

A creative workspace featuring a laptop with a Windows 8 interface, various pens and pencils in containers, spray paint cans, and a wire mesh organizer. The background is a blurred desk with more supplies.

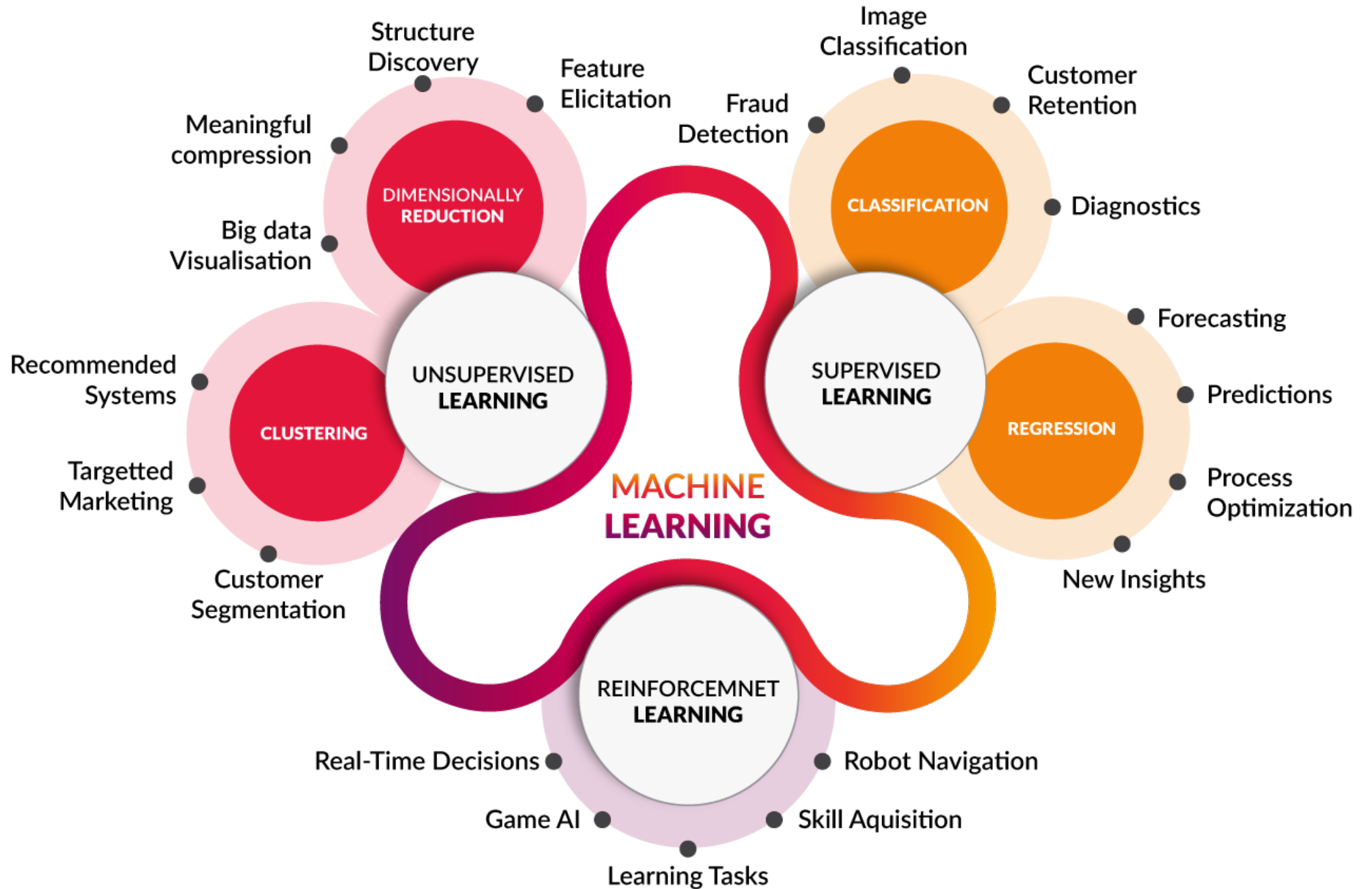
# ONTOLOGY DRIVEN MACHINE LEARNING

-- *Nguyen Hung Son* --



# The Lecture Outline

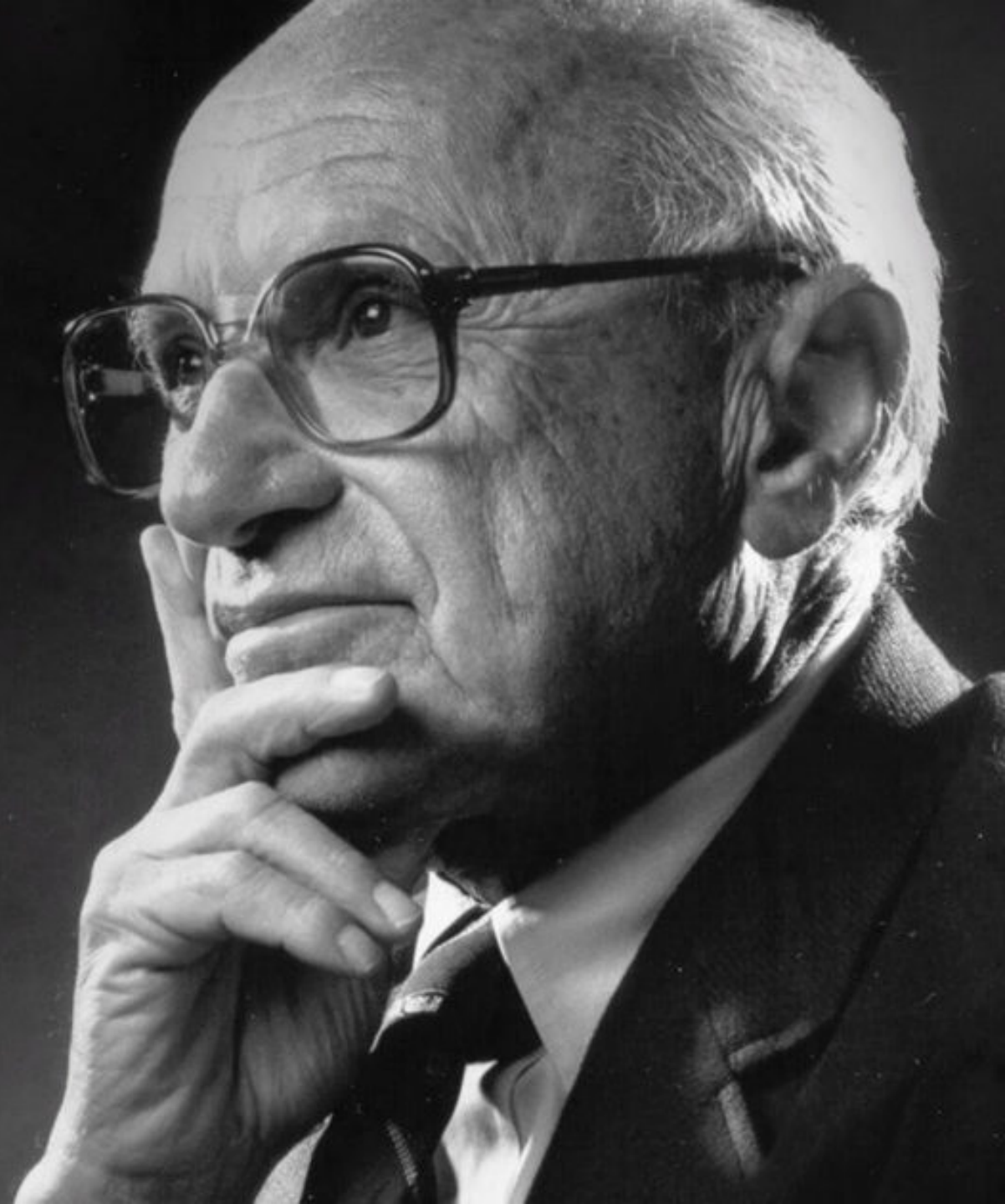
- Domain Knowledge in Machine Learning
- Knowledge representation methods: ontology, taxonomy, ...
- Ontology driven machine learning
  - Ontology and classification
  - Ontology and clustering
  - Semantic evaluation of clustering algorithms
  - Semantic search and text mining
- Machine learning and knowledge acquisition
- Concluding remarks



# Domain knowledge in Machine learning

*Milton Friedman favorite political aphorism:*

“There’s no  
such thing  
as a free lunch.”



