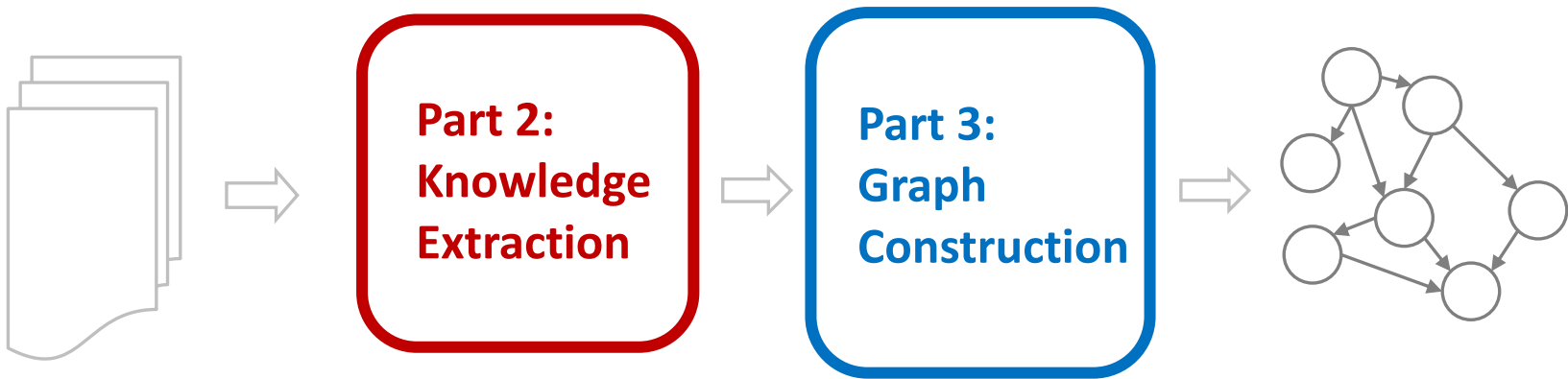
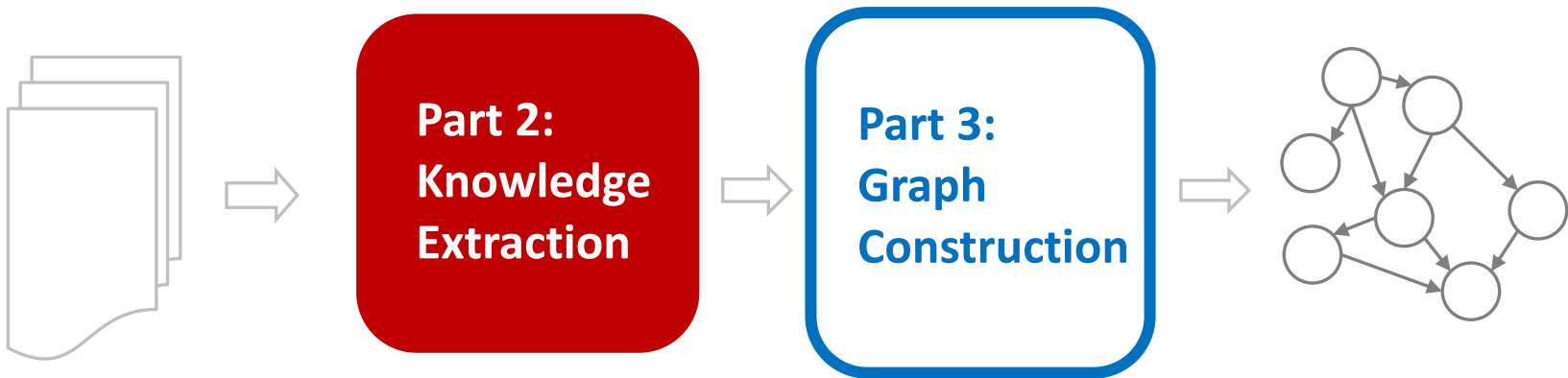


Part 1: Knowledge Graphs



Part 4: Critical Analysis

Part 1: Knowledge Graphs

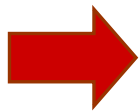
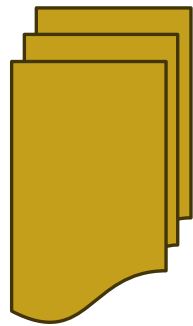


Part 4: Critical Analysis

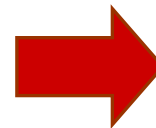
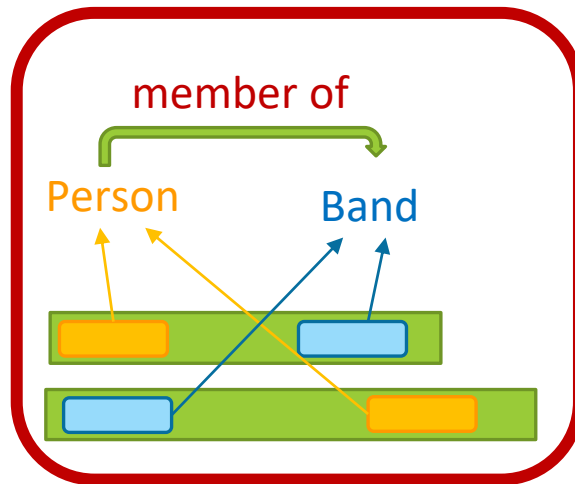
Information Extraction

Input

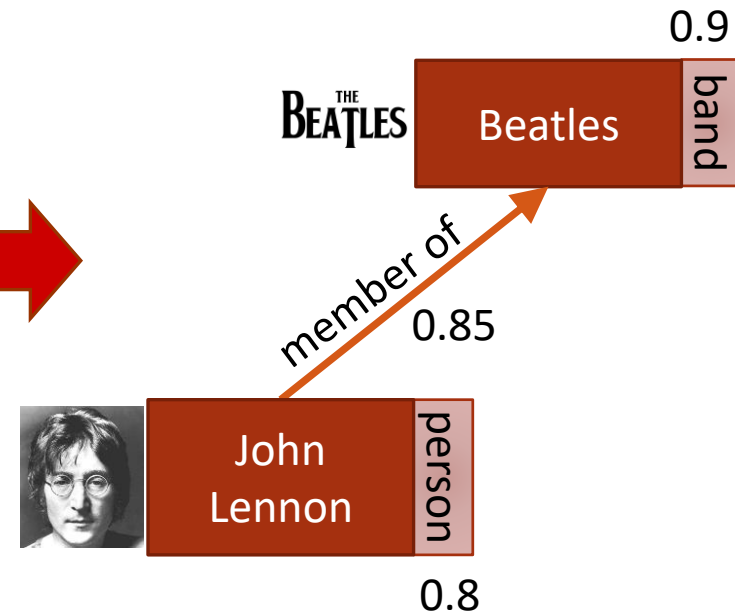
Text corpus



**Information
Extraction**



Output



Information Extraction

3 IMPORTANT SUB-PROBLEMS

(DEFINE DOMAIN, LEARN EXTRACTORS, SCORE EXTRACTIONS)

3 LEVELS OF SUPERVISION

(MANUAL, SEMI-SUPERVISED, UNSUPERVISED)

KNOWLEDGE FUSION WITH MULTIPLE EXTRACTORS

(CO-TRAINING, MULTI-VIEW LEARNING)

