# Text Mining CBS

Lecture 6: Relation extraction
Piek Vossen

# Literature

- NLTK book, chapter 7, section 6: https://www.nltk.org/book/ch07.html

- Background reading:

  - T. Mitchell et al. (2015) Never-Ending Learning, Association for the Advancement of Artificial Intelligence (www.aaai.org)

  - Kozareva, Z., Riloff, E., & Hovy, E. H. (2008). Semantic Class Learning from the Web with Hyponym Pattern Linkage Graphs

# Overview of relation extraction

- Relation extraction

  - between instances (property-value extraction or knowledge base completion):

    - Microsoft *hasCEO* Bill_Gates, Bill_Gates *isA* CEO

  - between concepts (ontology learning)

    - CEO *subclassOf* Manager, Manager *subclassOf* person

    - Company *make* Product, Pizzeria *subclassOf* Company, Pizza *subclassOf* Product

# Definition of the task of property-value extraction

- A relation is mostly a triple

  - **SUBJECT - PROPERTY/RELATION - VALUE/OBJECT**

  - RDF, but also columns in relational database

  - Barack_Obama **hasAge** 57;

  - Barack_Obama **hasPosition** president

# Definition of the task of property-value extraction

- 3 subtasks

  - identifying the arguments, which can be the subject or object of a relation

  - interpreting the arguments: identities of people and organisations (**subjects**) or values (**objects**) of relations

  - detecting the **property** or **relation**

- **Barack Hussein Obama II is** an American **attorney** and **politician** who **served as** the 44th **president** of the United States from 2009 to 2017. **Obama** was **born** in **Honolulu**, Hawaii. After **graduating** from **Columbia University** in **1983**, he **worked as a community organizer** in Chicago. https://en.wikipedia.org/wiki/Barack_Obama

- **Former President Barack Obama turned 57** on Saturday — and for the first time, his birthday is being celebrated as a commemorative holiday in his former home of Illinois. http://time.com/5358013/barack-obama-birthday-57/

# Property-value extraction builds on top of NERC and NED pipelines

- Knowing the entities you can almost guess the relation

  - Trump - 71, Barack - 57, Merkel - 63, Clinton (which one?)- 70

  - Trump - US, Merkel - Germany, Gates - Microsoft, Messi - Barcelona

- But not always:

  - Al **Gore**, Frank **Gore** (American Football), Lesley **Gore** (singer, songwriter, actress, activist, 68, died in 2015), **Gore** Vidal (writer, 80)

  - Clinton - Clinton

# Different task definitions

- Sentence-level extraction (text annotation, marking text segments where elements are expressed) versus instance-level extraction (triples with URIs)

- Closed property-value extraction:

  - Schema-based relation extraction: subjects, objects and properties/relations are pre-defined in an ontology

  - Slot filling or template-based: no semantic constraints on the type of object/value

- Open information extraction: relations/ontology needs to be discovered

- Knowledge base population/completion: fill gaps in a KB

# Property-value extraction in text, *NLP-ish*

- Text annotation: text segment in which a specific relation is being expressed, e.g. "[John Smith <PER>] was [born <**place-of-birth**>] in [Chicago <LOC>]".

- Slot filling: given a relation schema, extract the slots and relation from text

  - Message Understanding Conferences (MUC): https://cs.nyu.edu/cs/faculty/grishman/muc6.html

  - Knowledge Base Completion (TAC-KBC): https://tac.nist.gov

# Closed Relation Extraction & schema definitions

- **Predefined list of relations**, e.g. <founder-of>, <owner-of>

- number of arguments (usually binary (RDF triples), but also: *buyer-buys-goods-from-seller-for-price*)

- **Restrictions** on the type of slot **fillers**: PERSON, <founder-of>, ORGANISATION

- **Range** of the fillers:

  - only one answer: <birth-place>, <birth-date>

  - two answers: <child-of>

  - open list of answers: <married-to>, <founder-of>

# Schemas

- YAGO:

  - Knowledge Base derived from Wikipedia, WordNet and GeoNames: 10 million entities (like persons, organizations, cities, etc.) and contains more than 120 million facts about these entities.

  - https://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/

- http://schema.org: vocabularies to capture entities, relationships between entities and actions: http://schema.org/Movie

  - <div **itemscope itemtype** ="**http://schema.org/Movie**">

  - <h1 **itemprop**="**name**">Avatar</h1>

  - <span **itemprop**="**genre**">Science fiction</span>

  - </div>

# Advantages of using a schema: quality checks (!)

- **Disjunct properties**: *Ford* cannot be a person, company and a car at the same time while only one has a birth-date

- **Property dependencies**: if you have a <place-of-birth> you also have a <birth-data>

- **Constraints** on fillers to semantics types: certain relations are limited to politicians, football players, movie directors, etc.

- **Generalise** data from instances: if *some* movies have directors then also *other* movies are likely to have directors.

