes are entities
- Nodes are labeled with attributes (e.g., types)
- Typed edges between two nodes capture a relationship between entities



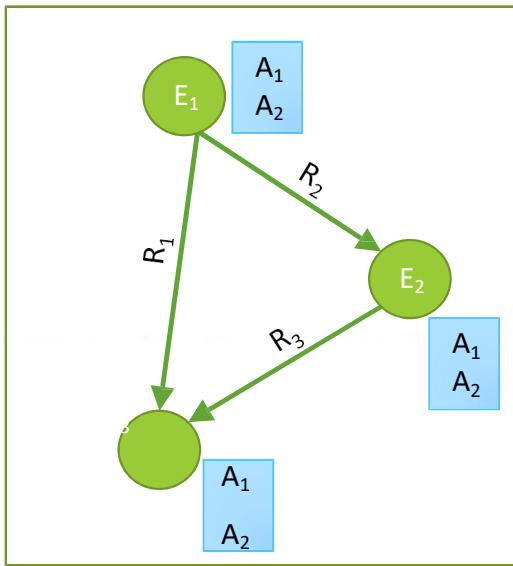
# Basic problems

- Who are the entities (nodes) in the graph?



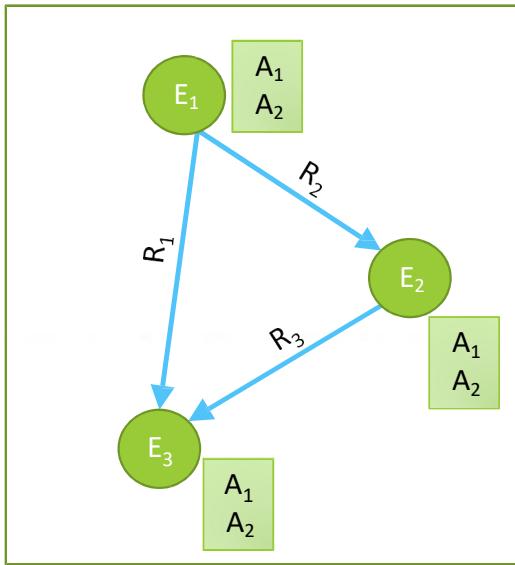
# Basic problems

- Who are the entities (nodes) in the graph?
- What are their attributes and types (labels)?



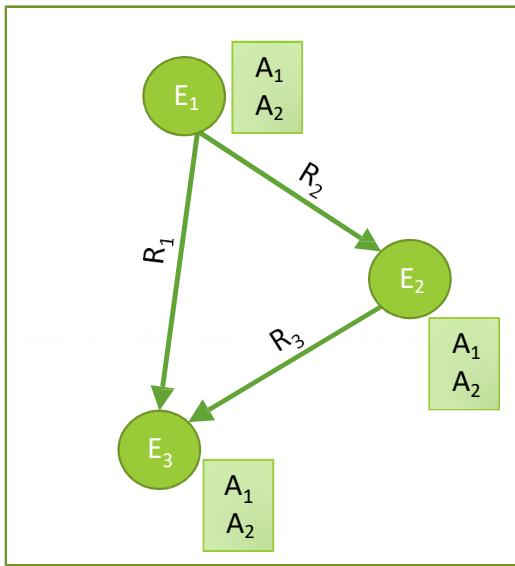
# Basic problems

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



# Basic problems

- **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

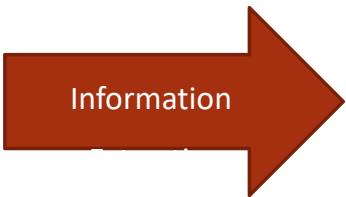
- Knowledge Extraction
- Who are the entities (nodes) in the graph?
  - Named Entity Recognition
  - Entity Coreference
- What are their attributes and types (labels)?
  - Named Entity Recognition
- How are they related (edges)?
  - Relation Extraction
  - Semantic Role Labeling



## Two perspectives

- Graph Construction
- Who are the entities (nodes) in the graph?
  - Entity Linking
  - Entity Resolution
- What are their attributes and types (labels)?
  - Collective Classification
- How are they related (edges)?
  - Link Prediction

# What is NLP?



Unstructured

Ambiguous

Lots and lots of it!

Humans can read them, but

... very slowly

... can't remember all

... can't answer questions

Structured

Precise, Actionable

Specific to the task

Can be used for downstream applications, such as creating Knowledge Graphs!

# Knowledge Extraction

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

[Text](#)

NLP

Lennon..  
John Lennon...  
Pers

the Pool  
Locatio

Mrs. Lennon..  
.. his mother..  
Person

his father  
he Alfred  
Person

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

NNP VBD VBD IN

CC

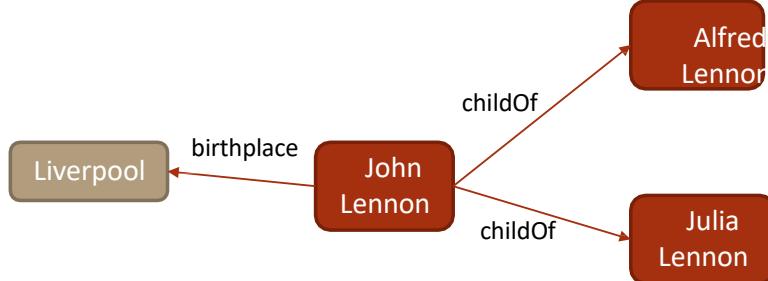
NNP

[Annotated text](#)

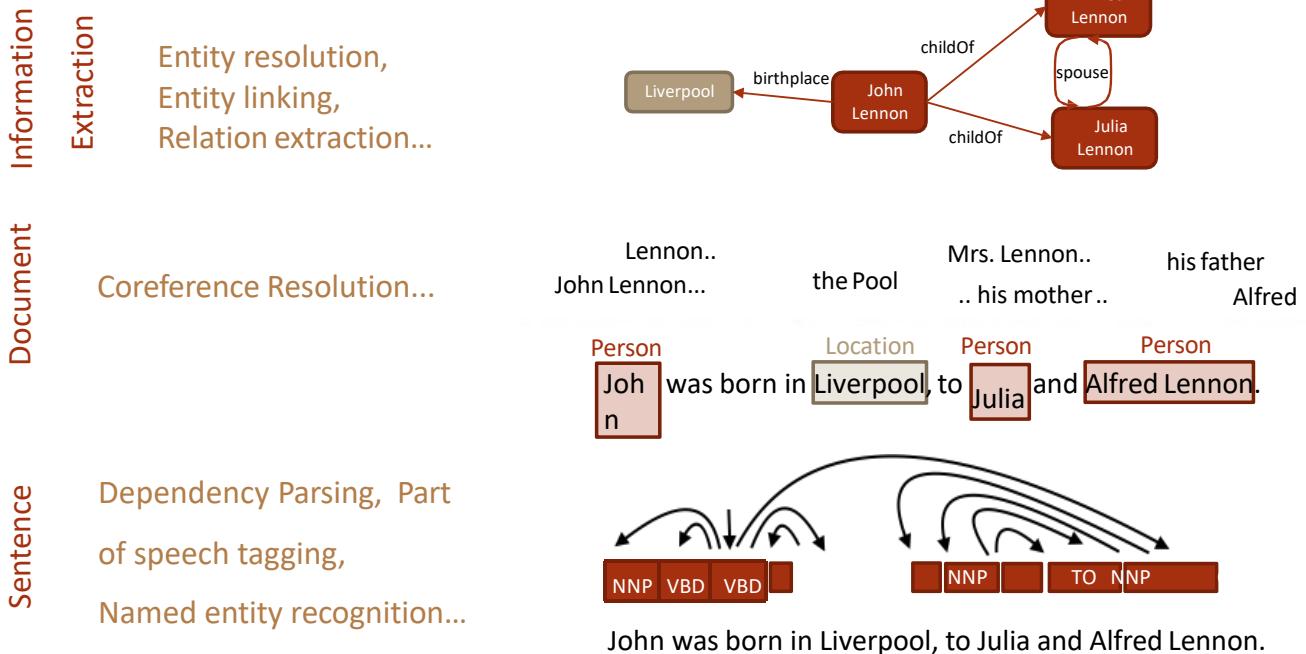
Information  
Extraction

Information  
Extraction

[Extraction graph](#)

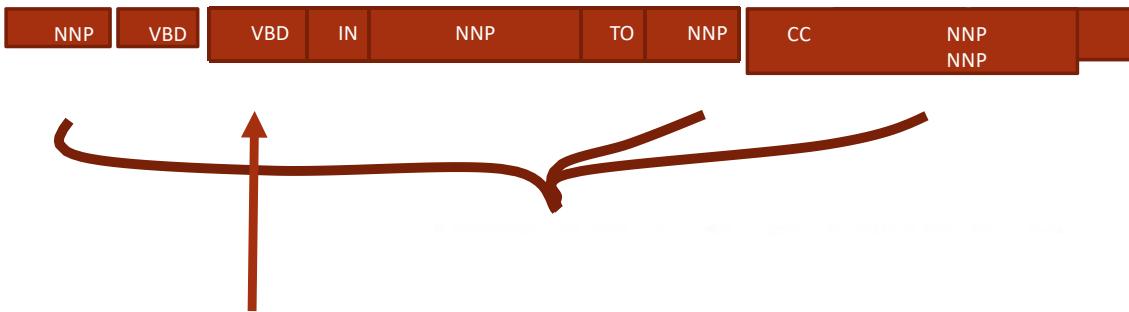


# Breaking it Down



# Tagging the Parts of Speech

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



Nouns are entities Verbs are relations

- Common approaches include CRFs, CNNs, LSTMs



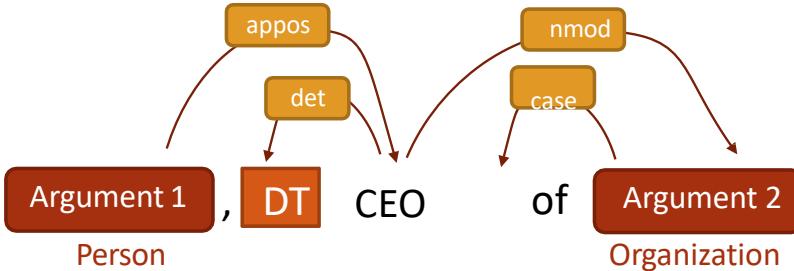
# Detecting Named Entities

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

- Structured prediction approaches
- Capture entity mentions and entity types

# 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...*



# Within-document Coreference

He... Mrs. Lennon.. Alfred  
Lennon.. .. his mother .. his  
the Pool father  
John Lennon... he  
John was born in Liverpool, to Julia and Alfred Lennon.