# No free lunch theorem

- Assume  $\mathcal{A}$  is a searching algorithm that looking for the maximum of a function  $f : S \rightarrow W$   
where  $S$  is a finite set of states,  $W$  is a finite subset of  $\mathbf{R}$ , and  $f \in \mathcal{F}$
- The work of algorithm  $\mathcal{A}$  after  $t$  steps can be identified by the sequence:  $V_{\mathcal{A}}(f, t) = \left[ (s_1, f(s_1)), (s_2, f(s_2)), \dots, (s_t, f(s_t)) \right]$
- The quality of algorithm  $\mathcal{A}$  can be measured by an evaluation function:  
$$M : \{V_{\mathcal{A}}(f, t) | \mathcal{A}, f, t\} \rightarrow \mathbb{R}$$
- for example:  $M(V_{\mathcal{A}}(f, t)) = \min\{i | f(s_i) = f_{\max}\}$

# No free lunch theorem

- **The class  $\mathcal{F}$  satisfies NFL condition:** if the following equation

$$\sum_{f \in \mathcal{F}} M(V_{\mathcal{A}}(f, |S|)) = \sum_{f \in \mathcal{F}} M(V_{\mathcal{A}'}(f, |S|))$$

holds for any measure  $M$  and any pair of algorithms  $\mathcal{A}, \mathcal{A}'$

- **$\mathcal{F}$  is closed under permutation:** for any permutation  $\sigma \in \text{Perm}(S)$  and  $f \in \mathcal{F}$  we have  $\sigma f \in \mathcal{F}$

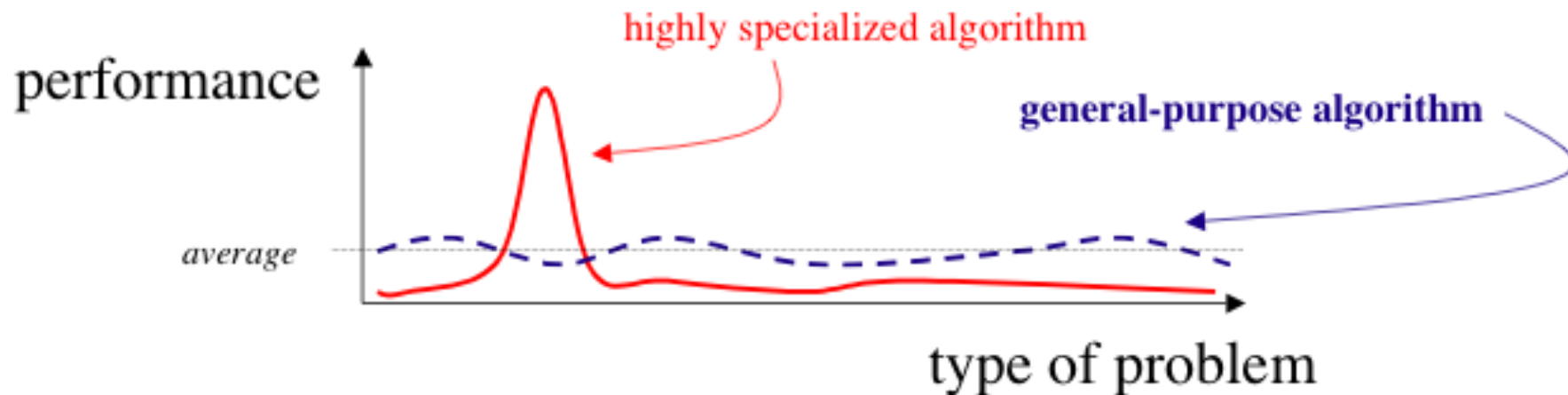
## **NFL theorem:**

*the class  $\mathcal{F}$  satisfies NFL condition iff  $\mathcal{F}$  is closed under permutation*