# Methods

- Rules and patterns, Supervised, Unsupervised (clustering)

- Bootstrap: start with seed relations to learn patterns and obtain new seeds, etc…

- Distant supervision: get seed data from a knowledge base

- Closed Relation Extraction: reason or filter the data using a schema

# Using Regular Expressions

Entities:

**NERC —>** PER: Jose Mourinho
POSITION: trainer
**NERC —>** ORG: Chelsea

Relation

Jose Mourinho
↓
Trainer
↓
Chelsea

FEBRUARY 27 2007 14:34h

## Mourinho Verbally Attacks Again

Text

f Me gusta    0 hare f    g+1  0

Hvala vam na povratnim informacijama! Nazad Pregledat ćemo ovaj oglas kako bismo poboljšali doživljaj u budućnosti.

The target like many times before was Arsenal and Arsene Wenger, and he did not miss Manchester United and Christian Ronaldo either.

Google

Jose Mourinho, trainer of Chelsea has captured the English media again with his well known "blunt" language and "moderate" sarcasm.

- Appelt, D.E., Hobbs, J.E., Bear, J., Israel, D. & Tyson, M. (1993), FASTUS: A finite-state processor for information extraction from real world text. IJCAI, pp.1172-1178.

# Never Ending Language Learning (NELL), Tom Mitchell, Carnegie Mellon University

- Semi-supervised learning method runs 24/7 to train hundreds of different extraction methods for a wide range of categories and relations.

- Accumulated a knowledge base of 3,109,311 asserted instances of 1,186 different categories and relations, since 2010