- Pairwise model for each noun/pronoun
  - Can consolidate information, provide context



# Entity Names: Two Main Problems

## Entities with Same Name

Same type of entities share names

Kevin Smith, John Smith, Springfield, ...

Things named after each other

Clinton, Washington, Paris, Amazon,

Princeton, Kingston, ...

Partial Reference

First names of people, Location instead  
of team name, Nick names

## Different Names for Entities

Nick Names

Bam Bam, Drumpf, ...

Typos/Misspellings

Baarak, Barak,  
Barrack, ...

Inconsistent References

MSFT, APPL, GOOG...



# Entity Linking Approach

Washington drops 10 points after game with UCLA Bruins.

## Candidate Generation

Washington DC, George Washington, Washington state,  
Lake Washington, Washington Huskies, Denzel Washington,  
University of Washington, Washington High School, ...

## Entity Types LOC/ORG

Washington DC, **George Washington**, Washington state,  
Lake Washington, Washington Huskies, **Denzel Washington**,  
University of Washington, Washington High School, ...

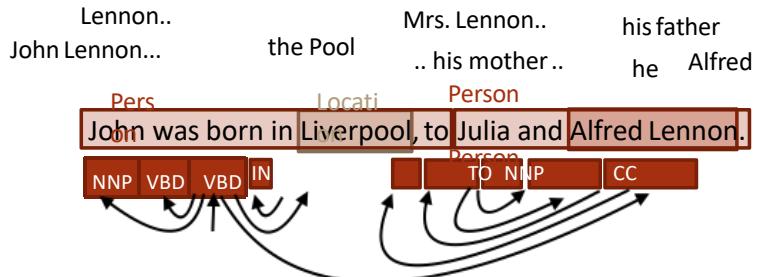
## Coreference UWashington, Huskies

**Washington DC, George Washington, Washington state,**  
**Lake Washington, Washington Huskies, Denzel Washington,**  
University of Washington, **Washington High School**, ...

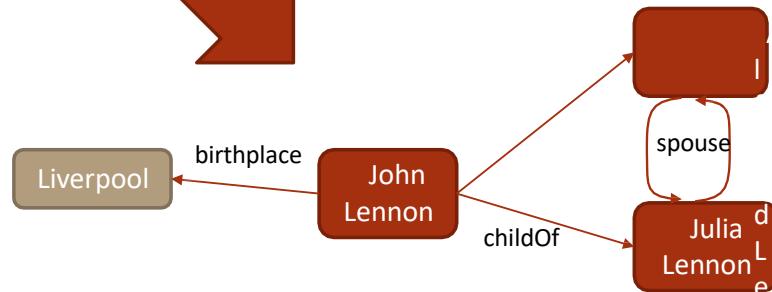
## Coherence UCLA Bruins, USC Trojans

**Washington DC, George Washington, Washington state,**  
**Lake Washington, Washington Huskies, Denzel Washington,**  
University of Washington, **Washington High School**, ...

# Information Extraction



Information Extraction



# Information Extraction

## 3 CONCRETE SUB-PROBLEMS

Defining domain

Learning extractors

Scoring the facts

## 3 LEVELS OF SUPERVISION

Supervised



Semi-supervised



Unsupervised



# Information Extraction

## 3 CONCRETE SUB-PROBLEMS

### Defining domain

Learning extractors

Scoring the facts



## 3 LEVELS OF SUPERVISION

Supervised



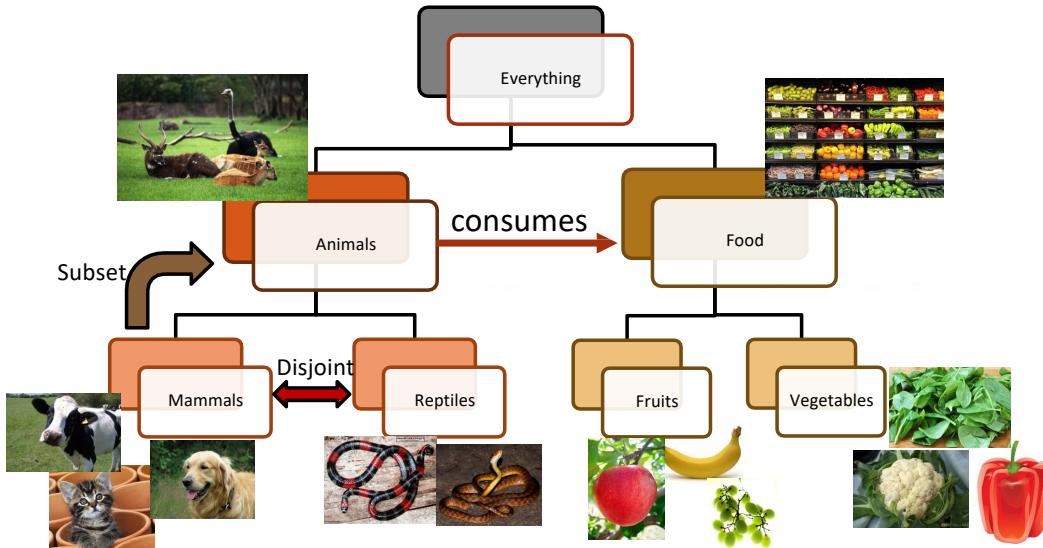
Semi-supervised



Unsupervised



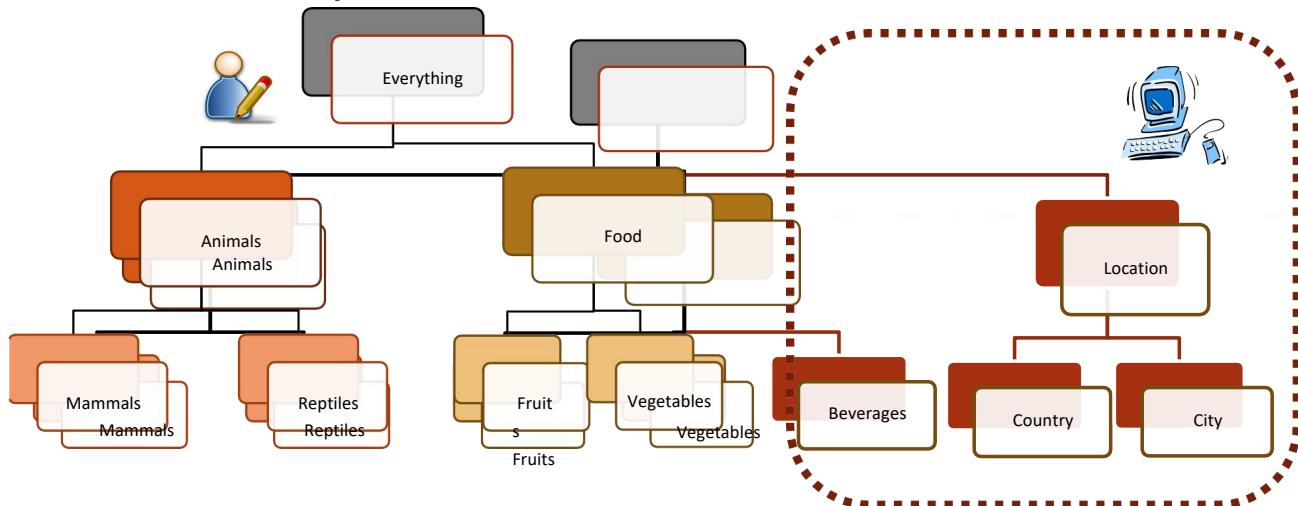
# Defining Domain: Manual



# Defining Domain: Semi-automatic



- Subset of types are manually defined
- SSL methods discover new types from unlabeled data





# Defining Domain: Automatic

---



- Any noun phrase is a candidate entity
  - Dog, cat, cow, reptile, mammal, apple, greens, mixed greens, lettuce, red leaf lettuce, romaine lettuce, iceberg lettuce...
- Any verb phrase is a candidate relation
  - Eats, feasts on, grazes, consumes,

# Information Extraction

## 3 CONCRETE SUB-PROBLEMS

Defining domain

## Learning extractors

Scoring candidate facts



## 3 LEVELS OF SUPERVISION

Supervised



Semi-supervised



Unsupervised



# Learning Extractors



- Supervised: high precision patterns
  - <PERSON> plays in <BAND>



- Semi-supervised: Bootstrapping to learn patterns
  - Create examples (**John Lennon, Beatles**), find patterns
  - Manually correct incorrect patterns



- Unsupervised: cluster phrases with constraints
  - Identify candidate verb phrases, find candidate arguments, cluster by NER types

# Information Extraction

## 3 CONCRETE SUB-PROBLEMS

Defining domain Learning  
extractors



## Scoring candidate facts

## 3 LEVELS OF SUPERVISION

Supervised



Semi-  
supervised



Unsupervised



# Scoring the candidate facts



- Human defined scoring function or  
Scoring function learnt using supervised ML with large amount of training data  
{expensive, high precision}



- Small amount of training data is available  
scoring refined over multiple iterations using both labeled and unlabeled data



- Completely automatic (Self-training)  
Confidence(extraction pattern)  $\propto$  (#unique instances it could extract)  
Score(candidate fact)  $\propto$  (#distinct extraction patterns that support it)  
{cheap, leads to semantic drift}

# IE systems in practice

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



# Knowledge Extraction: Key Points

---

- Built on the foundation of NLP techniques
  - Part-of-speech tagging, dependency parsing, named entity recognition, coreference resolution...
  - Challenging problems with very useful outputs
- Information extraction techniques use NLP to:
  - define the domain
  - extract entities and relations
  - score candidate outputs
- Trade-off between manual & automatic methods

# Graph Construction Issues

Extracted knowledge is:

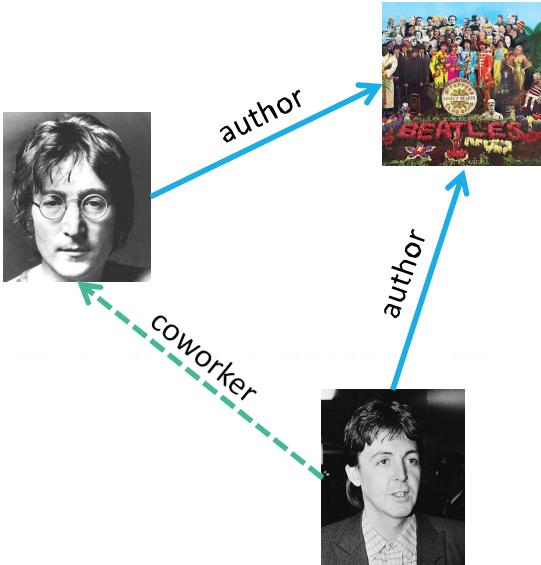
- ambiguous:
  - Ex: Beetles, beetles, Beatles
  - Ex: citizenOf, livedIn, bornIn