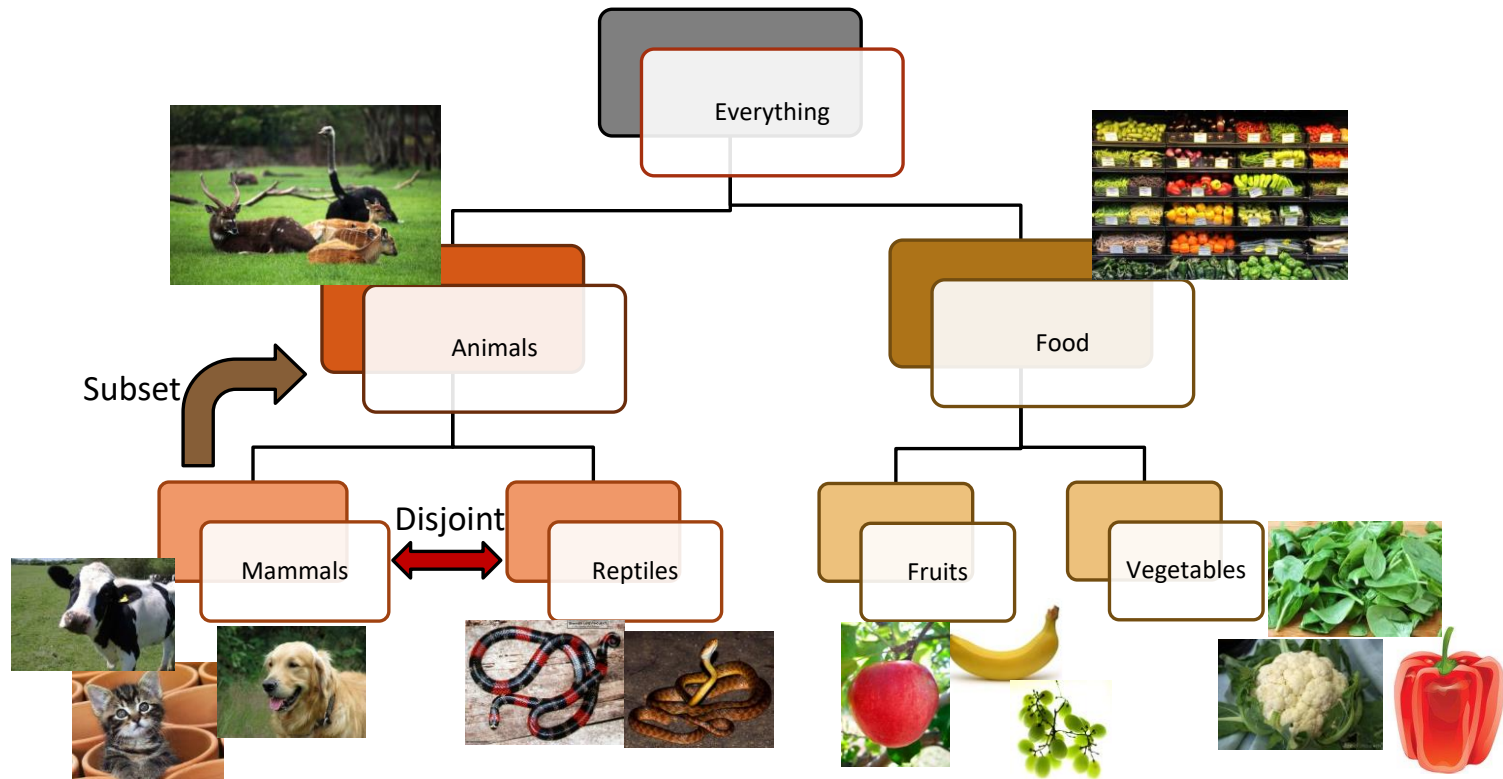
EXAMPLE IE SYSTEMS

Information Extraction



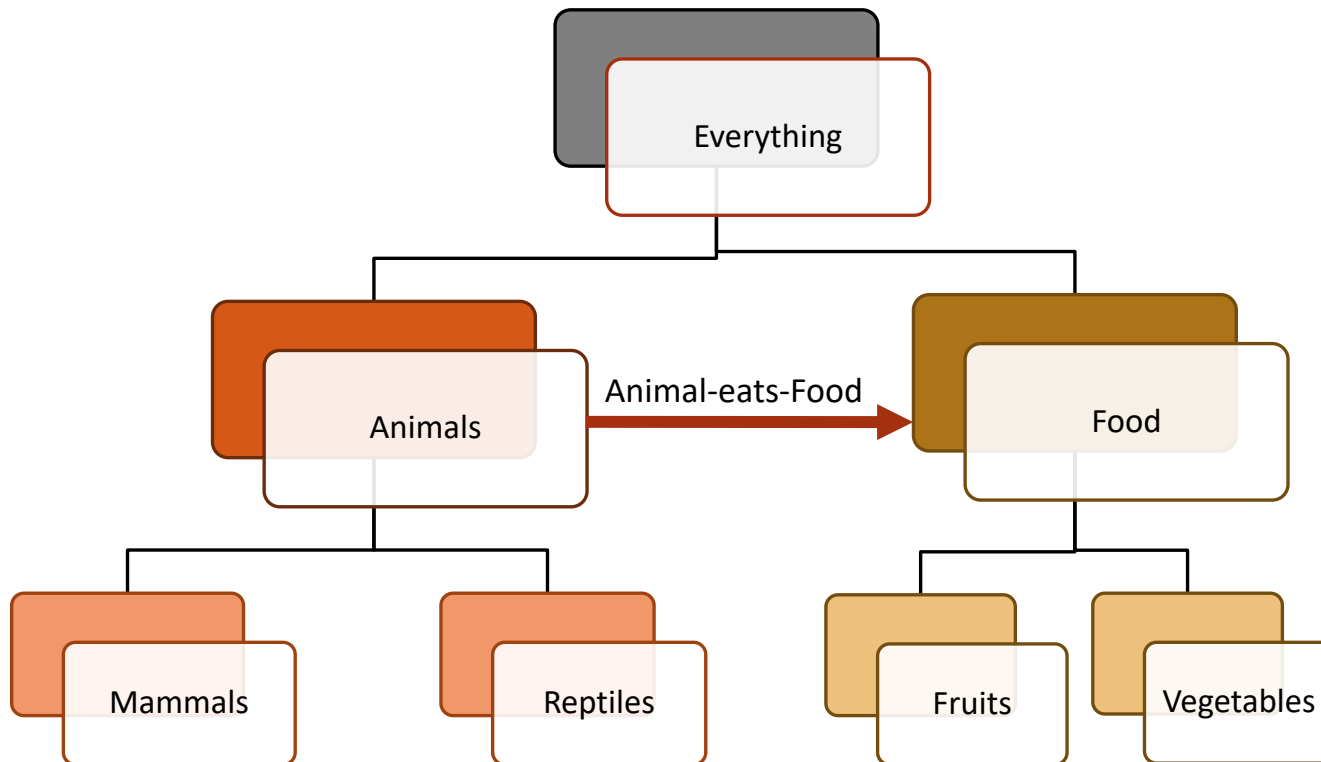
Defining domain:
types/relations of interest

Defining Domain: Manual



Defining Domain: Manual

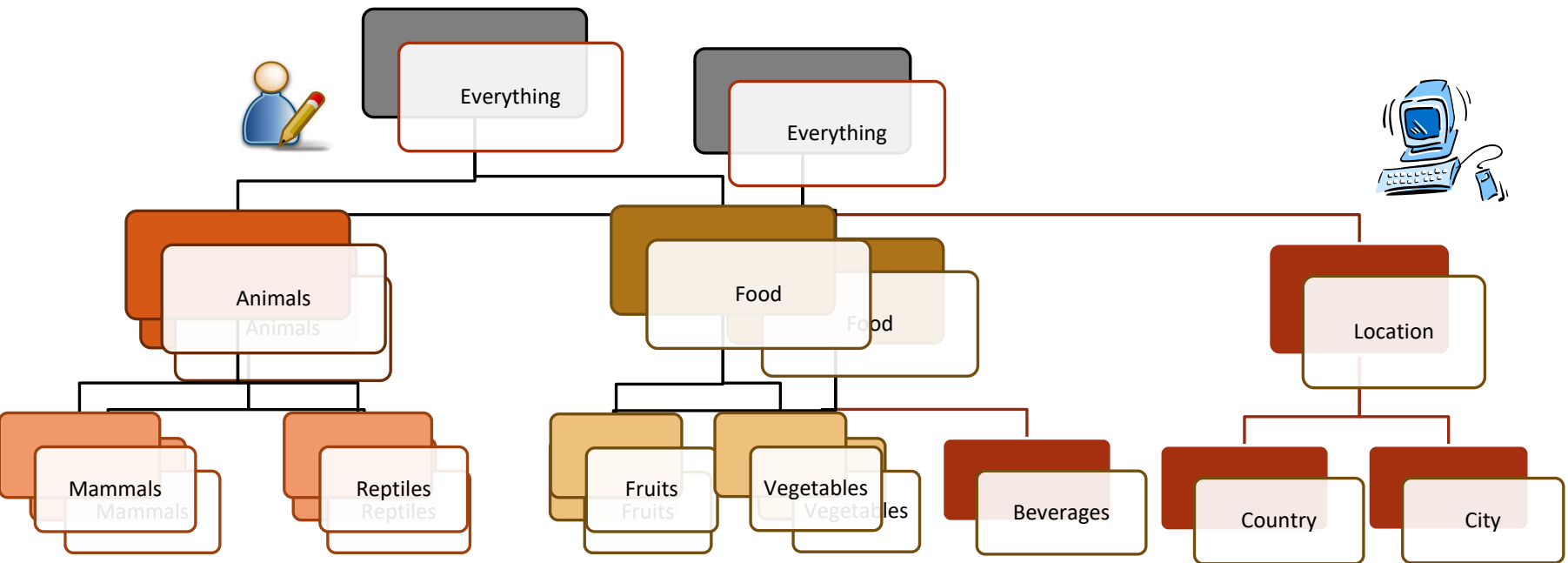
- **Highly semantic ontology**
- **Leads to high precision extractions**
- **Expensive to create**
- **Requires domain experts**



Defining Domain: Semi-automatic



- Subset of types are manually defined
- More types are discovered from data



Defining Domain: Semi-automatic



- Types and type hierarchy is manually defined
E.g. River, City, Food, Chemical, Disease, Bacteria
- Relations are automatically discovered
using clustering methods

Discovered relation	Patterns	Seed instances
River -in heart of- City	"in heart of" "in the center of" "which flows through"	"Seine, Paris", "Nile, Cairo" "Tiber river, Rome" "River arno, Florence"
Food -to produce- Chemical	"to produce" "to make" "to form"	"Salt, Chlorine" "Sugar, Carbon dioxide" "Protein , Serotonin"
Disease -caused by- Bacteria	"caused by" "is the causative agent of" "is the cause of"	"pneumonia, legionella" "mastitis, staphylococcus aureus" "gonorrhea, neisseria gonorrhoeae"

- Easier to derive types using existing resources
- Relations are discovered from the corpus
- Leads to moderate precision extractions
- Partially semantic ontology