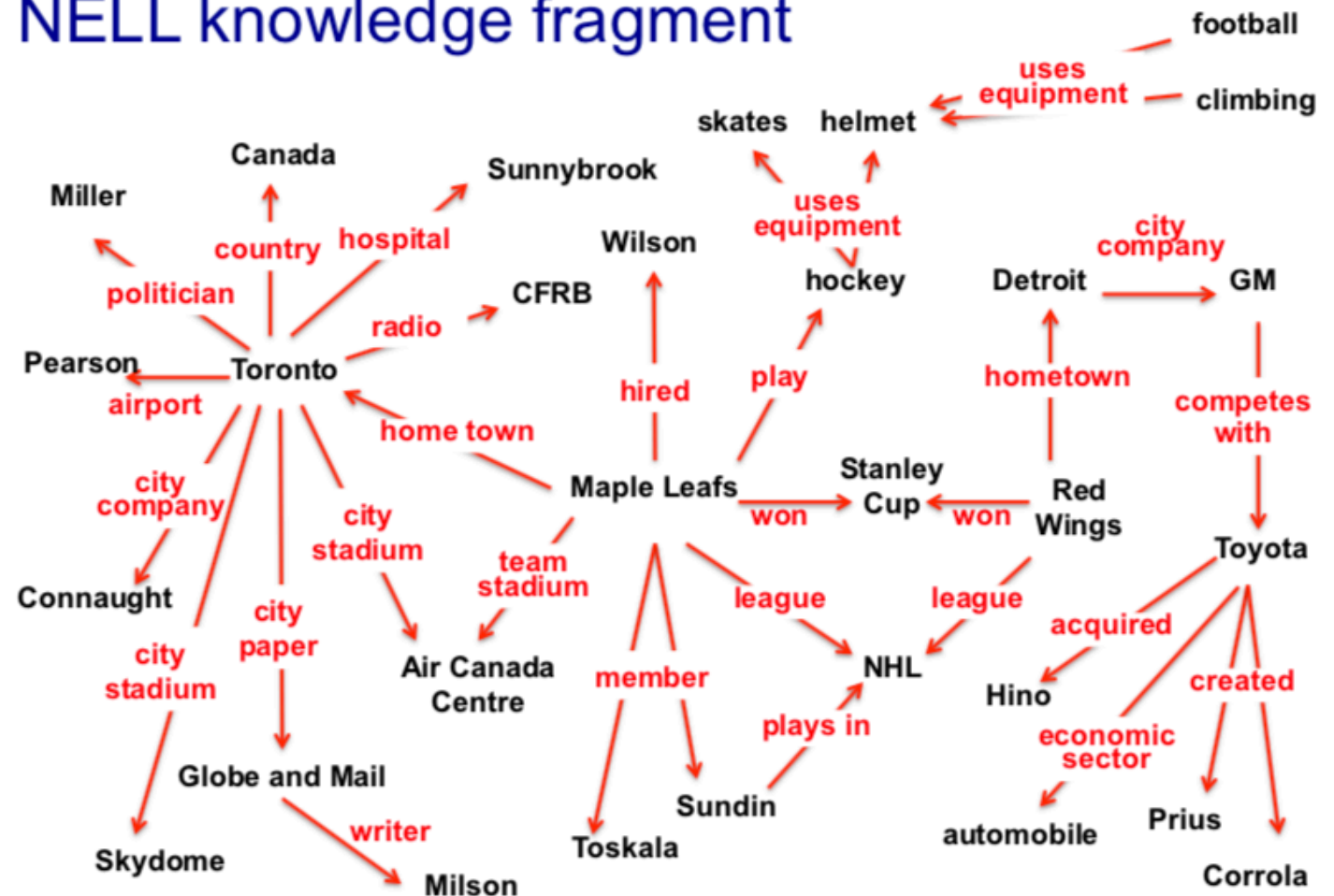- http://rtw.ml.cmu.edu/rtw/

- http://rtw.ml.cmu.edu/rtw//json0doc



Figure 1: **Fragment of the 80 million beliefs NELL has read from the web.** Each edge represents a belief triple (e.g., play(MapleLeafs, hockey), with an associated confidence and provenance not shown here. This figure contains only correct beliefs from NELL's KB – it has many incorrect beliefs as well since NELL is still learning.

# Never Ending Language Learning (NELL), Tom Mitchell, Carnegie Mellon University

- Initial ontology with hundreds of categories (e.g., person, sportsTeam, fruit, emotion) and relations (e.g., playsOnTeam (athlete ,sportsTeam), playsInstrument (musician, instrument)) and 10 to 15 seed examples.

- Train relation detectors on text and tables for the given relations using variety of features and exploiting redundancy on the web

- Index 500 million web pages and access to the remainder of the web through search engine APIs,

# Concept Drift

- Precision declines over time

- Drift of bootstrapping methods, e.g. "work" can mean many different things (function, have an effect, provoke) and noise will change the meaning of "work" based on the data.

- Human intervention (active learning): NELL every few weeks, 5 minutes fix blatant errors



NELL KB assertions vs. time

# Supervised relation extraction pipeline

**Extraction Template**

| Relation | Subject | Object | List? |
|---|---|---|---|
| founder-of | PER | ORG | yes |
| birthdate | PER | DATE | no |

**features**

| Tokens | POS | NE | ... |
|---|---|---|---|
| Bill | NNP | PER | |
| Gates | NNP | PER | |
| founded | VBD | O | |
| Microsoft | NNP | ORG | |
| . | . | O | |

Sentence Splitting, Tokenisation, POS tagging, NERC, TimeEx Normalisation, Parsing

**Annotated Corpus**

**Training Instances**

| Relation | Sentence |
|---|---|
| founder-of | <s>Bill Gates</s> founded <o>Microsoft</o>. |
| birthdate | <s>John Smith</s> was born in Chicago on <o>New Years Day 1985</o>. |

**Testing Instances**

| founder-of | <o>Google</o> was founded by <s>Sergey Brin</s>. |
|---|---|
| birthdate | <s>Amanda Smith</s>'s birthdate is <o>12th March 1956</o>. |

Develop (on train set)

**Relation Extractor** — Predict (on test set) →

**Results**

| Relation | Subject | Object |
|---|---|---|
| founder-of | Sergey Brin | Google |
| birthdate | Amanda Smith | 1956-03-12 |