# Graph Construction Issues

Extracted knowledge is:

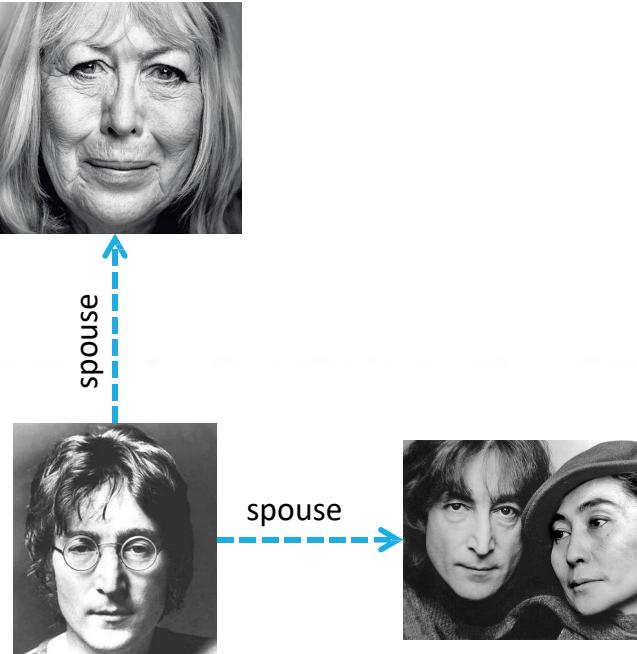
- ambiguous
- incomplete
  - Ex: missing relationships
  - Ex: missing labels
  - Ex: missing entities



# Graph Construction Issues

Extracted knowledge is:

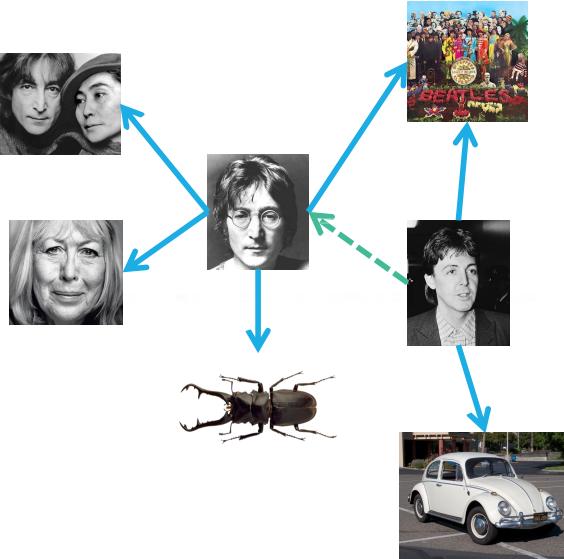
- ambiguous
- incomplete
- inconsistent
  - Ex: Cynthia Lennon, Yoko Ono
  - Ex: exclusive labels (alive, dead)
  - Ex: domain-range constraints



# Graph Construction Issues

Extracted knowledge is:

- ambiguous
- incomplete
- inconsistent

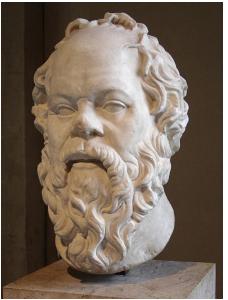




# Graph Construction approach

- Graph construction **cleans** and **completes** extraction graph
- Incorporate ontological constraints and relational patterns
- Discover statistical relationships within knowledge graph

# Beyond Pure Reasoning



- Classical AI approach to knowledge: reasoning

$\text{Lbl}(\text{Socrates}, \text{Man}) \ \& \ \text{Sub}(\text{Man}, \text{Mortal}) \rightarrow \text{Lbl}(\text{Socrates}, \text{Mortal})$

# Beyond Pure Reasoning



- Classical AI approach to knowledge: reasoning

$\text{Lbl}(\text{Socrates}, \text{Man}) \ \& \ \text{Sub}(\text{Man}, \text{Mortal}) \rightarrow \text{Lbl}(\text{Socrates}, \text{Mortal})$

- Reasoning difficult when extracted knowledge has errors

# Beyond Pure Reasoning



- Classical AI approach to knowledge: reasoning

$\text{Lbl}(\text{Socrates}, \text{Man}) \& \text{Sub}(\text{Man}, \text{Mortal}) \rightarrow \text{Lbl}(\text{Socrates}, \text{Mortal})$

- Reasoning difficult when extracted knowledge has errors
- Solution: probabilistic models

$P(\text{Lbl}(\text{Socrates}, \text{Mortal}) | \text{Lbl}(\text{Socrates}, \text{Man})) = 0.9$



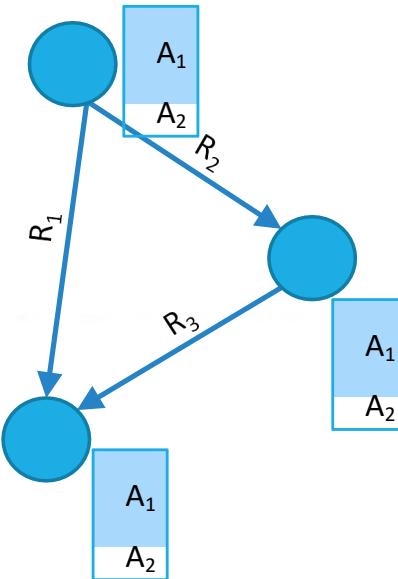
# Graphical Models: Overview

- Define **joint probability distribution** on knowledge graphs
- Each candidate fact in the knowledge graph is a **variable**
- Statistical signals, ontological knowledge and rules parameterize the **dependencies** between variables
- Find most likely knowledge graph by **optimization/sampling**

# Knowledge Graph Identification

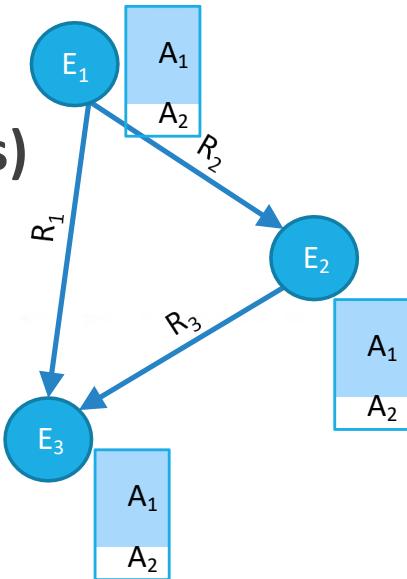
Define a graphical model to perform all three of these tasks simultaneously!

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



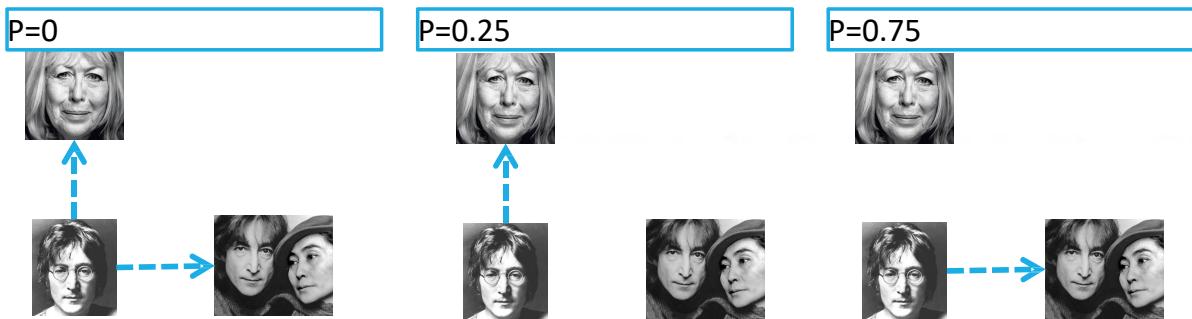
# Knowledge Graph Identification

P(Who, What, How | Extractions)



# Probabilistic Models

- Use dependencies between facts in KG
- Probability defined *jointly* over facts





# What determines probability?

- Statistical signals from text extractors and classifiers



# What determines probability?

- **Statistical signals from text extractors and classifiers**
  - $P(R(\text{John}, \text{Spouse}, \text{Yoko}))=0.75;$   
 $P(R(\text{John}, \text{Spouse}, \text{Cynthia}))=0.25$
  - LevenshteinSimilarity(Beatles, Beetles) = 0.9



# What determines probability?

- Statistical signals from text extractors and classifiers
- Ontological knowledge about domain



# What determines probability?

- Statistical signals from text extractors and classifiers
- **Ontological knowledge about domain**
  - $\text{Functional}(\text{Spouse}) \ \& \ R(A, \text{Spouse}, B) \rightarrow !R(A, \text{Spouse}, C)$
  - $\text{Range}(\text{Spouse}, \text{Person}) \ \& \ R(A, \text{Spouse}, B) \rightarrow \text{Type}(B, \text{Person})$



# What determines probability?

- Statistical signals from text extractors and classifiers
- Ontological knowledge about domain
- Rules and patterns mined from data



# What determines probability?

- Statistical signals from text extractors and classifiers
- Ontological knowledge about domain
- Rules and patterns mined from data
  - $R(A, \text{Spouse}, B) \& R(A, \text{Lives}, L) \rightarrow R(B, \text{Lives}, L)$
  - $R(A, \text{Spouse}, B) \& R(A, \text{Child}, C) \rightarrow R(B, \text{Child}, C)$



# What determines probability?

- Statistical signals from text extractors and classifiers
  - $P(R(\text{John}, \text{Spouse}, \text{Yoko})) = 0.75$ ;  
 $P(R(\text{John}, \text{Spouse}, \text{Cynthia})) = 0.25$
  - LevenshteinSimilarity(Beatles, Beetles) = 0.9
- Ontological knowledge about domain
  - Functional(Spouse) & R(A, Spouse, B) -> !R(A, Spouse, C)
  - Range(Spouse, Person) & R(A, Spouse, B) -> Type(B, Person)
- Rules and patterns mined from data
  - R(A, Spouse, B) & R(A, Lives, L) -> R(B, Lives, L)
  - R(A, Spouse, B) & R(A, Child, C) -> R(B, Child, C)