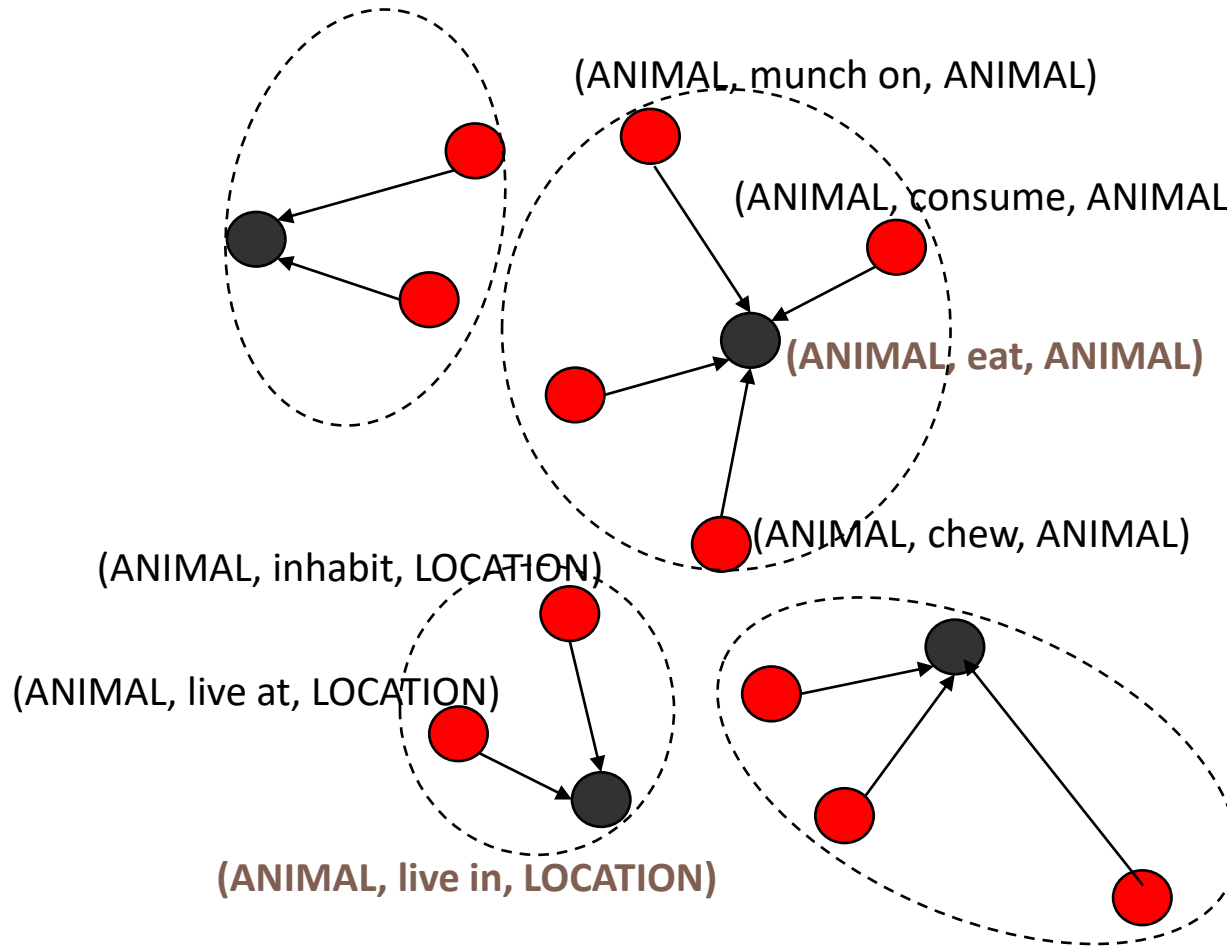
Defining Domain: Automatic



- Any noun phrase is a candidate entity
 - Any verb phrase is a candidate relation
- **Cheapest way to induce types/relations from corpus**
 - **Little/no expert annotations needed**
 - **Limited semantics**
 - **Leads to noisy extractions**

Unsupervised relation induction

(Relation clustering)



Extractors for each relation of interest

Learning Extractors: Manual



- Human defined high-precision extraction patterns for each relation



Person-member of-Band



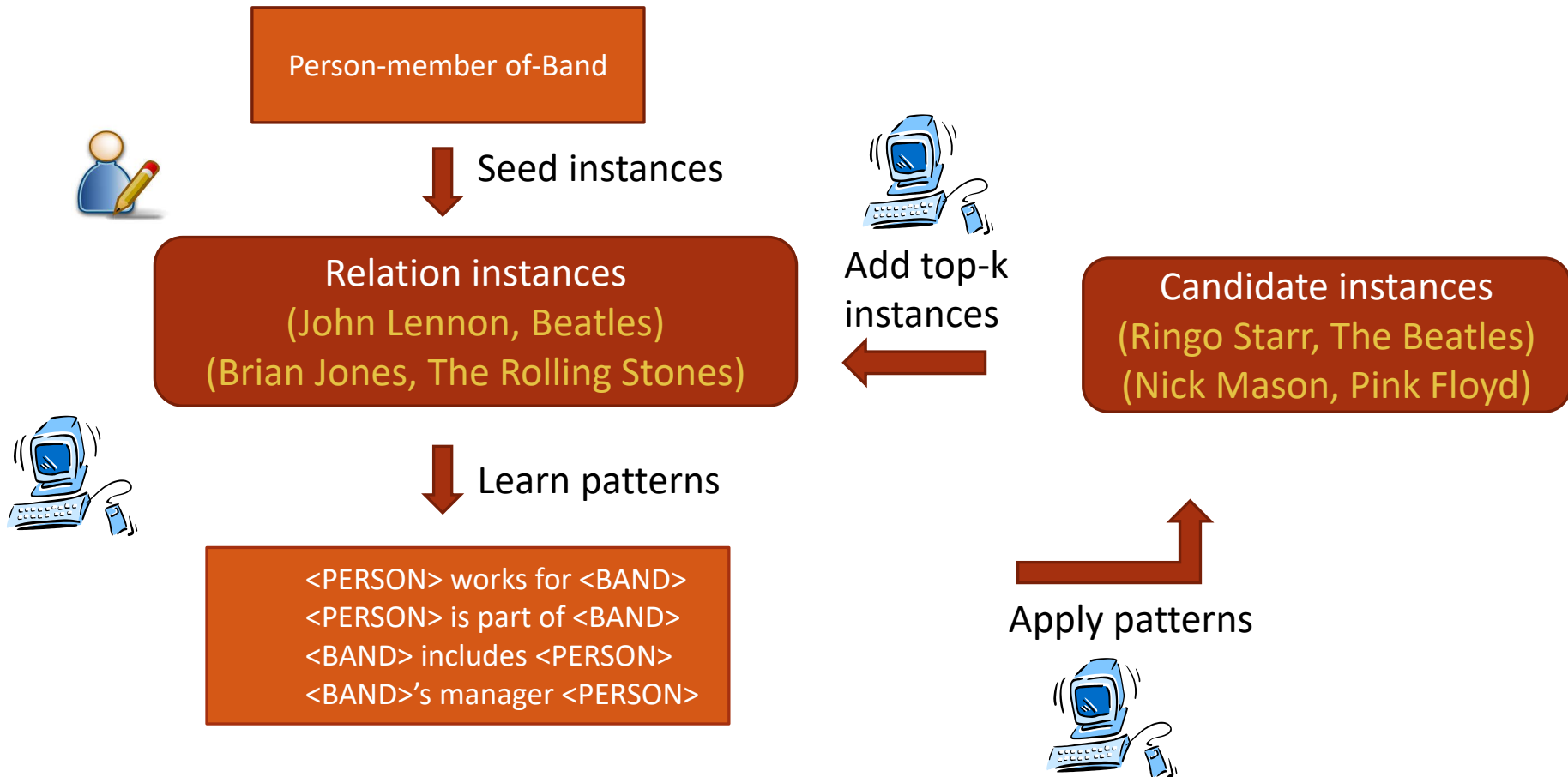
<PERSON> works for <BAND>
<PERSON> is part of <BAND>



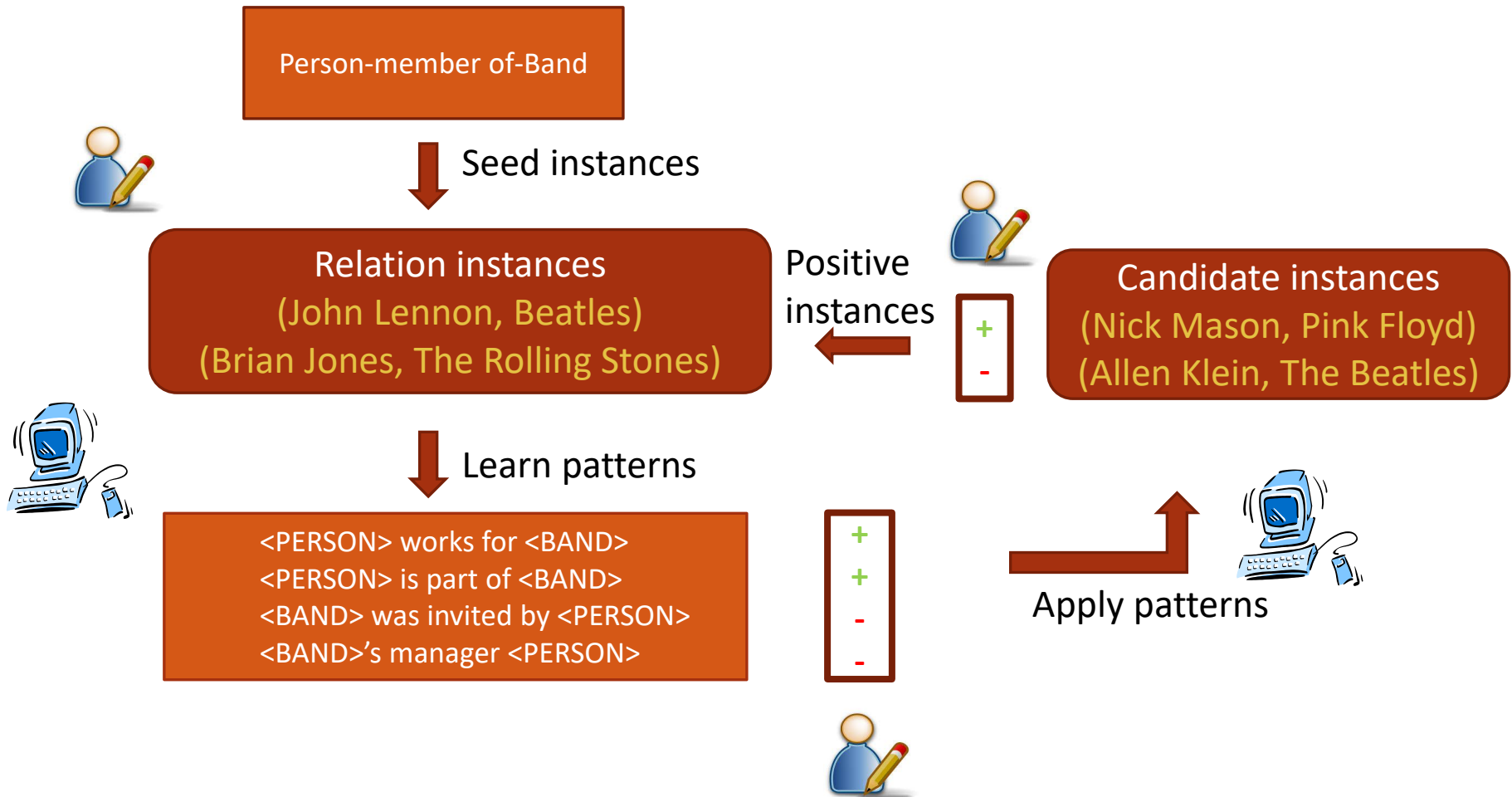
Extract relation instances
(John Lennon, The Beatles)
(Brian Jones, The Rolling Stones)