**features**

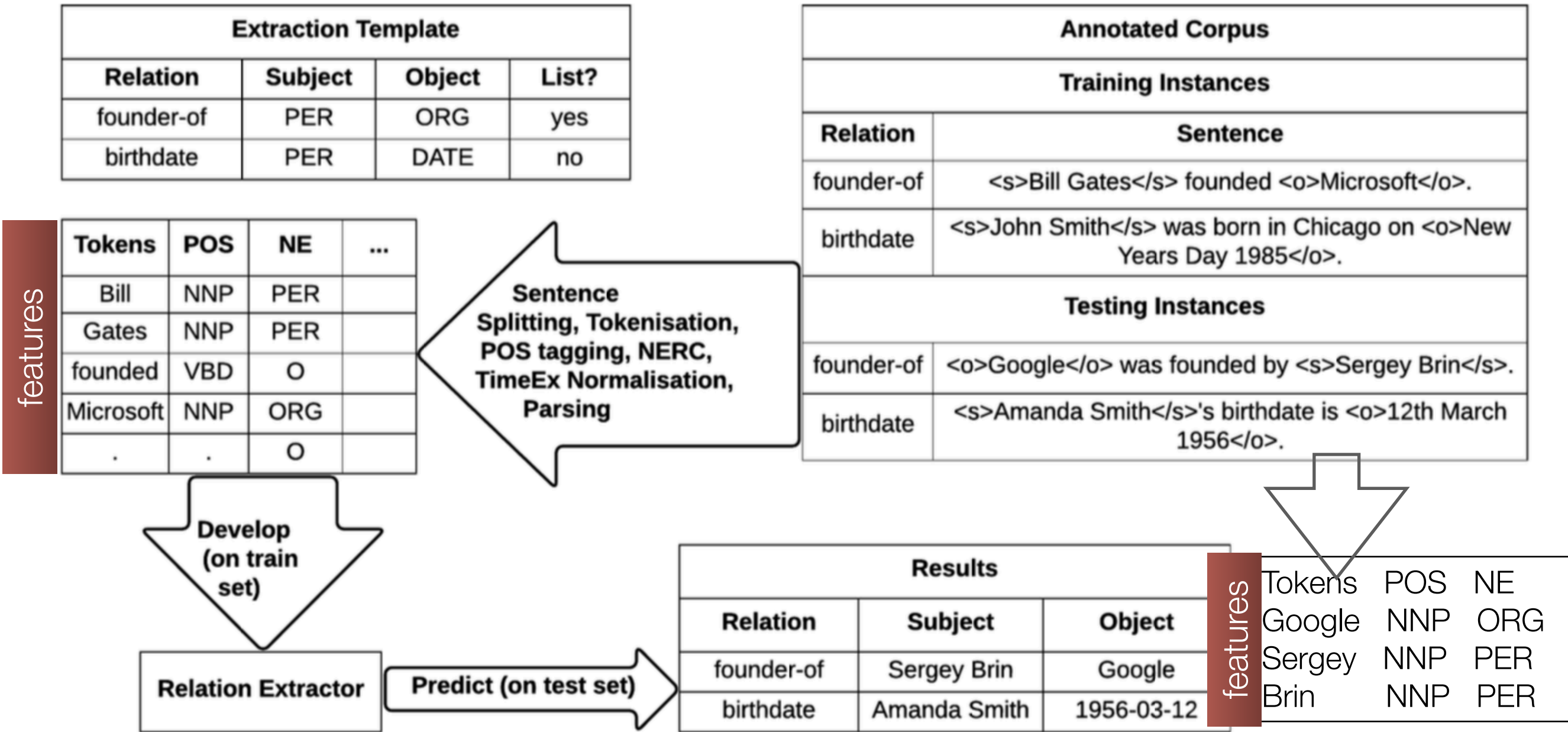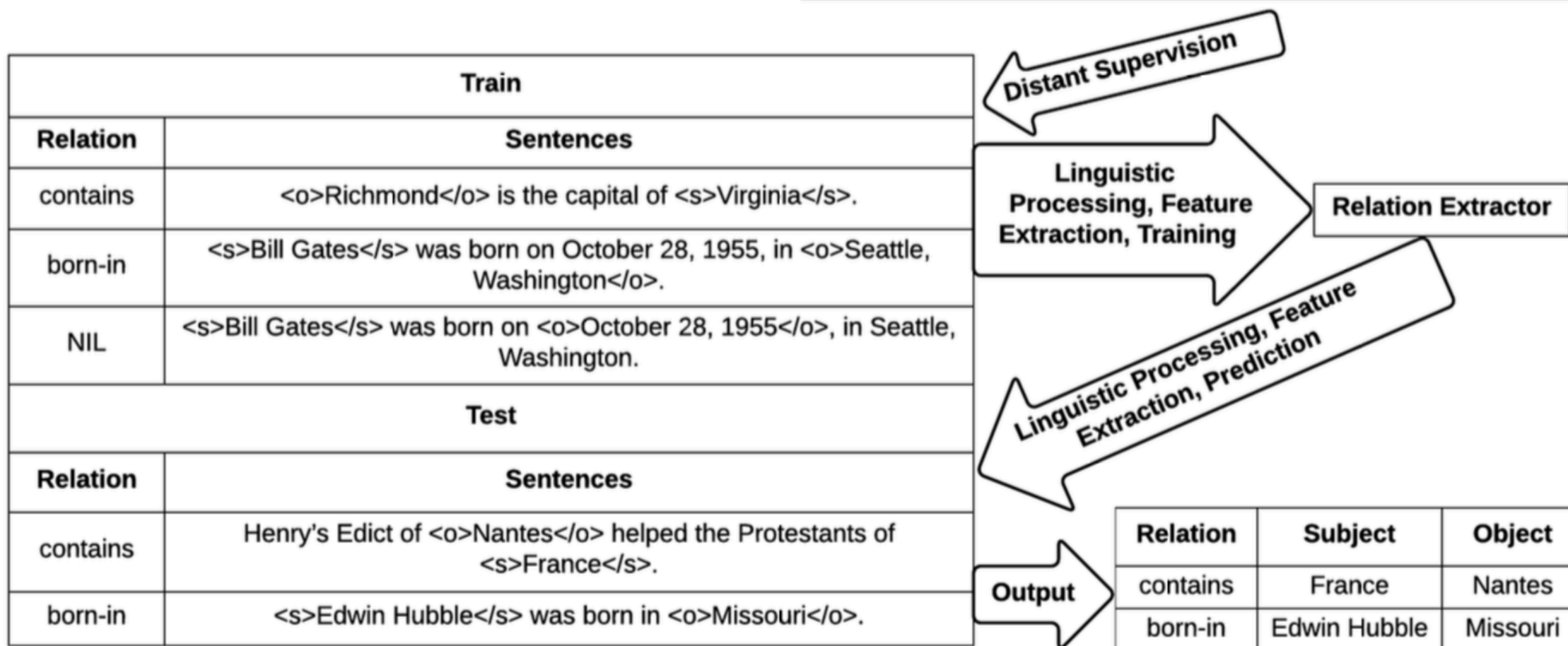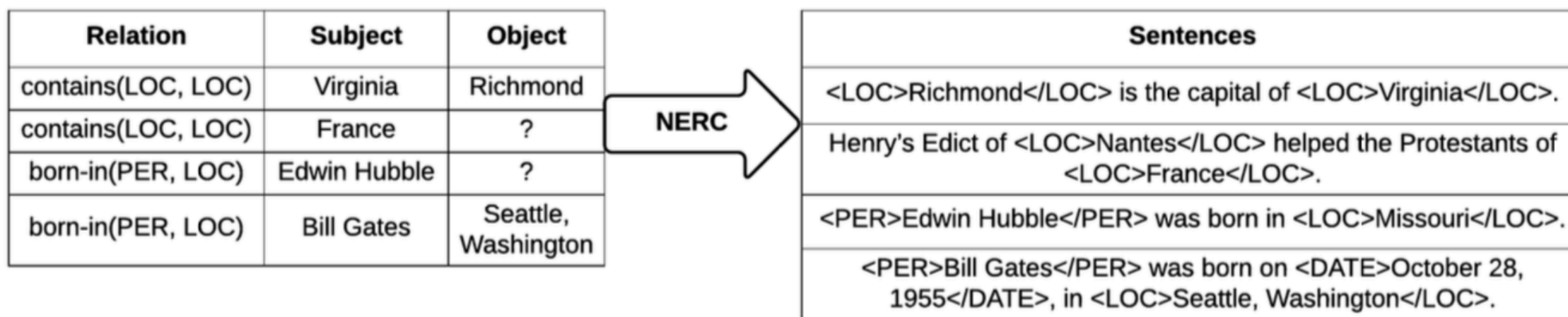| Tokens | POS | NE |
|---|---|---|
| Google | NNP | ORG |
| Sergey | NNP | PER |
| Brin | NNP | PER |