# Graphical Models: Pros/Cons

## BENEFITS

- Define probability distribution over KGs
- Easily specified via rules
- Fuse knowledge from many different sources

## DRAWBACKS

- Requires optimization over all KG facts - overkill
- Dependent on rules from ontology/expert
- Require probabilistic semantics - unavailable



# Random Walk Overview

- Given: a query of an **entity** and **relation**
- Starting at the entity, **randomly walk** the KG
- Random walk ends when reaching an appropriate **goal**
- Learned **parameters** bias choices in the random walk
- Output **relative probabilities** of goal states

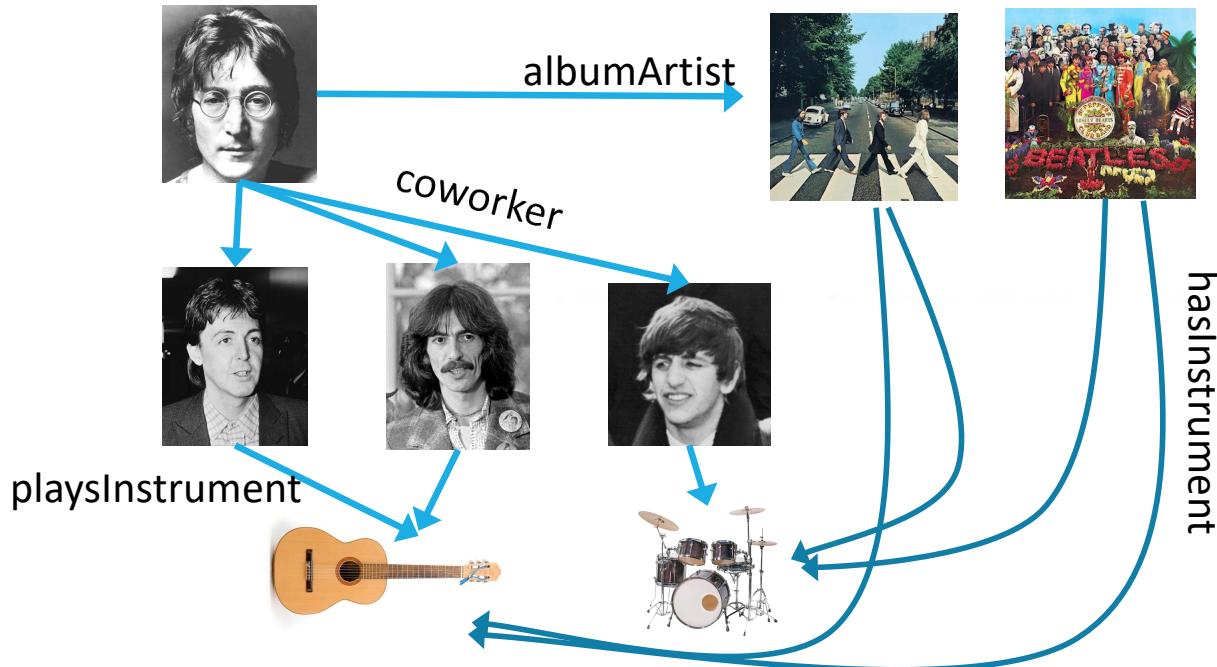


# Random Walk Illustration

Query: R(Lennon, PlaysInstrument, ?)

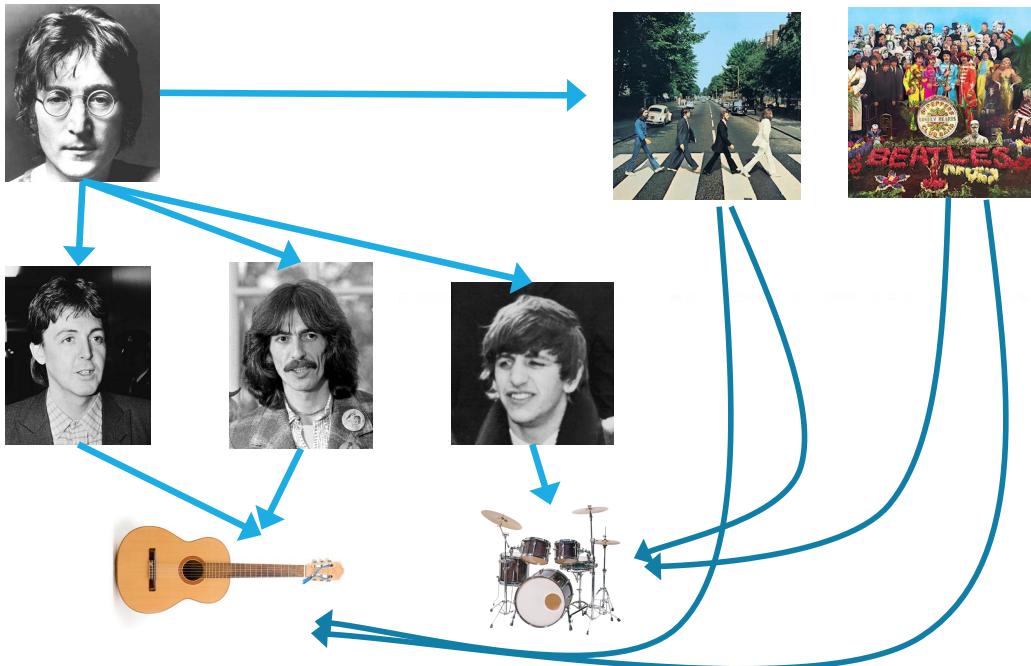
# Random Walk Illustration

Query:  $R(\text{Lennon}, \text{PlaysInstrument}, ?)$



# Random Walk Illustration

Query: R(Lennon, PlaysInstrument, ?)





# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



Path  
Weight of path

$P(Q | \theta = \langle \text{coworker}, \text{playsInstrument} \rangle) W$

# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



$P(Q | \theta = \langle \text{coworker}, \text{playsInstrument} \rangle) W$

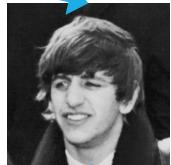


# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)

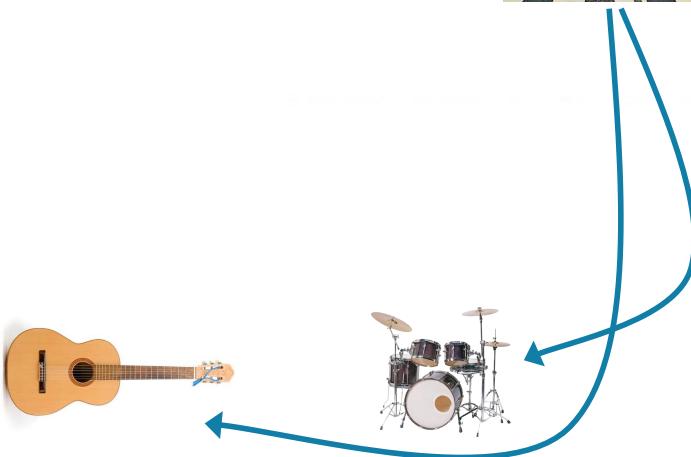
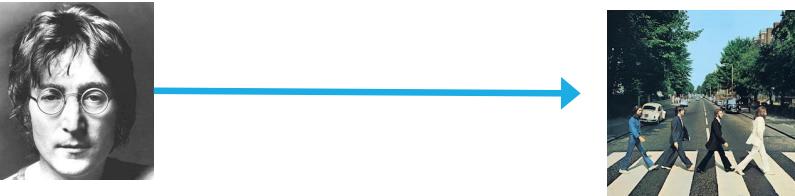


$P(Q | \theta = \langle \text{coworker}, \text{playsInstrument} \rangle) W$



# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)



$P(Q | \theta = \langle \text{albumArtist}, \text{hasInstrument} \rangle) W$



# Random Walk Illustration

Query Q: R(Lennon, PlaysInstrument, ?)

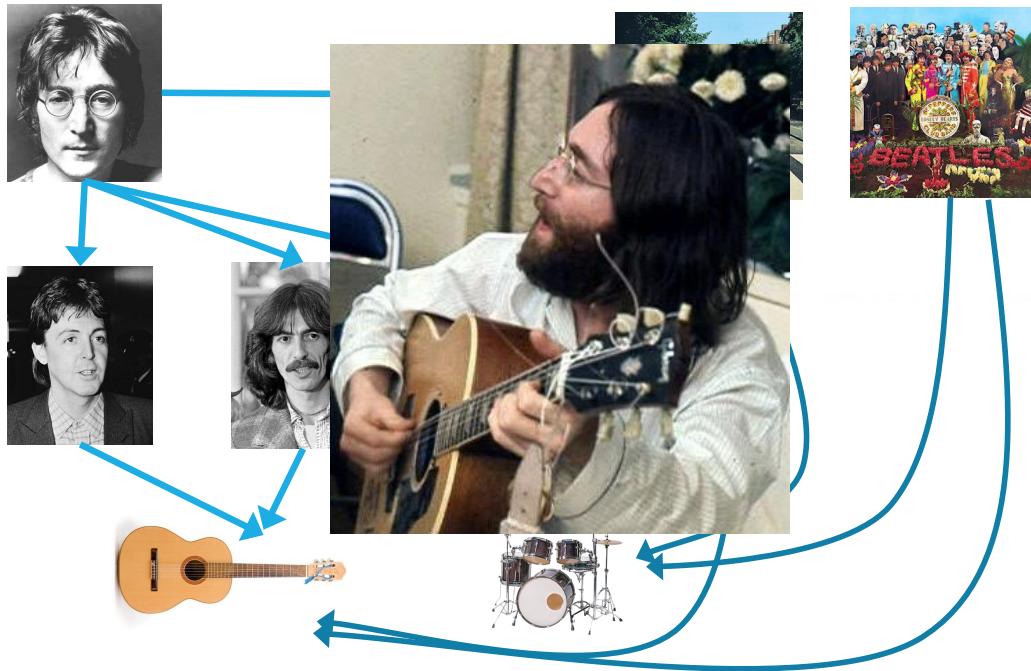


$P(Q | \theta = \langle \text{albumArtist}, \text{hasInstrument} \rangle) W$



# Random Walk Illustration

Query: R(Lennon, PlaysInstrument, ?)