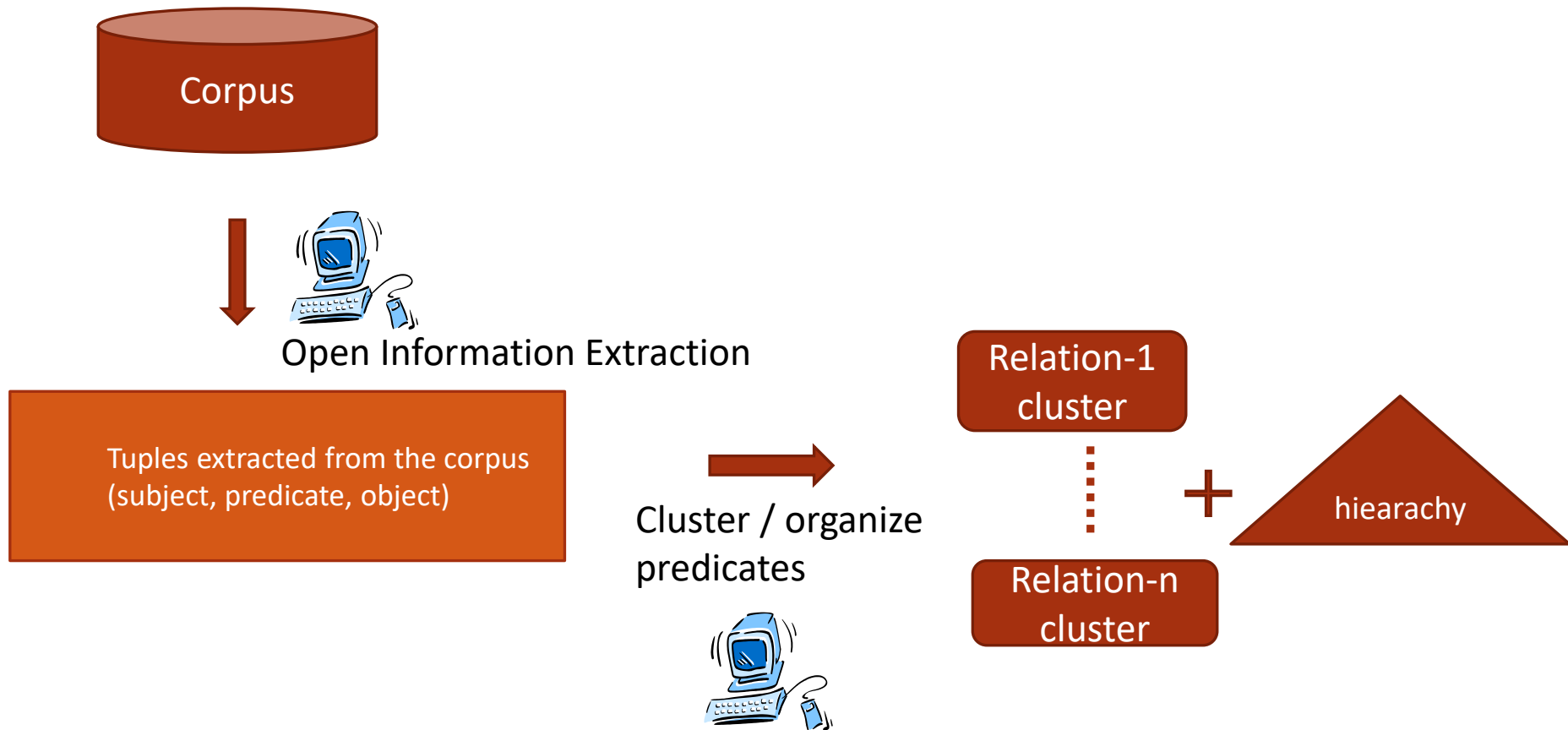
Learning Extractors: Semi-supervised



Learning Extractors : Interactive



Learning Extractors : Unsupervised



Scoring the candidate extractions

Scoring the candidate extractions



- Human defined scoring function
(expensive, high precision, low recall)



- Expert comes up with features
Crowdsourced true/false evaluation of training data
Scoring function is learnt using standard ML



- Completely automatic (Self-training)
Updated set of instances → weights of extraction
patterns → more instances →
(cheap, leads to semantic drift)

Effect of supervision on extractions

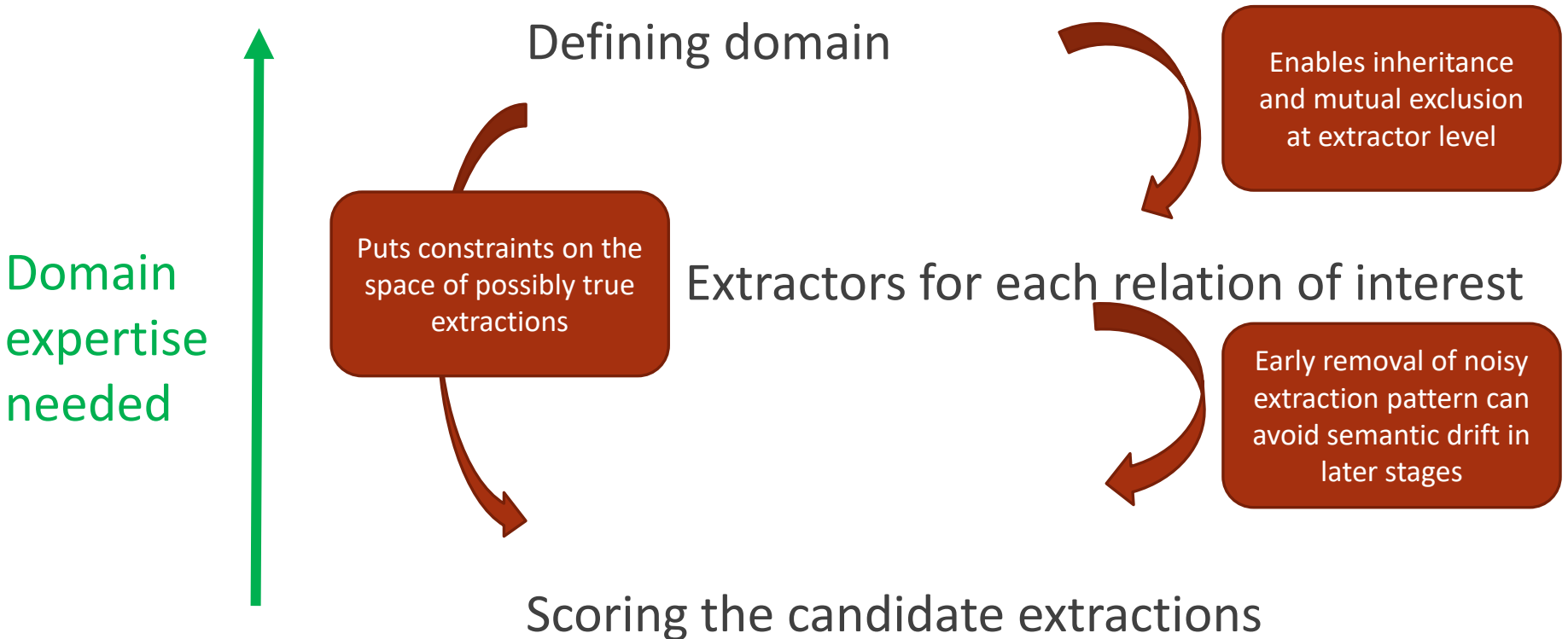
Precision,
Human efforts



Recall,
Speed



Impact of early supervision



Examples: Information Extraction Techniques

(1) Narrow domain patterns

Defining domain		Learning extractors		Scoring extractions	
					

(1) Narrow domain patterns

Use the collection of rules as the system itself

if “X was born in Y” then predict “X birthplace Y”

High precision: when it fires, it's correct

Easy to explain predictions

Easy to fix mistakes



However...

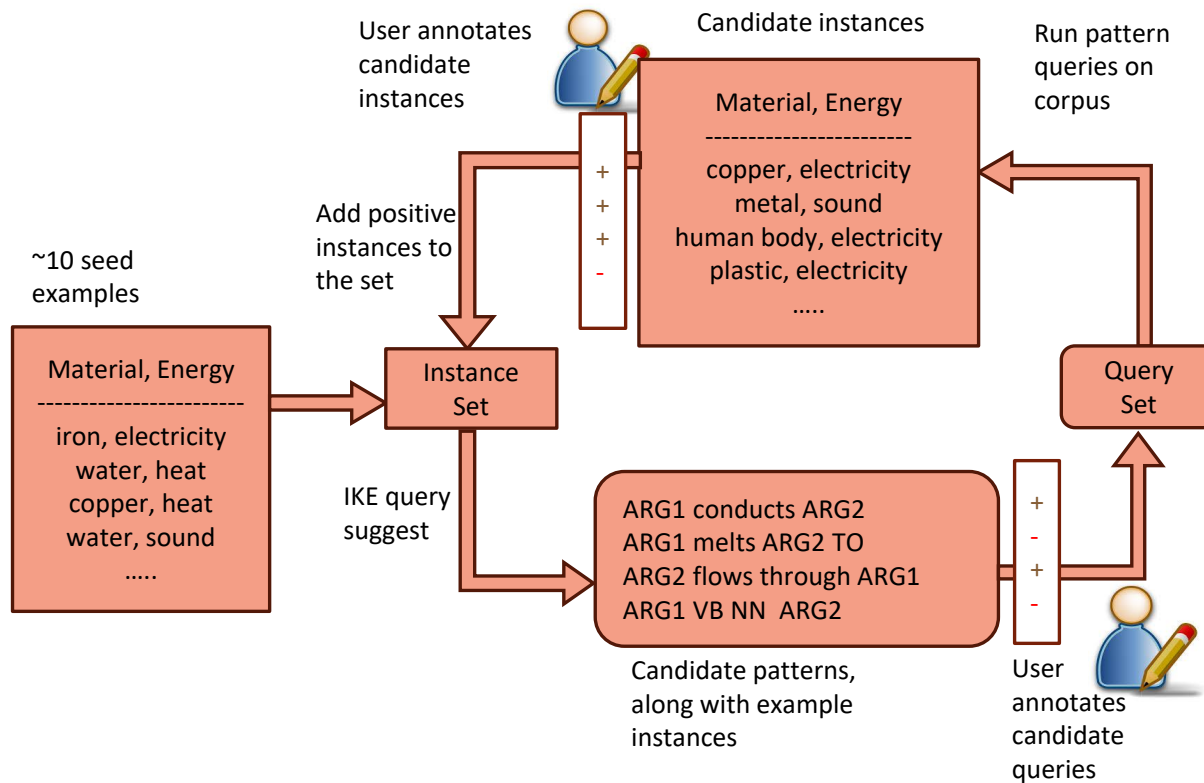
Only work when the rules fire

Do not generalize!