**Figure 4.1:** Typical Relation Extraction Pipeline

*Feature sets need to match across training & test*

# Distant supervision

If two entities participate in a relation, any sentence that contains those two entities might express that relation.

| Relation | Subject | Object |
|---|---|---|
| contains(LOC, LOC) | Virginia | Richmond |
| contains(LOC, LOC) | France | ? |
| born-in(PER, LOC) | Edwin Hubble | ? |
| born-in(PER, LOC) | Bill Gates | Seattle, Washington |

**NERC** →

| Sentences |
|---|
| \<LOC\>Richmond\</LOC\> is the capital of \<LOC\>Virginia\</LOC\>. |
| Henry's Edict of \<LOC\>Nantes\</LOC\> helped the Protestants of \<LOC\>France\</LOC\>. |
| \<PER\>Edwin Hubble\</PER\> was born in \<LOC\>Missouri\</LOC\>. |
| \<PER\>Bill Gates\</PER\> was born on \<DATE\>October 28, 1955\</DATE\>, in \<LOC\>Seattle, Washington\</LOC\>. |

Training data for free

**Train**

| Relation | Sentences |
|---|---|
| contains | \<o\>Richmond\</o\> is the capital of \<s\>Virginia\</s\>. |
| born-in | \<s\>Bill Gates\</s\> was born on October 28, 1955, in \<o\>Seattle, Washington\</o\>. |
| NIL | \<s\>Bill Gates\</s\> was born on \<o\>October 28, 1955\</o\>, in Seattle, Washington. |

**Test**

| Relation | Sentences |
|---|---|
| contains | Henry's Edict of \<o\>Nantes\</o\> helped the Protestants of \<s\>France\</s\>. |
| born-in | \<s\>Edwin Hubble\</s\> was born in \<o\>Missouri\</o\>. |

Distant Supervision

Linguistic Processing, Feature Extraction, Training → **Relation Extractor**

Linguistic Processing, Feature Extraction, Prediction

**Output** →

| Relation | Subject | Object |
|---|---|---|
| contains | France | Nantes |
| born-in | Edwin Hubble | Missouri |

**Figure 4.2:** [81] Distant Supervision Method Overview

# Open Information Extraction

- Relations are not predefined but to be discovered from large collections of text:
  M1<ORG> <***relatedTo***> M2<PER>

- Double processing, compare Double Propagation for sentiment and targets:

  - Process all data to annotate POS and NP Chunks

  - Supervised classifier trained on small subset with relevant/not-relevant or positive/negative cases

- Unsupervised clustering of sentences in which similar entities (embeddings) co-occur:

  - [Trump, Obama, Bush, Clinton] <***relatedTo***> [White House, government]

- Map detected clusters of relations to schema a posteriori —> <***worksAt***>

- Incoherent (contains X and omits Y) and uninformative relations (X has colour Y): syntactic constraints, sufficient variety of arguments, etc..