# Random Walks: Pros/Cons

## BENEFITS

- KG query estimation independent of KG size
- Model training produces interpretable, logical rules
- Robust to noisy extractions through probabilistic form

## DRAWBACKS

- Full KG completion task inefficient
- Training data difficult to obtain at scale
- Input must follow probabilistic semantics



# Two classes of Probabilistic Models

---

- GRAPHICAL MODELS
  - Possible facts in KG are variables
  - Logical rules relate facts
  - Probability rules
  - Universally-quantified



## Two classes of Probabilistic Models

---

- RANDOM WALK METHODS
  - Possible facts posed as queries
  - Random walks of the KG constitute “proofs”
  - Probability path lengths/transitions
  - Locally grounded



# Probabilistic Models: Downsides

## Limitation to Logical Relations

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

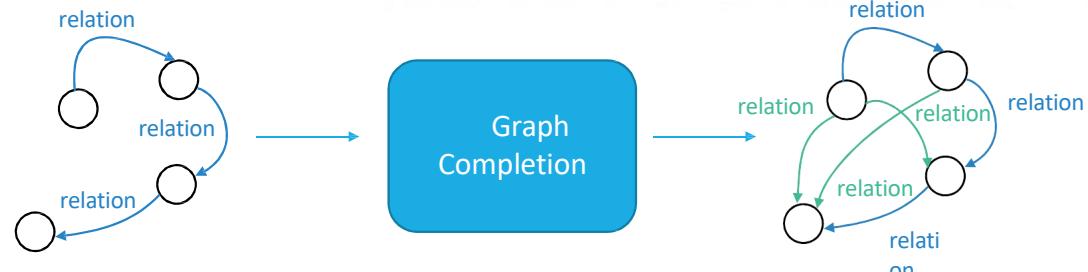
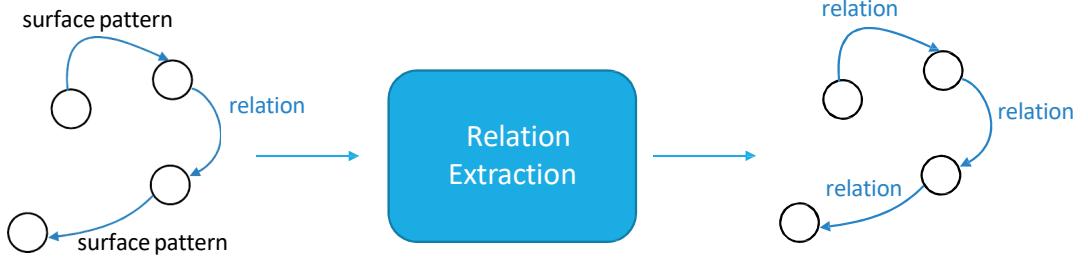
## Embeddings

- **Everything as dense vectors**
- **Can capture many relations**
- **Learned from data**
- **Complexity depends on latent dimensions**
- **Learning using stochastic gradient, back-propagation**
- **Querying is often cheap**
- **GPU-parallelism friendly**

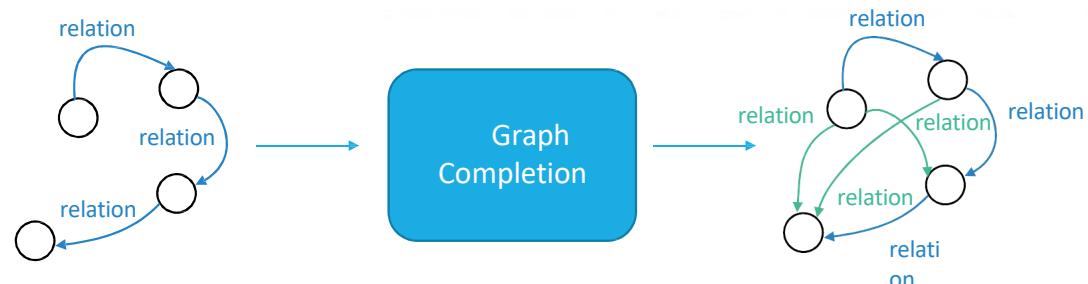
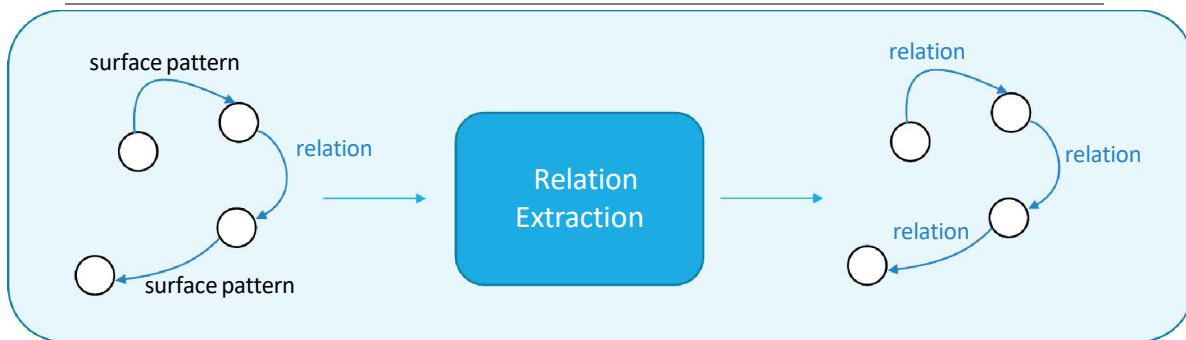
## Computational Complexity of Algorithms

- Complexity depends on explicit dimensionality
  - Often NP-Hard, in size of data
  - More rules, more expensive inference
- Query-time inference is sometimes NP-Hard
- Not trivial to parallelize, or use GPUs

# Two Related Tasks



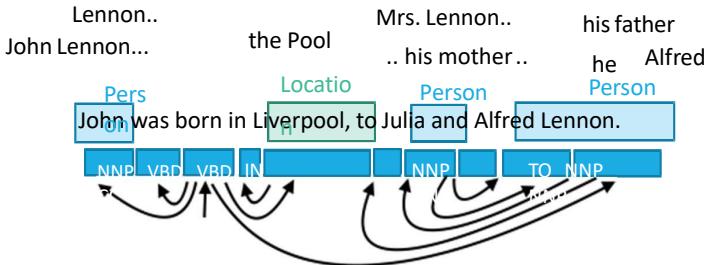
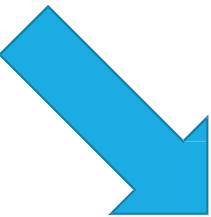
# Two Related Tasks



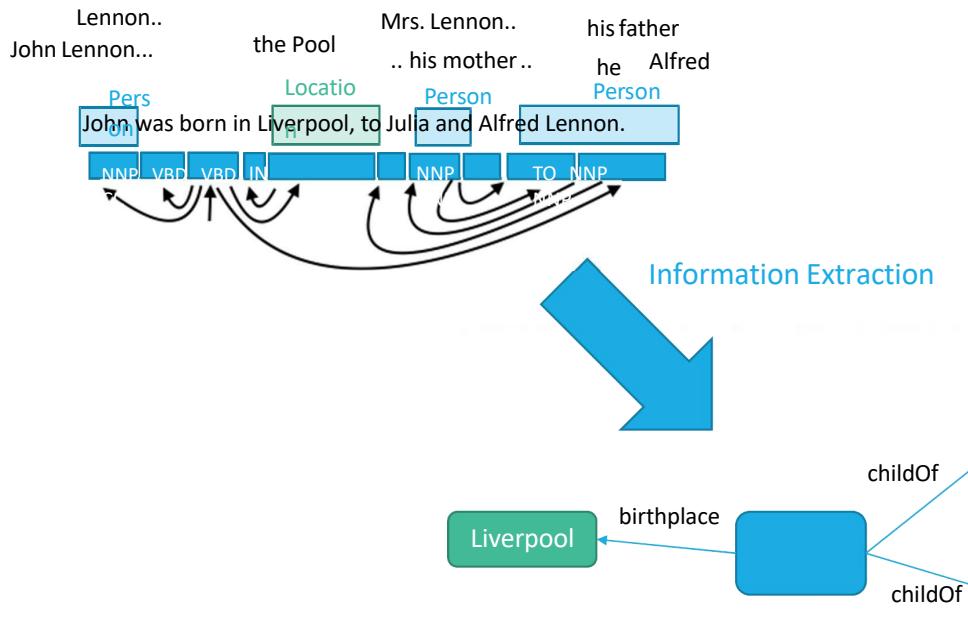
# What is NLP?