(2) Interactive Bootstrapping (IKE)



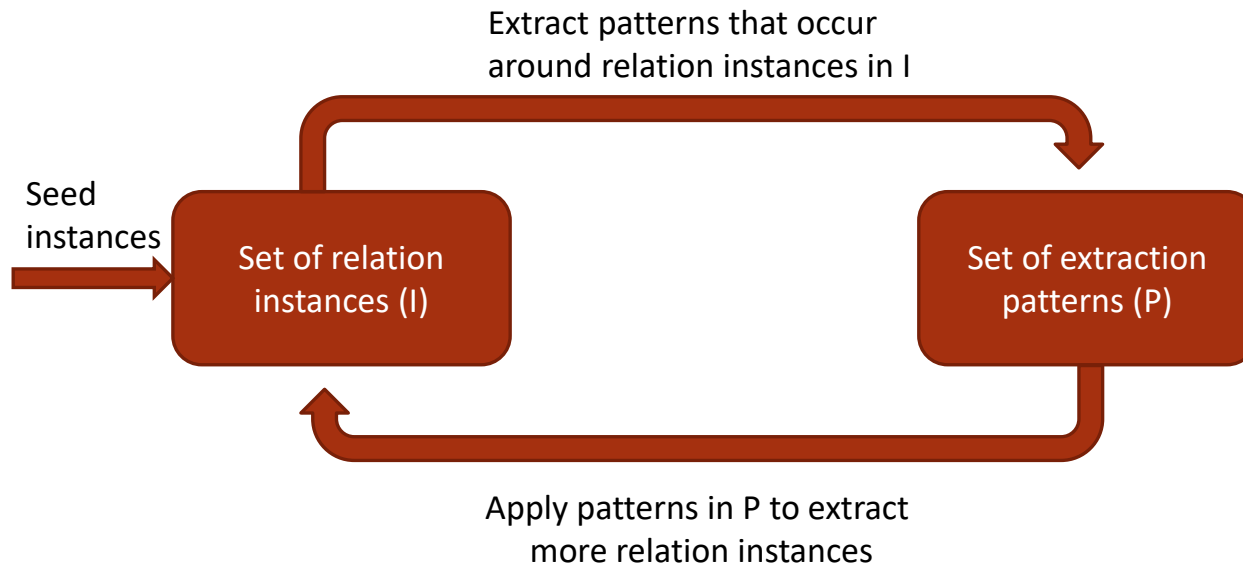
(2) Interactive Bootstrapping (IKE)



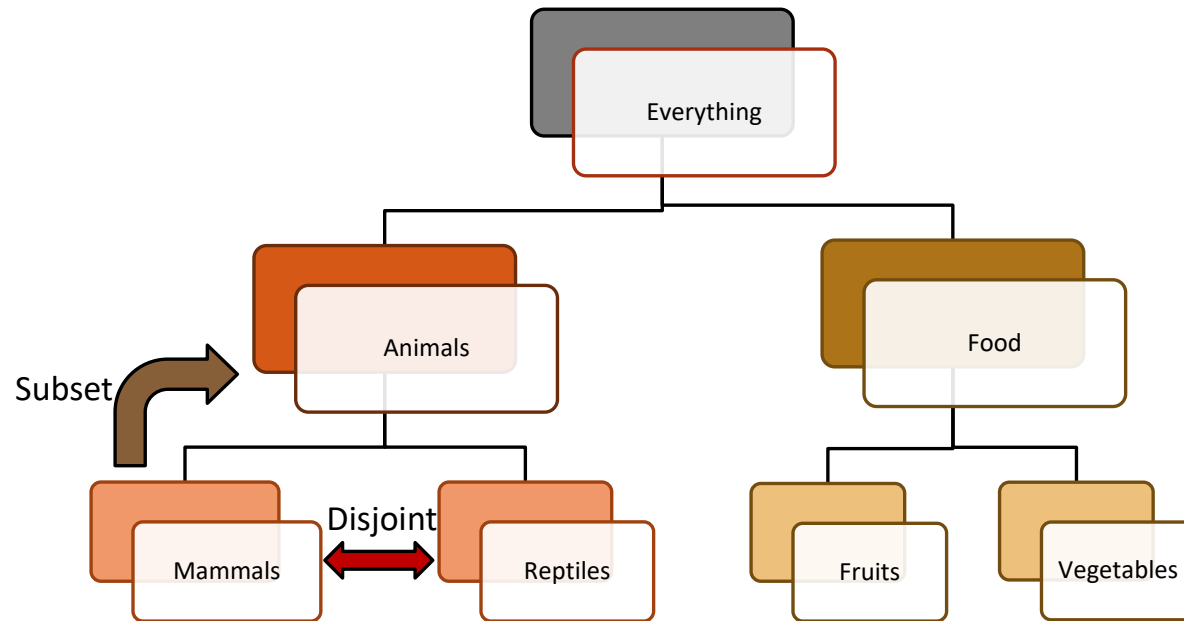
(3) Ontology based extraction



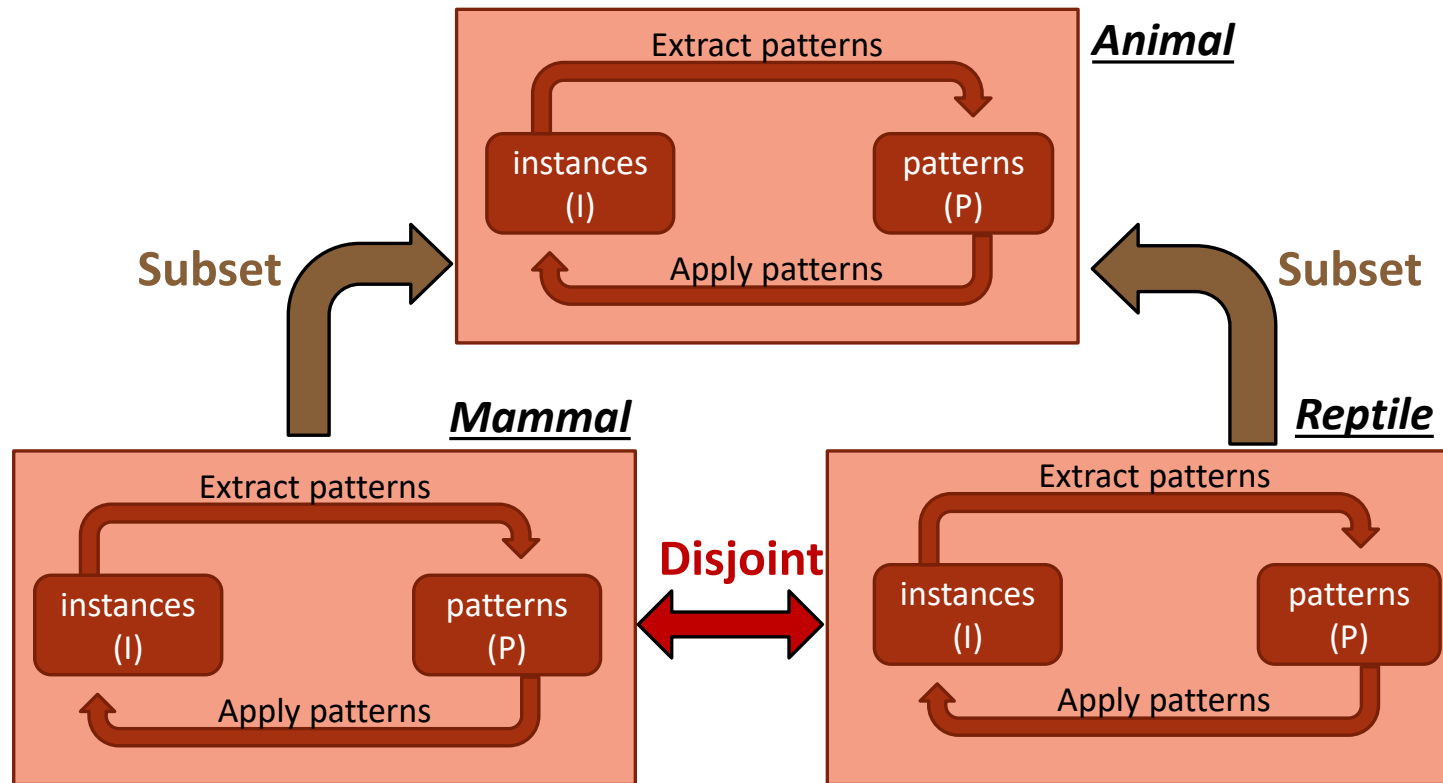
Semi-supervised learning (bootstrapping)



Coupling Constraints (Ontology)