# Overview of relation extraction methods

| Method | Input | Output | Description | Advantages | Disadvantages |
|---|---|---|---|---|---|
| Boot-strapping | Unlabelled text, relation schema, rules and/or examples | Extraction rules, relations | Using a small set of extraction rules, extract examples, keep prominent ones, iteratively learn more extraction rule and examples | Easy to add new rules, can also be supplied by user | Often low recall and/or manual refinement needed for high precision |
| Rule-based | Unlabelled text, relation schema, rules and NE gazetteers | Relations | Using extraction rules and gazetteers, extract relations | Easy to add new rules, can also be supplied by user | Often low recall, much manual effort to develop |
| Supervised | Labelled text, relation schema | Relations | Using a schema and labelled training data, train model | Currently highest precision and recall for schema-specific relation extraction | Up-front effort of labelling data, risk of overfitting training set |
| Open IE | Unlabelled text | Groups of relations | Discover groups of relations from text using clustering, keep prominent ones | No knowledge about text needed | Difficult to make sense of groups and map to relation schemas |
| Distantly Supervised | Unlabelled text, relation schema, examples | Extraction model, relations | Using a schema and examples of relations, automatically annotate training data, train a model to extract more relations | Extracting relations with high recall and precision | Initial examples required |
| Universal Schema | Several partly populated knowledge bases | Unified knowledge | Take several KBs defined by different schemas, partly populated with relations, predict union of KBs | Integrate relations defined by different schemas after extraction | For small KBs it can be faster to do this manually |

**Table 4.2:** Comparison of different minimally supervised relation extraction methods
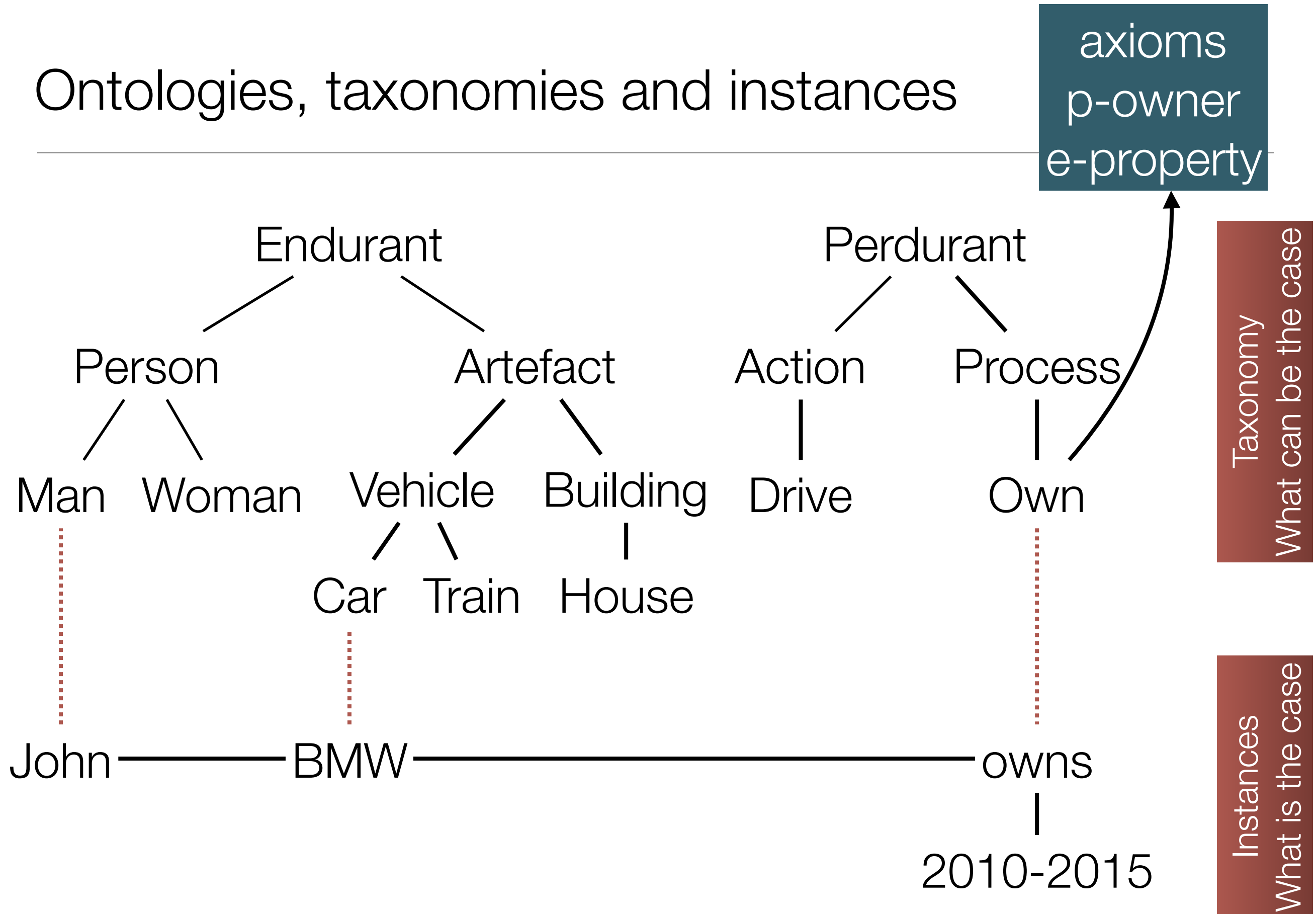
Best performance