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

Natural Language  
Processing

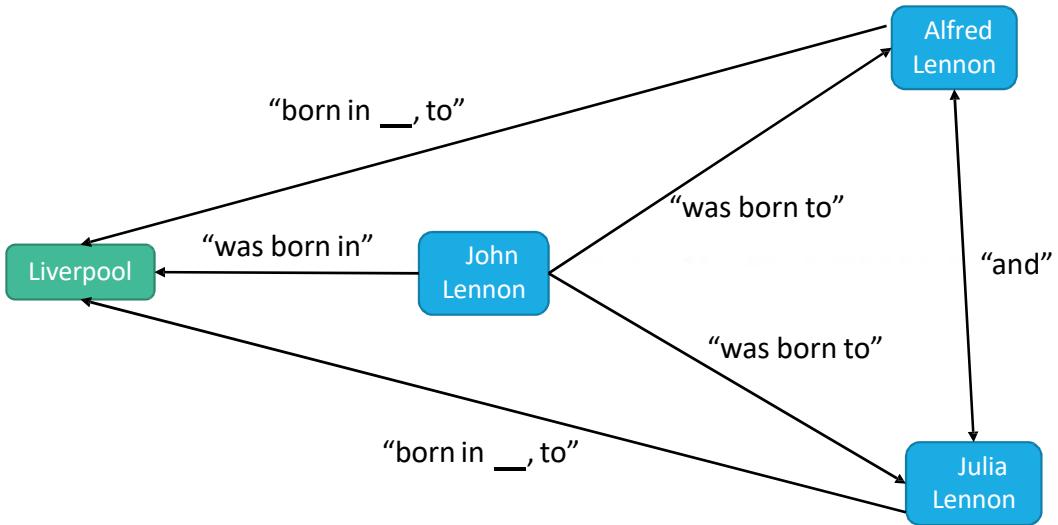


# What is Information Extraction?



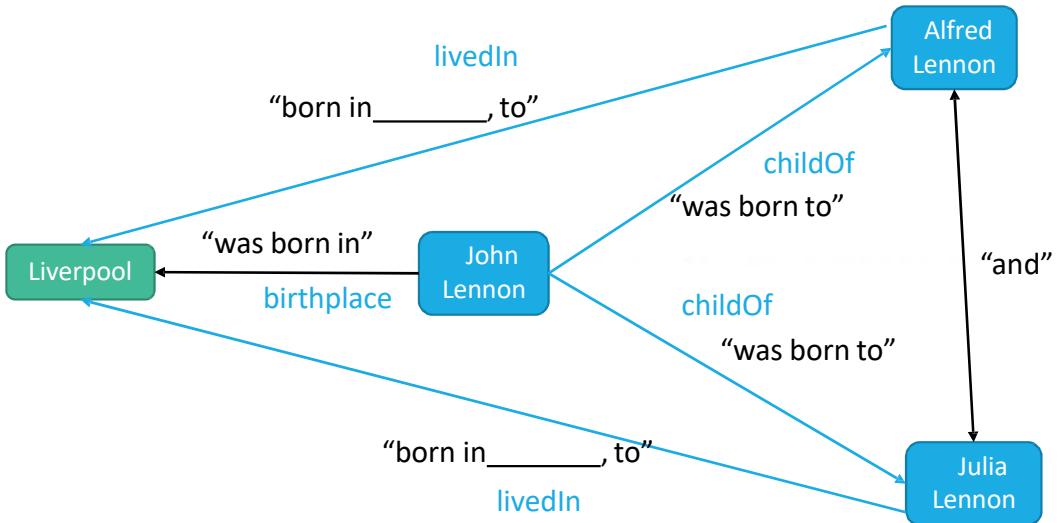
# Relation Extraction From Text

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

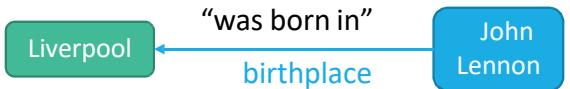


# Relation Extraction From Text

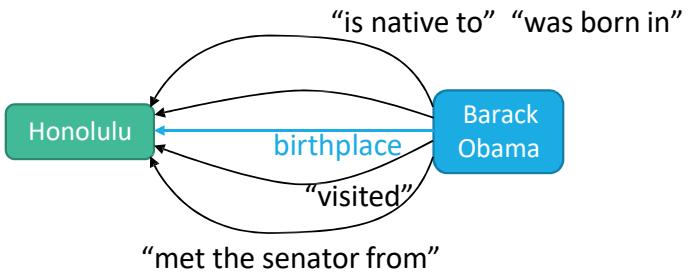
John was born in Liverpool, to Julia and Alfred Lennon.



# “Distant” Supervision



No direct supervision gives us this information. **Supervised**:  
Too expensive to label sentences **Rule-based**: Too much variety in language  
Both only work for a small set of relations, i.e. 10s, not 100s

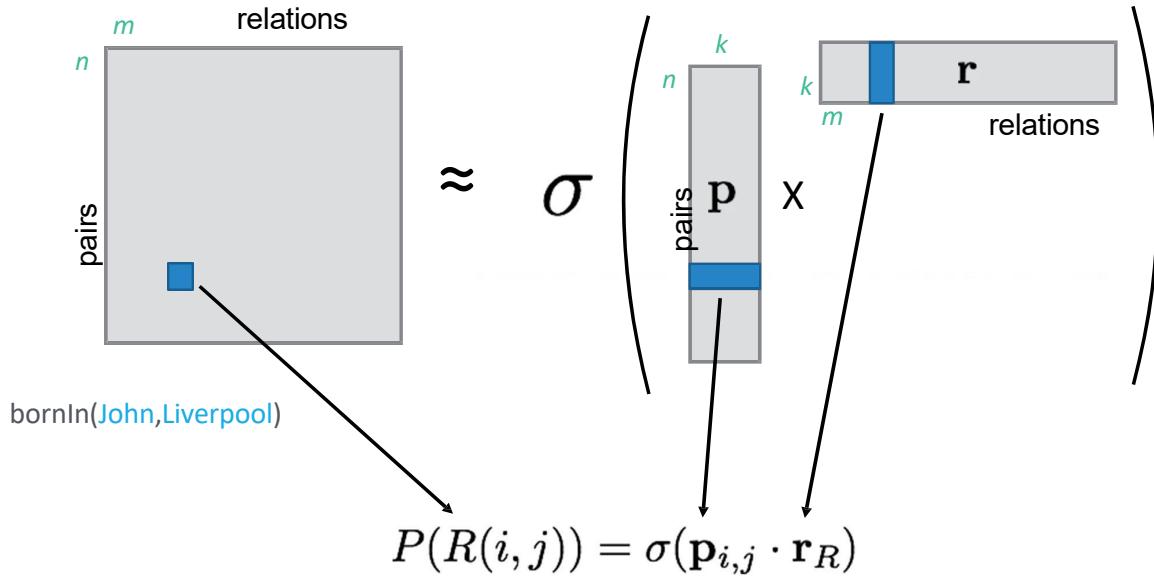


# Relation Extraction as a Matrix

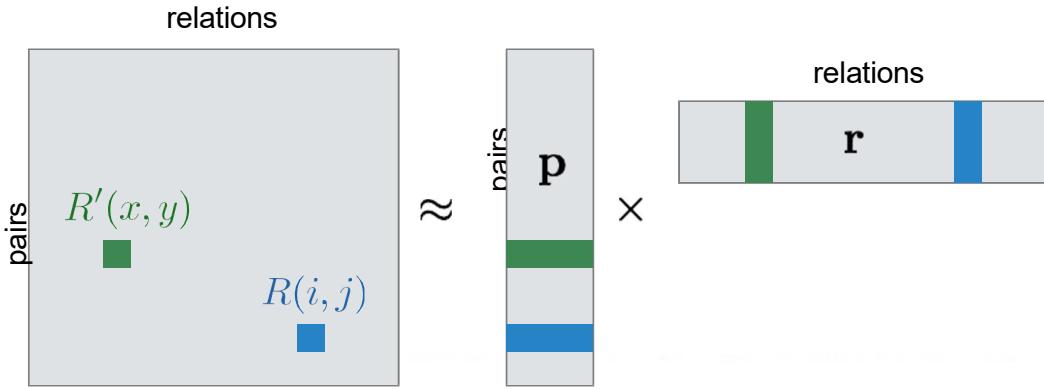
John was born in Liverpool, to Julia and Alfred Lennon.

	<i>was born in</i> <i>r-hsubjpass-born&lt;-nmod:in</i>	<i>was born to</i>	<i>and</i>	<i>birthplace(X,Y)</i>	<i>spouse(X,Y)</i>
John Lennon, Liverpool	1				?
John Lennon, Julia Lennon		1			
John Lennon, Alfred Lennon			1		
Julia Lennon, Alfred Lennon				?	
Barack Obama, Hawaii	1			1	
Barack Obama, Michelle Obama			1		1

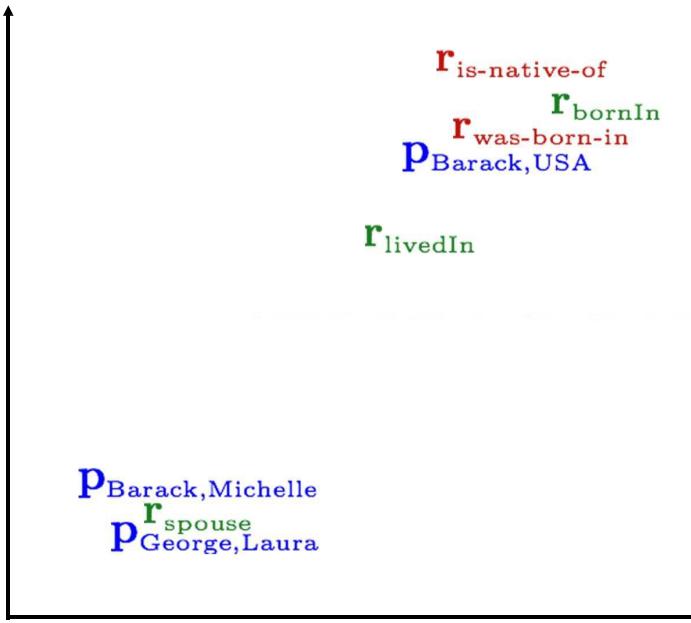
# Matrix Factorization



# Training: Stochastic Updates



# Relation Embeddings





# Embeddings ~ Logical Relations

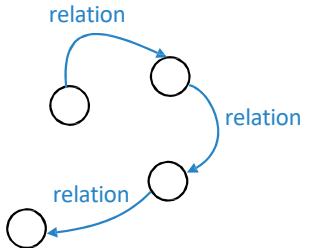
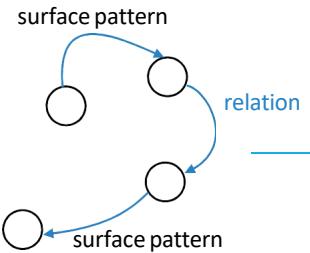
## Relation Embeddings, $w$

- Similar embedding for 2 relations denote they are paraphrases
  - `is married to`, `spouseOf(X,Y)`, `/person/spouse`
- One embedding can be contained by another
  - $w(\text{topEmployeeOf}) \subset w(\text{employeeOf})$
  - $\text{topEmployeeOf}(X,Y) \rightarrow \text{employeeOf}(X,Y)$
- Can capture logical patterns, without needing to specify them!