Coupled bootstrap learning



(4) Open Domain IE



(4) Open domain IE

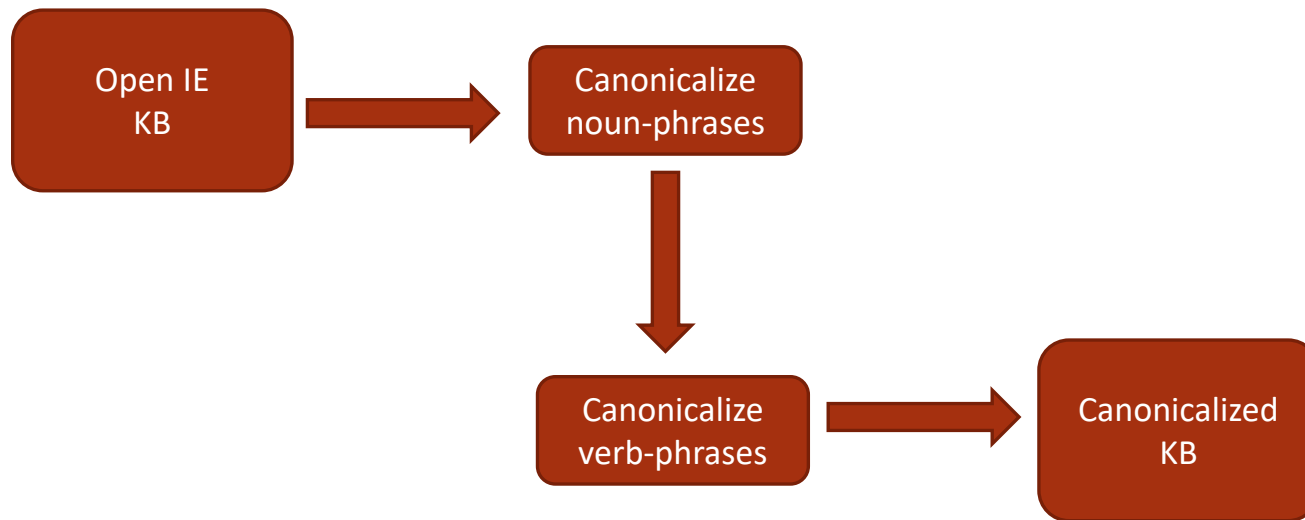
- Any noun phrase is a candidate entity
- Any verb phrase is a candidate relation
- Sentence: “John Lennon was an English music artist who gained worldwide fame as one of the members of the Beatles.”
Open IE extractions:
 - 0.95 (John Lennon; **was**; an English music artist)
 - 0.94 (an English music artist; **gained**; L:worldwide; fame; as one of the members of the Beatles)

(5) Hybrid approach

Adding structure to Open KB



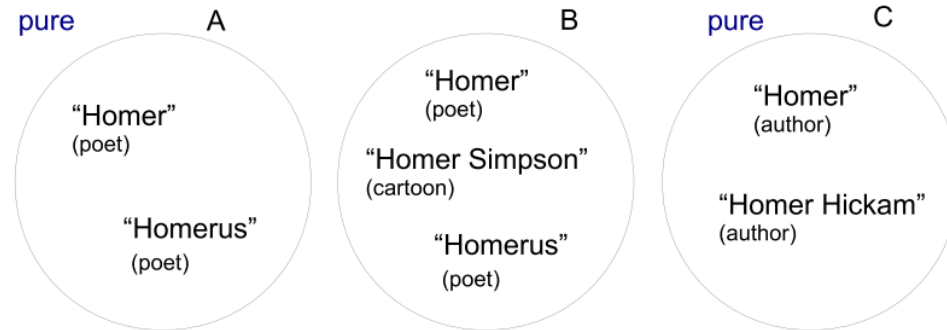
(5) Hybrid approach



(5) Hybrid approach

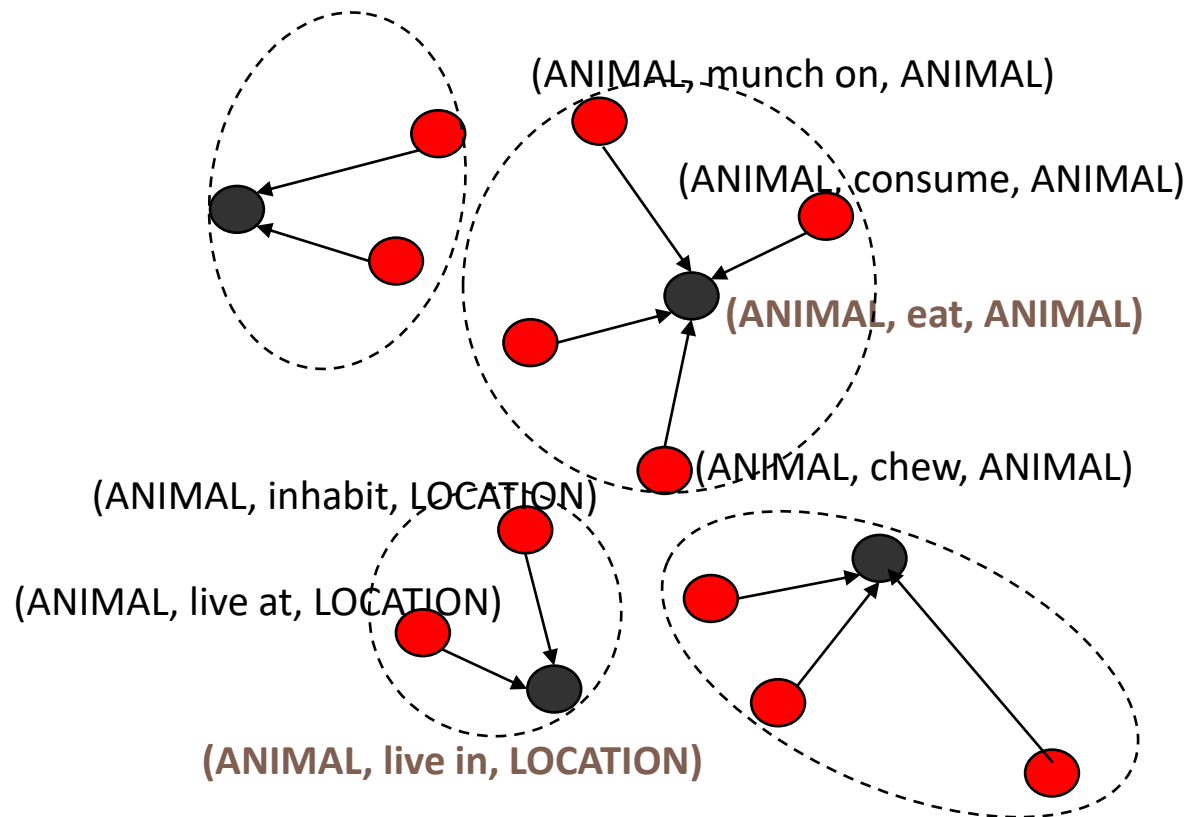
(adding structure to Open KB)

- ***Canonicalizing noun phrases***



Canonical schema induction

Verb phrases	Freebase relation
be an abbreviation-for, be known as, stand for, be an acronym for	-
be spoken in, be the official language of, be the national language of	location.country.official_language
be bought, acquire	organization.organization.acquired_by



Knowledge fusion with multiple extractors

VOTING (AND VS OR OF EXTRACTORS)

CO-TRAINING (MULTIPLE EXTRACTION METHODS)

MULTI-VIEW LEARNING (MULTIPLE DATA SOURCES)

MACHINE LEARNING FOR KNOWLEDGE FUSION

Information Extraction

Single extractor



Fusing multiple extractors

Multiple weak extractors

- **Extractor 1:** text patterns to extract ISA relations
e.g. coupled pattern learner in NELL
- **Extractor 2:** learning wrappers for HTML pages to extract ISA relations
from structured text

Voting Schemes

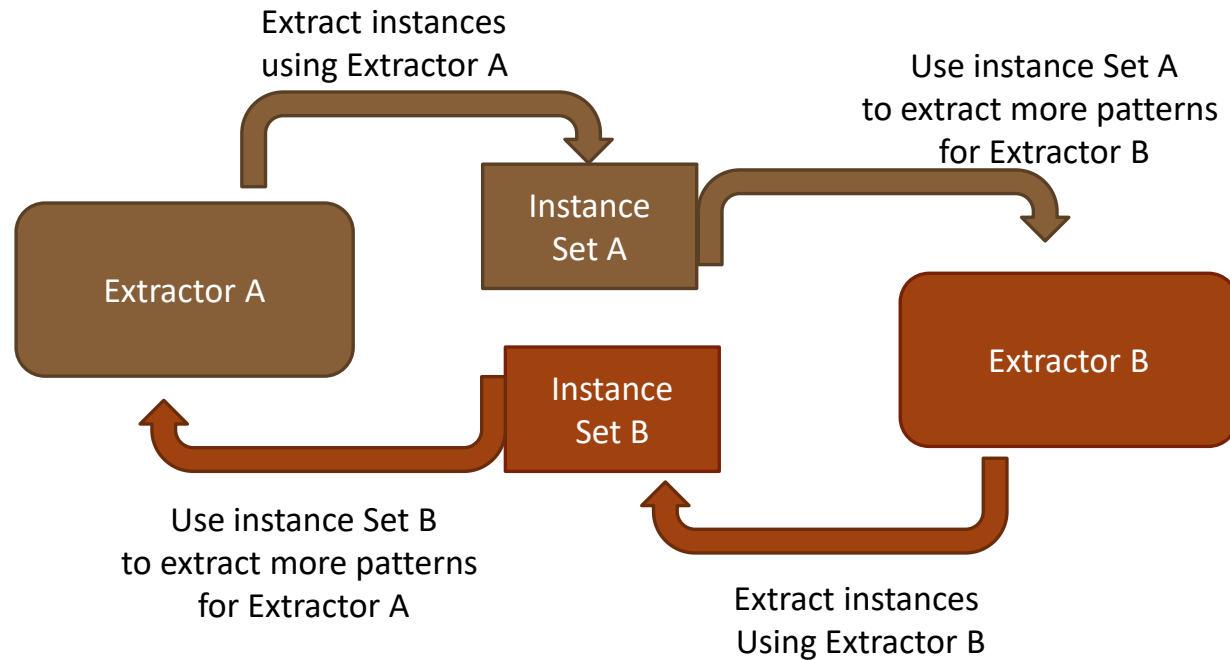
- ***AND of two extractors:***

- For a candidate extraction to be promoted to a fact in KB, both the extractors should support the fact
- $\text{score}(\text{fact}) = \text{score_extractor1}(\text{fact}) * \text{score_extractor2}(\text{fact})$