Trend

# Performance per relation

- **Factors**:
  - Amount of training data
  - Quality of background data
  - Type of relation:
    - abstract - concrete
    - range of values
    - frequency of the relation
  - Popularity and redundancy
  - Range of possible relations (PER, **?**R, PER)

| Method | P | R | F1 |
|---|---|---|---|
| employee of | 32 | 46 | 38 |
| top members | 26 | 60 | 36 |
| (org:)alt names | 48 | 39 | 43 |
| title | 26 | 35 | 30 |
| spouse | 54 | 85 | 66 |
| origin | 43 | 70 | 53 |
| cause of death | 93 | 39 | 55 |
| children | 62 | 18 | 27 |
| date of death | 64 | 39 | 48 |
| age | 97 | 90 | 93 |

# Ontologies, taxonomies and instances

# Ontology Learning (OL) from text

- Differences with property-value (PV) extraction:

  - PV: non-taxonomic relations between instances (entities):

    - Trump <worksAt> White House

  - OL: taxonomic relations between concepts (nouns):

    - *apples* and *bananas* are *fruits; people workAt organisations*

# Ontology Learning (OL) from text

- Why learn ontologies from text?

  - knowledge bases are incomplete, get out-of-date (products change, new products)

  - too expensive to build, revise and maintain manually

  - more empirical and captures cultural and temporal bound conceptualisations: *pet, toy, weapon, food, drugs*

  - extremely valuable: reasoning, retrieval, generalising patterns and data from instances to types (fighting sparse data problem), as features in text representations for machine learning, embeddings in neural networks
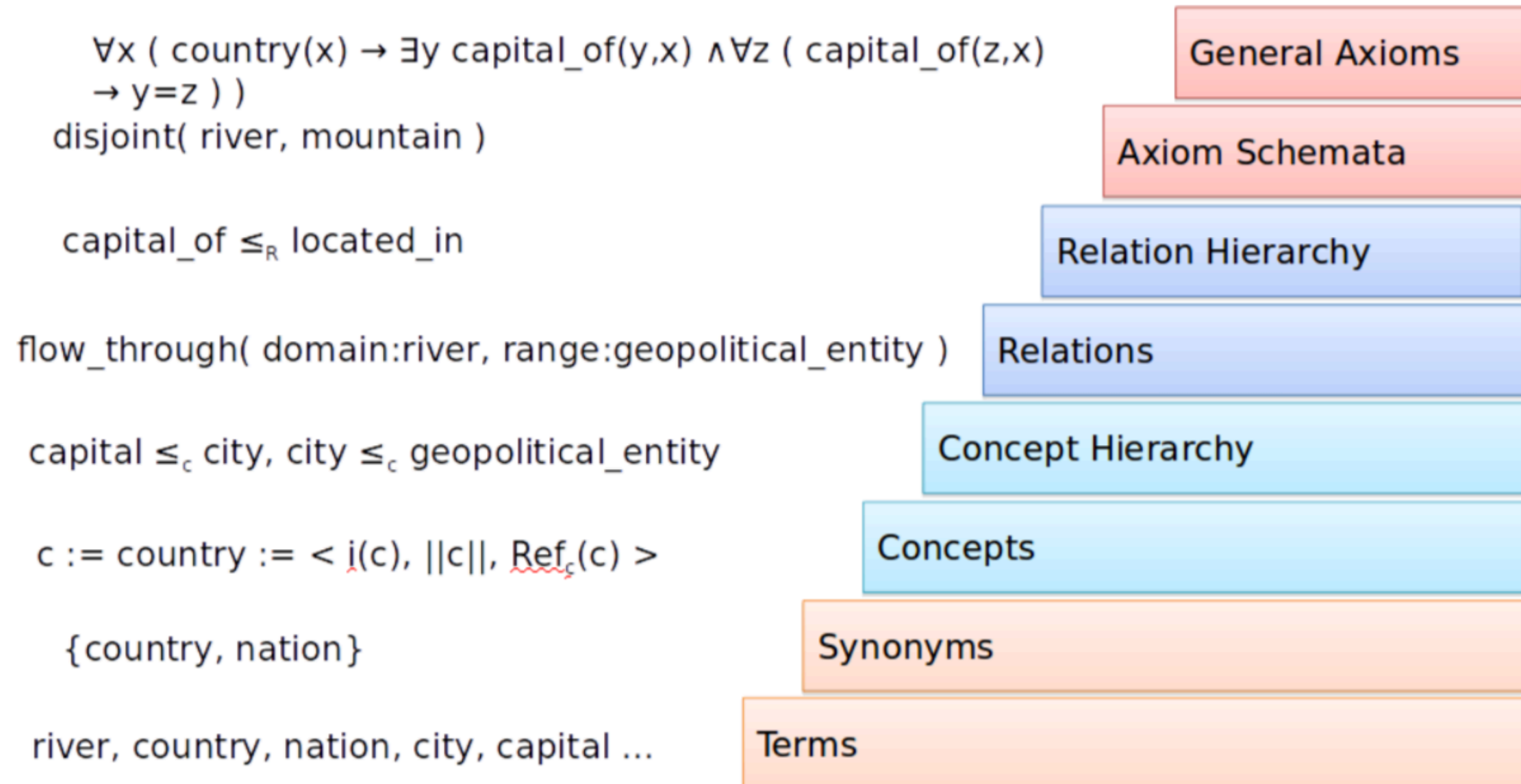
# Ontology learning cake