## Entity Pair Embeddings, $v$

Similar entity pairs denote similar relations between them  
Entity pairs may describe multiple “relations”  
independent `foundedBy` and `employeeOf` relations

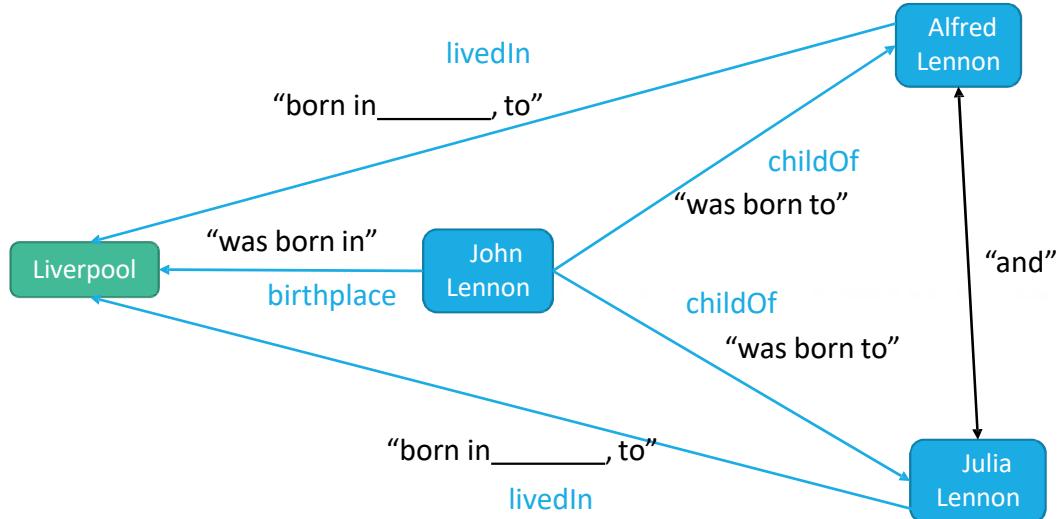
# Two Related Tasks



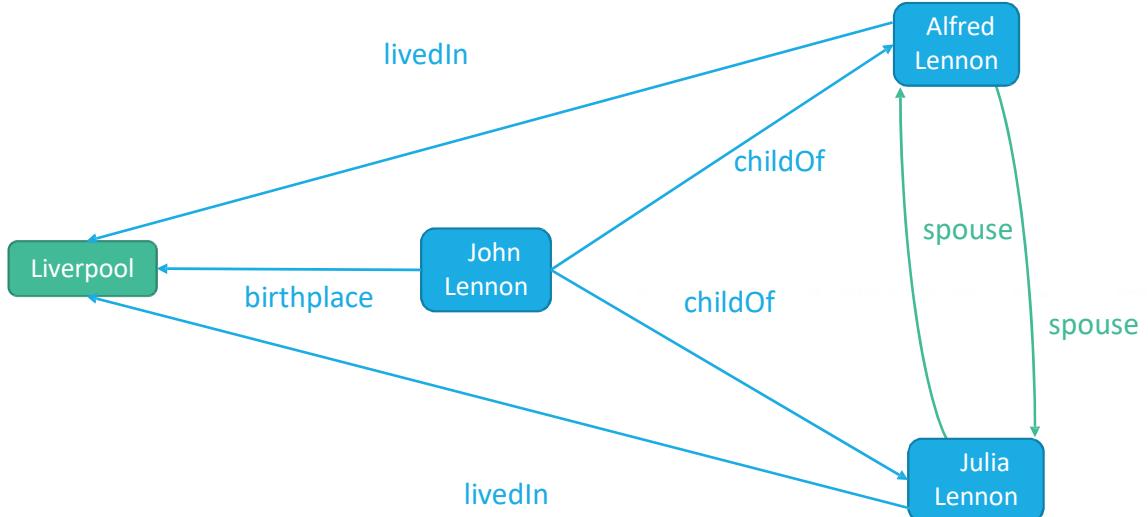
Relation  
Extraction

Graph  
Completion

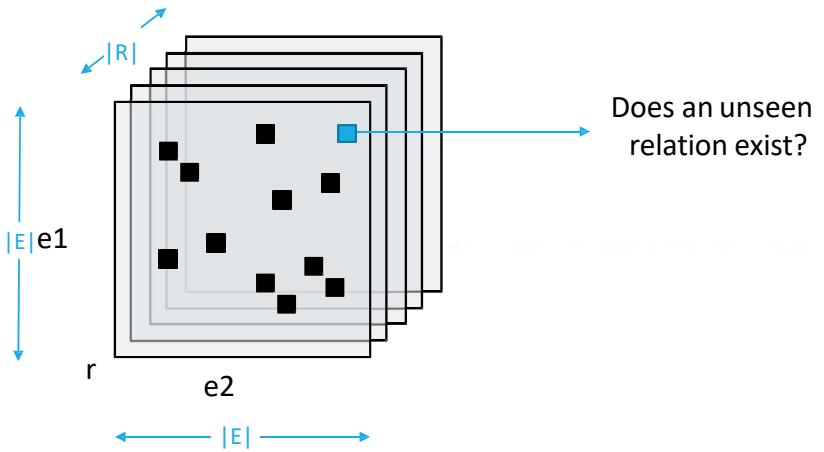
# Graph Completion



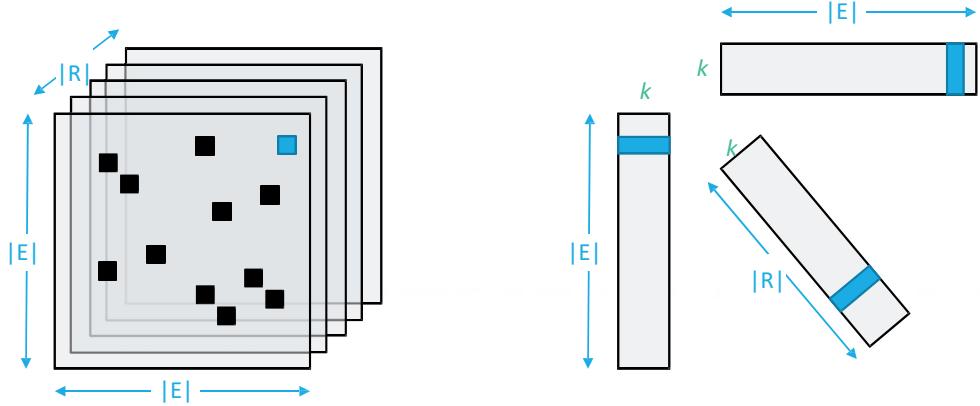
# Graph Completion



# Tensor Formulation of KG



# Factorize that Tensor



$$S(r(a, b)) = f(\mathbf{v}_r, \mathbf{v}_a, \mathbf{v}_b)$$



# Many Different Factorizations

CANDECOMP/PARAFAC-Decomposition

$$S(r(a, b)) = \sum_k R_{r,k} \cdot e_{a,k} \cdot e_{b,k}$$

Tucker2 and RESCAL Decompositions

$$S(r(a, b)) = (\mathbf{R}_r \times \mathbf{e}_a) \times \mathbf{e}_b$$

Model E

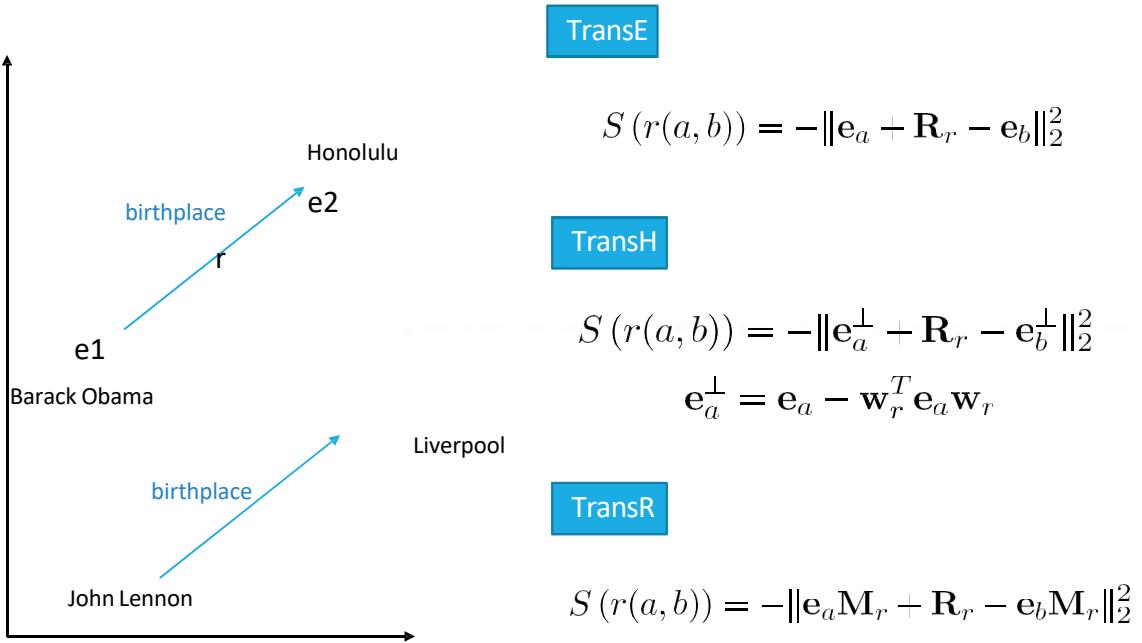
$$S(r(a, b)) = \mathbf{R}_{r,1} \cdot \mathbf{e}_a + \mathbf{R}_{r,2} \cdot \mathbf{e}_b$$

Holographic Embeddings

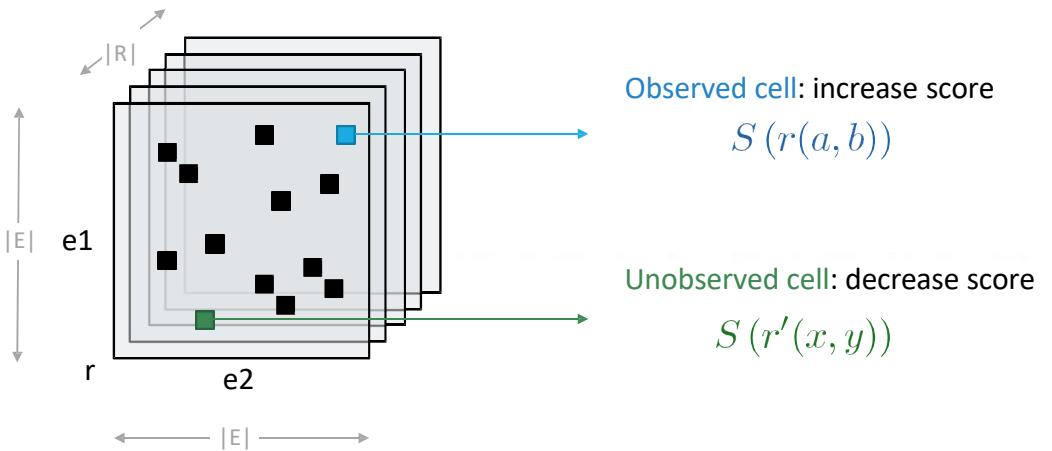
$$S(r(a, b)) = \mathbf{R}_r \times (\mathbf{e}_a \star \mathbf{e}_b)$$

Not tensor  
factorization  
(per se)