- ***OR of two extractors***

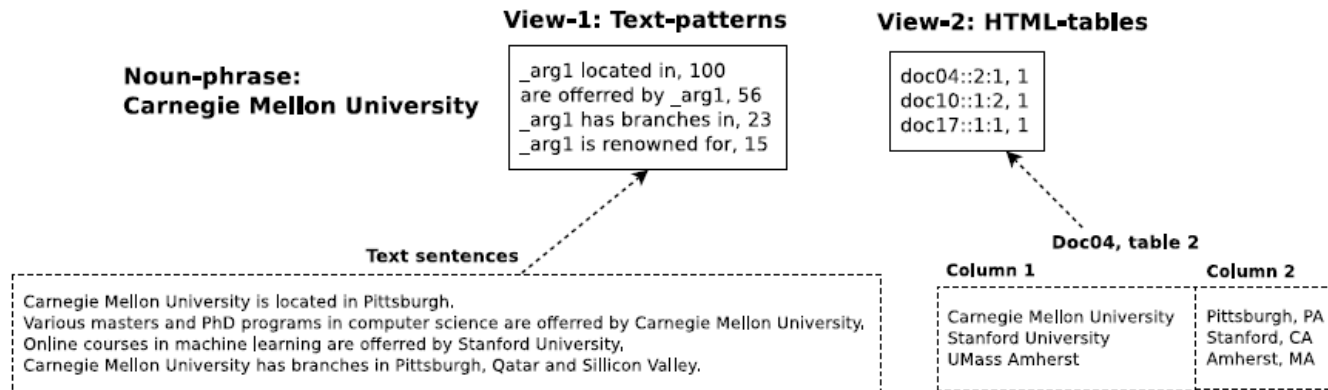
- For a candidate extraction to be promoted to a fact in KB, both the extractors should support the fact
- $\text{score}(\text{fact}) = \text{score_extractor1}(\text{fact}) * \text{score_extractor2}(\text{fact})$

Co-training

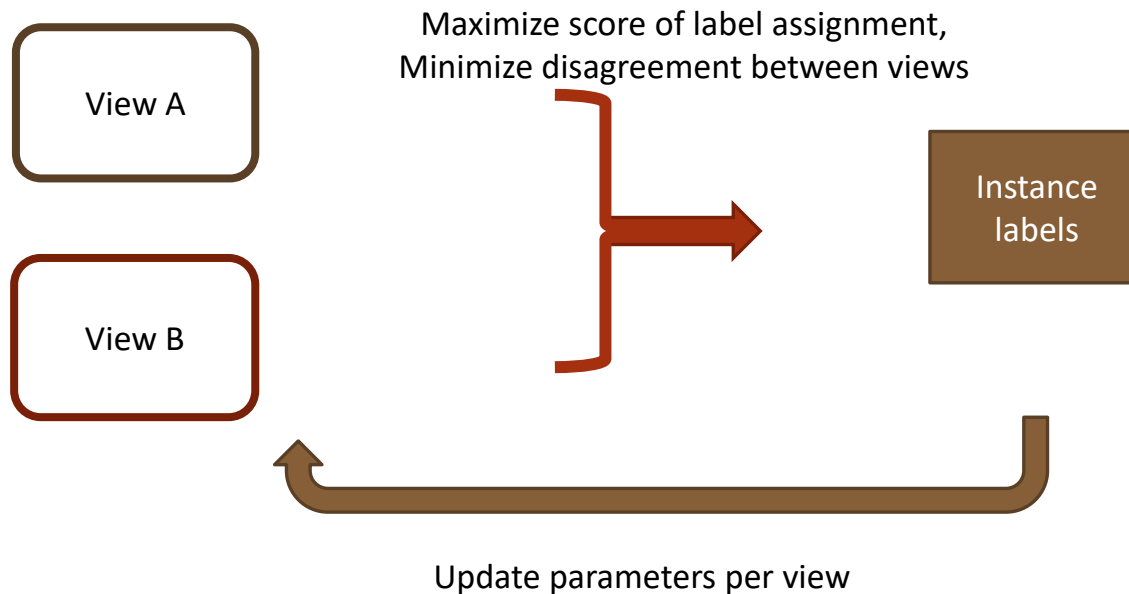


Multiple data-views

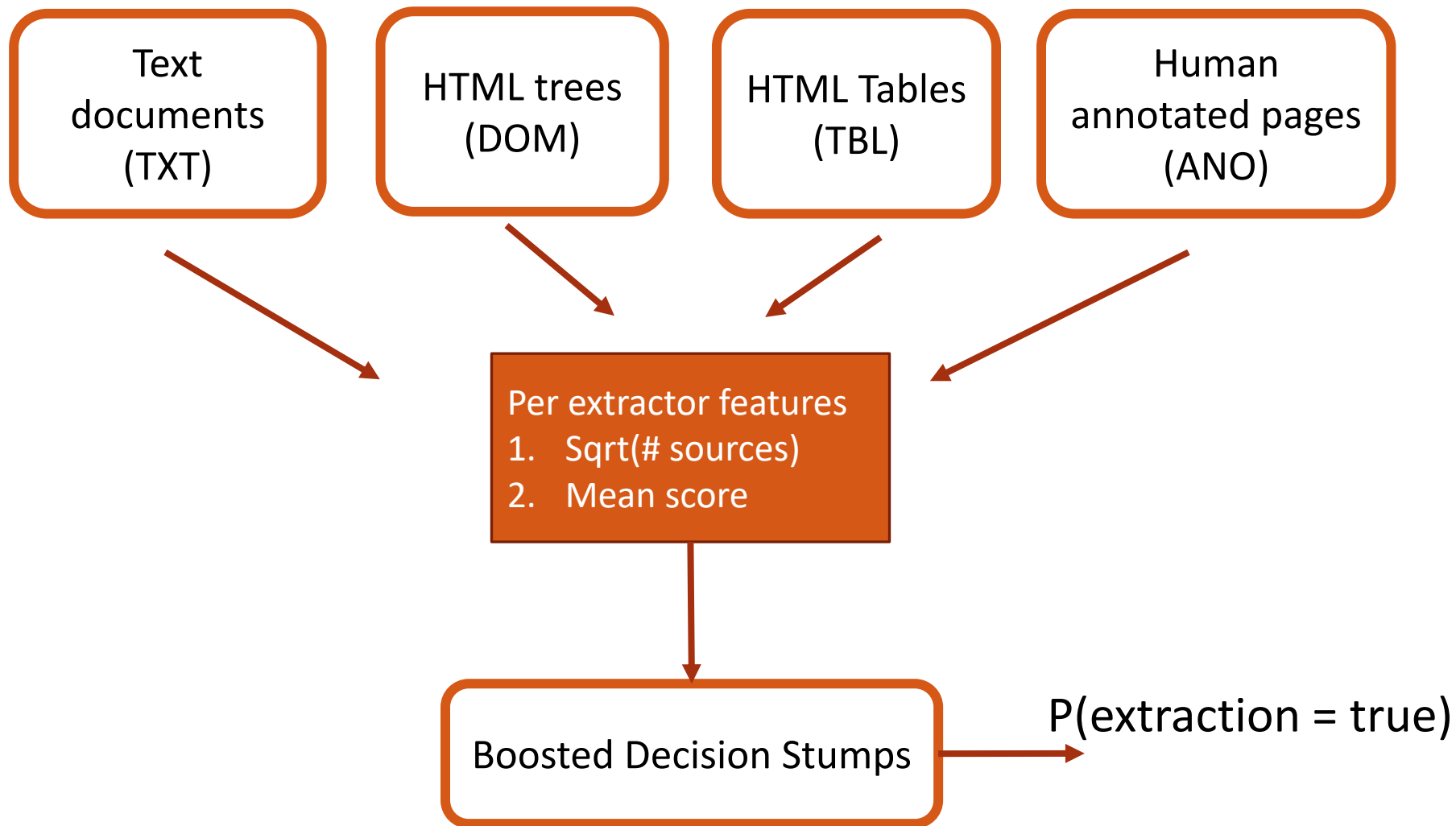
- NP “Carnegie Mellon University” can be represented in two different ways based on its occurrence in text documents and HTML tables.



Multi-view learning



Knowledge vault: fusing the extractors



Example IE Systems

OPEN IE

NELL

KNOWLEDGE VAULT

Open IE (KnowItAll)



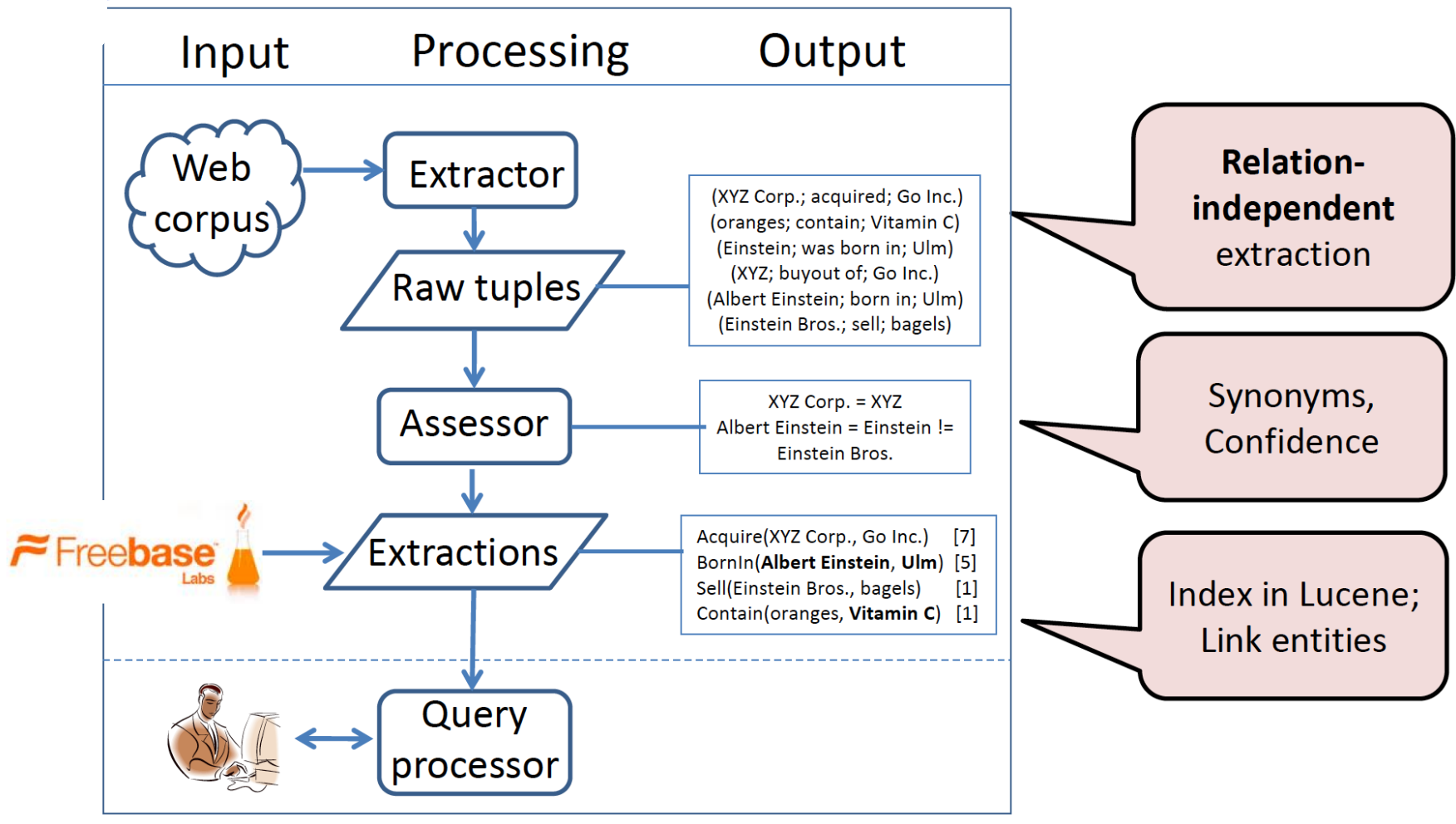
Open Information Extraction



Open IE (KnowItAll)