∀x ( country(x) → ∃y capital_of(y,x) ∧ ∀z ( capital_of(z,x) → y=z ) )
disjoint( river, mountain )

**General Axioms**

**Axiom Schemata**

capital_of ≤_R located_in

**Relation Hierarchy**

flow_through( domain:river, range:geopolitical_entity )

**Relations**

capital ≤_c city, city ≤_c geopolitical_entity

**Concept Hierarchy**

c := country := < i(c), ||c||, Ref_c(c) >

**Concepts**

{country, nation}

**Synonyms**

river, country, nation, city, capital ...

**Terms**

**Figure 6.1:** Ontology learning layer cake (reproduced from Cimiano, P: Ontology Learning and Population from Text: Algorithms, Evaluation and Application, Springer-Verlag, New York, 2006)

# Methods

- What are the terms to learn?:

$$P(D_i \mid t) = \frac{count(t, D_i)}{count(t)},$$

  - Domain text **relevance**: words (single word terms) occurring significantly more frequent than expected (normalised term entropy)

  - **Co-occurrence** statistics: multiword terms "apple juice", "sailing boat", (e.g. point wise mutual information)

$$\mathrm{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)}$$

- What is the relation between terms:

  - **Similarity** & **relatedness**: distributional models or word embeddings

  - **Hyponymy** & **meronymy**: Logical patterns for subtype or isa relations (manually or automatically learned)

# What terms to learn?
## Multiword terms

- Phrases

  - NN, AN, N of N, [A] N Preposition [A] N

  - *apple juice (NN), heavy metal (AN), frequently asked questions (BVN), toxic medication (AN), medication for toxication (NpN)*

$$\mathrm{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)}$$

- Co-occurrence statistics

  frequency

  - Pointwise Mutual Information (pmi)

  - Association Ratio, Jaccard, Log likelihood

# What is the relation between terms?
## Hearst Patterns

- Hearst, Marti A. "Automatic acquisition of hyponyms from large text corpora." Proceedings of the 14th conference on Computational linguistics-Volume 2. Association for Computational Linguistics, 1992.

- What does Gelidium mean?

- "Αγαρ ισ α συβστανχε πρεπαρεδ φρομ α μιξτυρε οφ red algae, such as **Gelidium**, φορ λαβορατορψ ορ ινδυστριαλ υσε"

- "Agar is a substance prepared from a mixture of red algae, such as **Gelidium**, for laboratory or industrial use"

# What is the relation between terms?
## Hearst patterns

| Hearst pattern | Example occurrences |
|---|---|
| X and other  Y | ...temples, treasuries, and other important civic buildings. |
| X or other  Y | Bruises, wounds, broken bones or other injuries... |
| Y such as X | The bow lute, such as the Bambara ndang... |
| Such  Y as X | ...such authors as Herrick, Goldsmith, and Shakespeare. |
| Y including X | ...common-law countries, including Canada and England... |
| Y , especially X | European countries, especially France, England, and Spain... |

# Manually created patterns

- Pros

  - Tend to be high-precision

  - Can be adapted to specific domains

- Cons

  - Tend to be low-recall

  - A lot of work to define all possible patterns

  - Many different patterns are needed for every relation

# Recap

- Relation extraction mainly consists of detecting binary relations between entities

- Rule-based and supervised machine learning methods can be integrated in a bootstrapping method starting from seed data and obtaining more examples

- Open versus closed information extracting, where the former uses clustering and the latter uses a schema

- Schemas can play an important role to define and restrict relations but also combine evidence

- Distant supervision exploits existing data to boost machine learning

- Performance varies depending on various factors, e.g. training data quality and volume, type of relation, difference between train and test especially in time (!)

- Ontology learning boils down to term detection and taxonomic relation detection between terms

- Statistical methods and pattern-based methods are used: the former are more robust and higher recall, the latter have more precision and lower recall