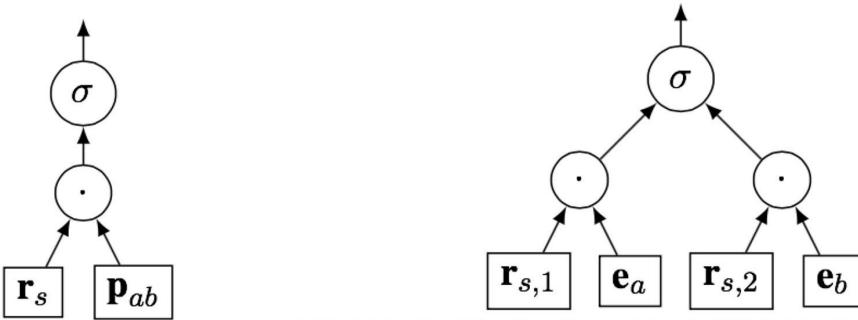
# Translation Embeddings



# Parameter Estimation

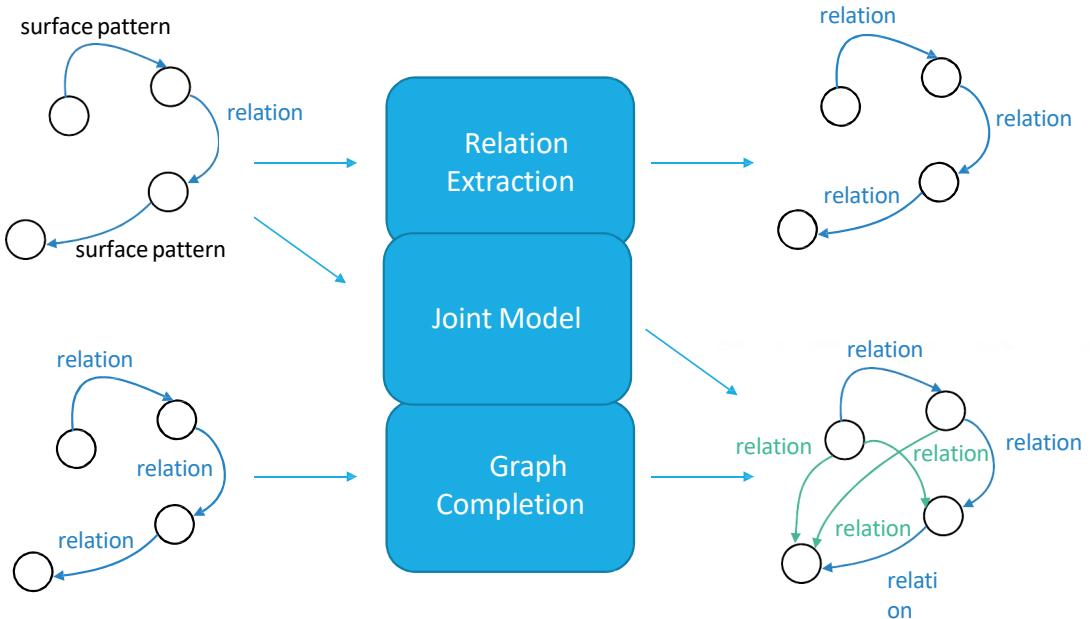


# Matrix vs Tensor Factorization



- Vectors for each entity pair
- Can only predict for entity pairs that appear in text together
- No sharing for same entity in different entity pairs
- Vectors for each entity
- Assume entity pairs are “low-rank”
  - But many relations are not!
  - Spouse: you can have only  $\sim 1$
- Cannot learn pair specific information

# Joint Extraction+Completion





# Review: Embedding Techniques

## Two Related Tasks:

- Relation Extraction from Text
- Graph (or Link) Completion

## Relation Extraction:

- Matrix Factorization  
Approaches

## Graph Completion:

- Tensor Factorization  
Approaches



# Overview

---



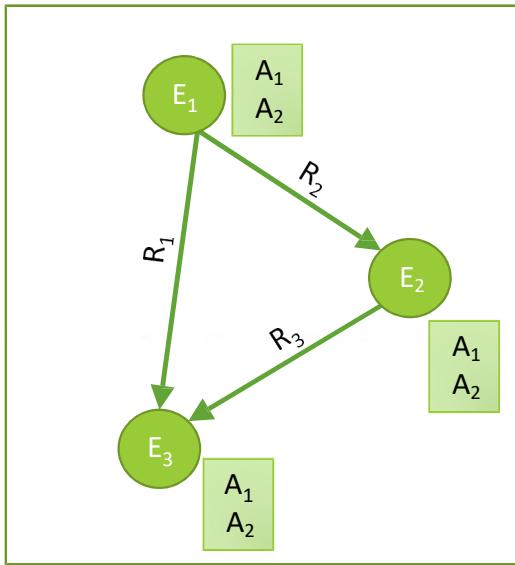
## Why do we need Knowledgegraphs?

---

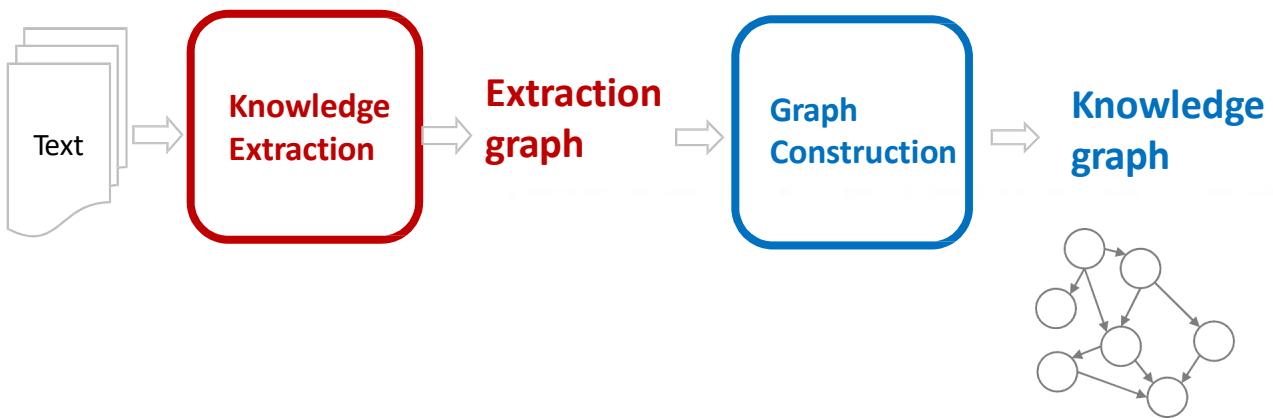
- Humans can explore large database in intuitive ways
- AI agents get access to human common sense knowledge

# 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
<b>Who are the entities? (nodes)</b>	<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>
<b>What are their attributes? (labels)</b>	<ul style="list-style-type: none"><li>• Entity Typing</li></ul>	<ul style="list-style-type: none"><li>• Collective classification</li></ul>
<b>How are they related? (edges)</b>	<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

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

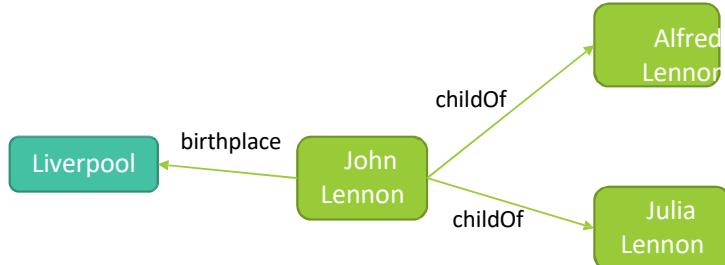
Pers	Locatio	Person	Person
John	Liverpool	Julia	Alfred
NNP	VRD	VBD	IN
		NNP	TO NNP CC

The diagram illustrates the flow of air over a curved surface. The air flows from left to right, indicated by arrows pointing along the surface. As the surface curves upwards, the air follows the curve initially. However, as the angle of attack increases, the air's path becomes increasingly vertical relative to the surface. This results in a point where the air's velocity parallel to the surface becomes zero, marking the point of boundary layer separation. After separation, the air no longer follows the surface's curvature, leading to a wake behind the curve.

## Annotated text

# Information Extraction

## Extraction graph



# Information Extraction

## Single extractor

Defining domain

Learning extractors

Scoring candidate facts

Supervised



Semi-supervised

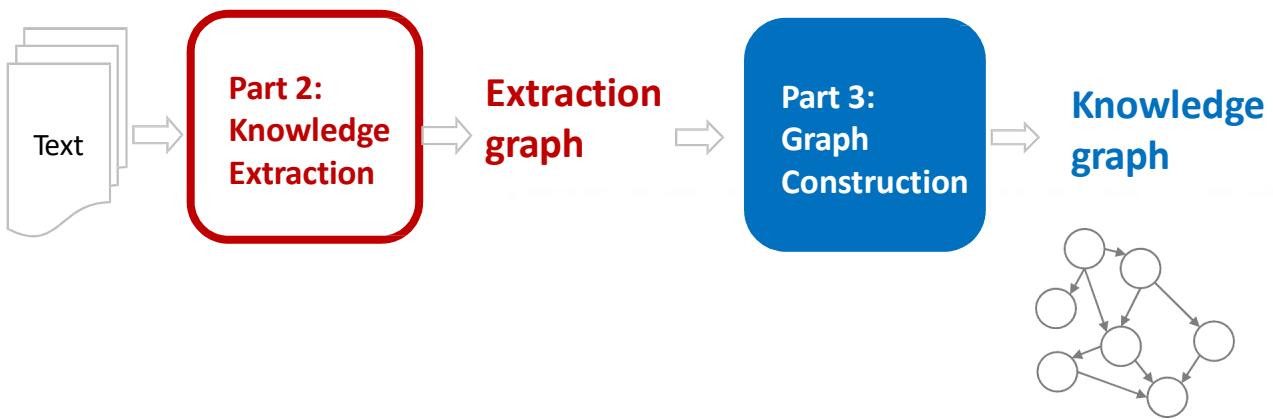


Unsupervised



## Fusing multiple extractors

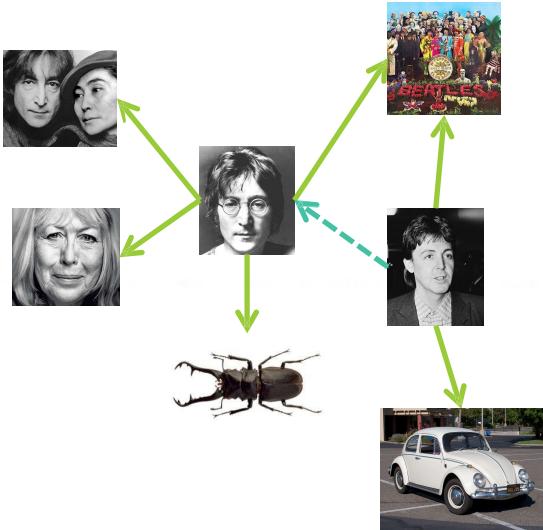
# Knowledge Graph Construction



# 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



# Relation Embeddings