Open Information Extraction



Never Ending Language Learning (NELL)

Ontology based extraction

NELL Knowledge Base Browser

CMU Read the Web Project

categories

relations

**Defining
domain**



**Learning
extractors**



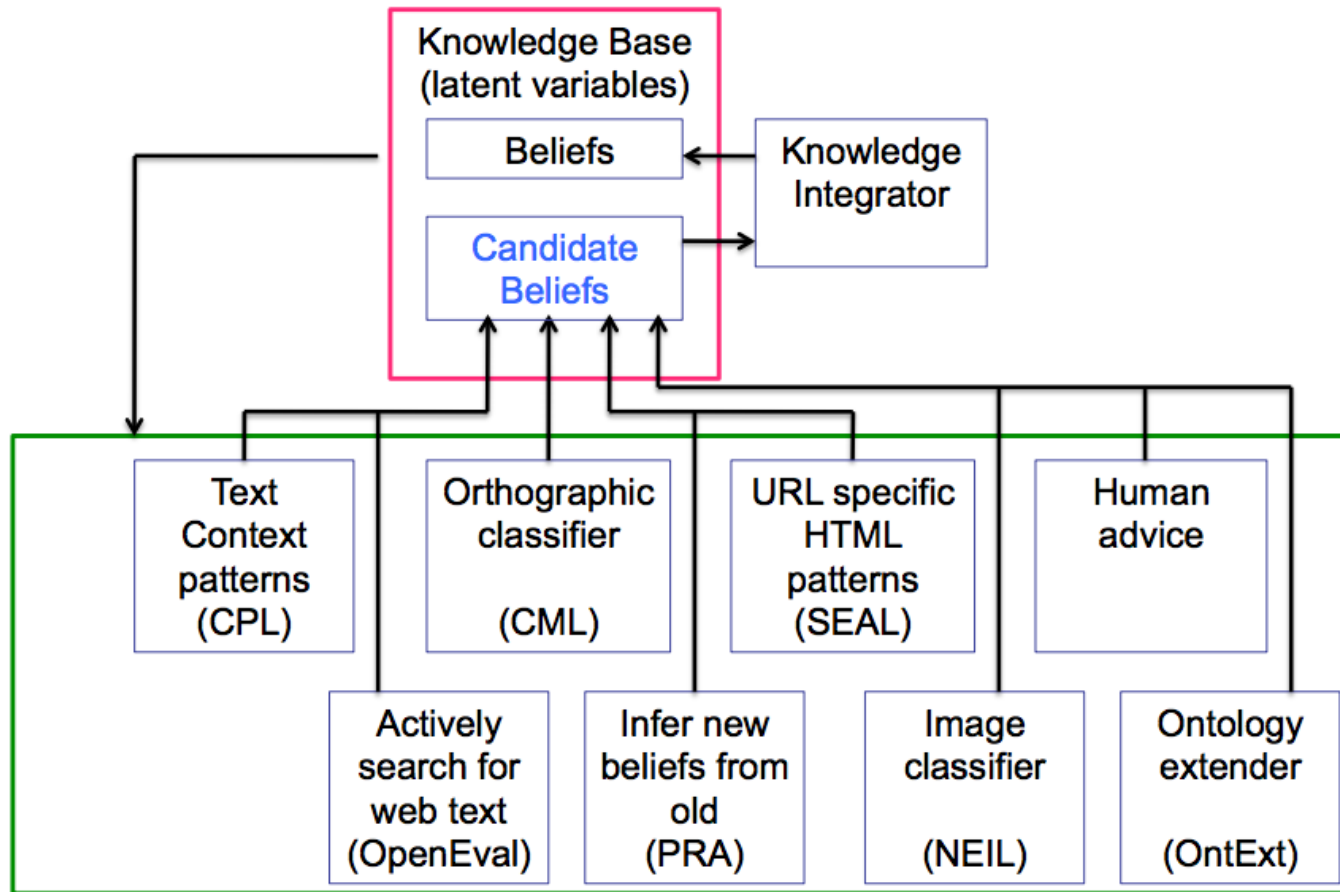
**Scoring
extractions**



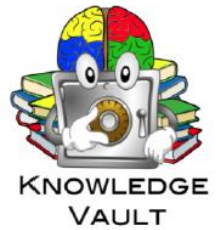
**Fusing
extractors**

Voting +
local constraint
satisfaction

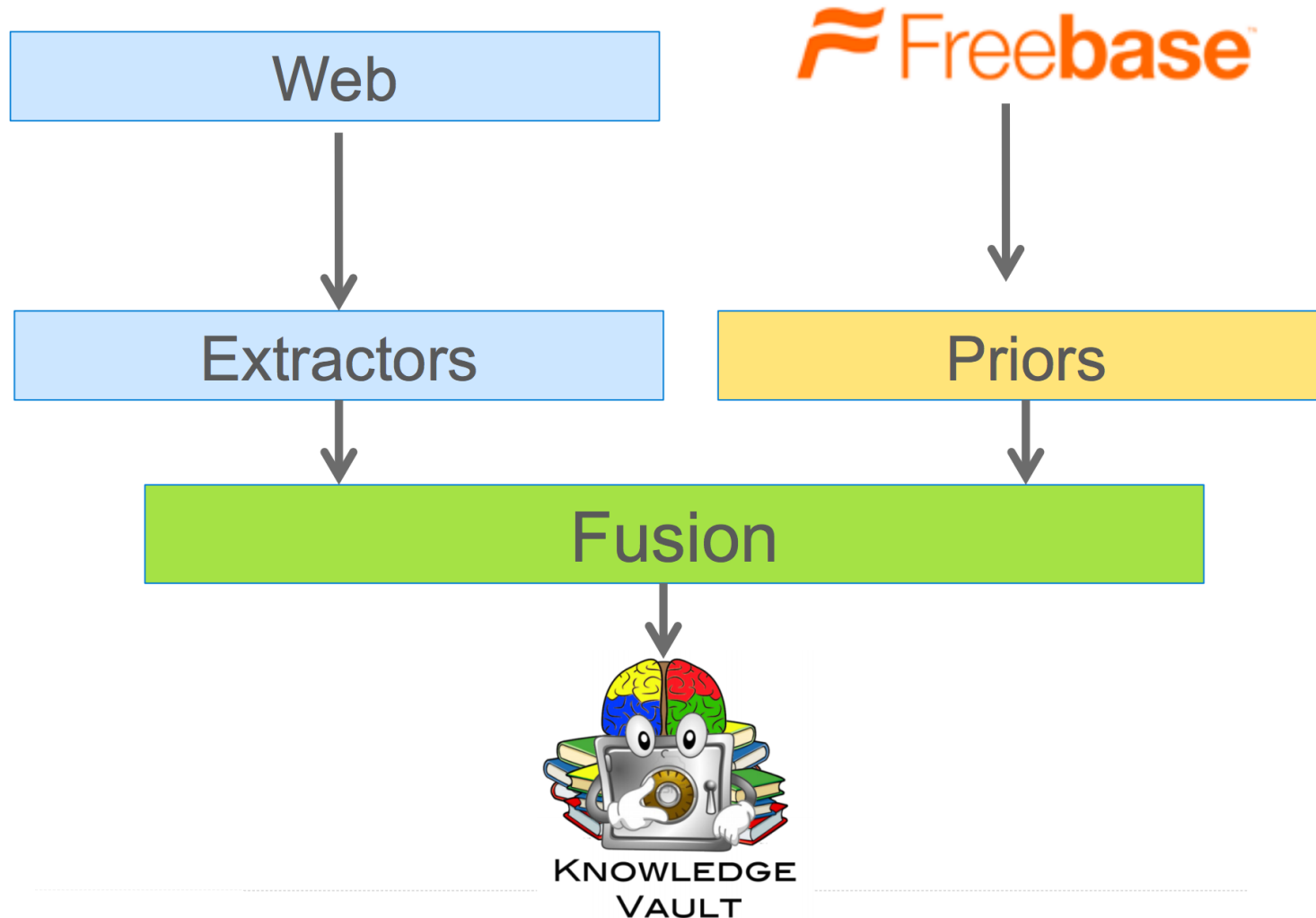
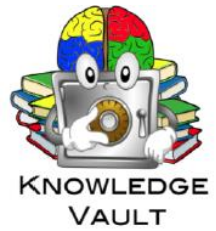
Never Ending Language Learning (NELL)



Knowledge Vault



Knowledge Vault



Summary: Information Extraction

3 IMPORTANT SUB-PROBLEMS

(DEFINE DOMAIN, LEARN EXTRACTORS, SCORE EXTRACTIONS)

3 LEVELS OF SUPERVISION

(MANUAL, SEMI-SUPERVISED, UNSUPERVISED)

KNOWLEDGE FUSION WITH MULTIPLE EXTRACTORS

(CO-TRAINING, MULTI-VIEW LEARNING)

EXAMPLE IE SYSTEMS

Thank You



SEE YOU AFTER THE COFFEE BREAK!