# No free lunch

- For example: the class of all functions from  $S$  to  $W$  is closed under permutation
- The probability that a random class of functions from  $S$  to  $W$  is closed under permutation equals

$$\frac{2^{\binom{|S|+|W|-1}{|S|}} - 1}{2^{|S||W|} - 1}$$



# No free lunch theorem for learning

Wolpert (1996) shows that in a noise-free scenario where the loss function is the misclassification rate, if one is interested in off-training-set error, then there are no a priori distinctions between learning algorithms.

More formally, where

- $d$  = training set;
- $m$  = number of elements in training set;
- $f$  = 'target' input-output relationships;
- $h$  = hypothesis (the algorithm's guess for  $f$  made in response to  $d$ ); and
- $C$  = off-training-set 'loss' associated with  $f$  and  $h$  ('generalization error')
- all algorithms are equivalent, on average, by any of the following measures of risk:  $E(C|d)$ ,  $E(C|m)$ ,  $E(C|f,d)$ , or  $E(C|f,m)$ .



## Therefore ...

- No search or learning algorithm can be the best on all possible learning or optimization problems.
- In fact, every algorithm is the best algorithm for the same number of problems.
- But only some problems are of interest.
- For example:  
*a random search algorithm is perfect for a completely random problem (the "white noise" problem), but for any search or optimization problem with structure, random search is not so good.*

## KNOWLEDGE PRESENTATION METHODS:

classification,  
taxonomy,  
ontology,  
thesaurus.

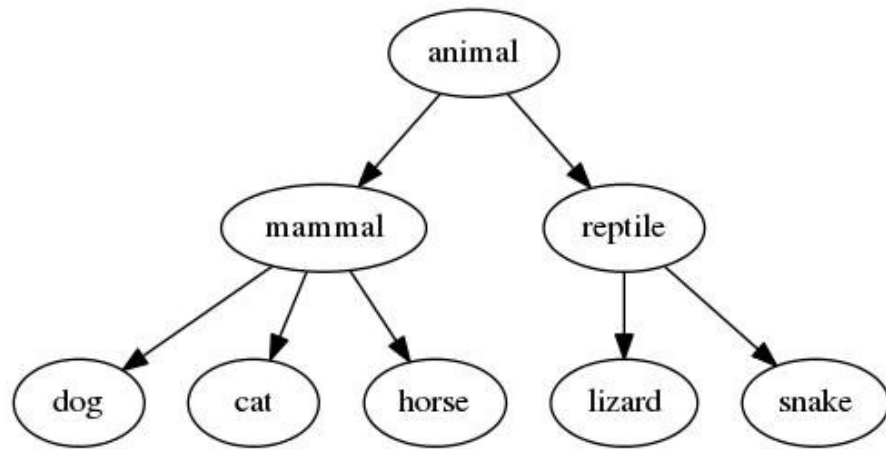
# Merriam-Webster definition

## Classification

- *systematic arrangement in groups or categories according to established criteria*

## Taxonomy

- *orderly classification of plants and animals according to their presumed natural relationships.*



## Ontology

- *Old: a branch of metaphysics concerned with the nature and relations of being or a particular theory about the nature of being or the kinds of existents.*
- *New: machine-readable set of definitions that create a taxonomy of classes and subclasses and relationships between them*

## Thesaurus:

- *a thesaurus deals only with words, alternatives for those words, synonyms, translations, et cetera*
- *can be used by a classification, a taxonomy and an ontology*

# Ontology vs taxonomy

## Taxonomy:

- Sub-concept relation only
- A proper taxonomy is a strict hierarchy (one parent), e.g.  
*Natural Science (500)*
  - > *Zoological Sciences (590)*
    - > *Other Invertebrates (595)*
      - > *Insects (595.7)*
        - > *Lepidoptera (595.78)*
          - > *Butterflies (595.789)*.
- Relaxed model: multi-parent but still noncyclic e.g.

Computer Accessories

Cell Phone Accessories

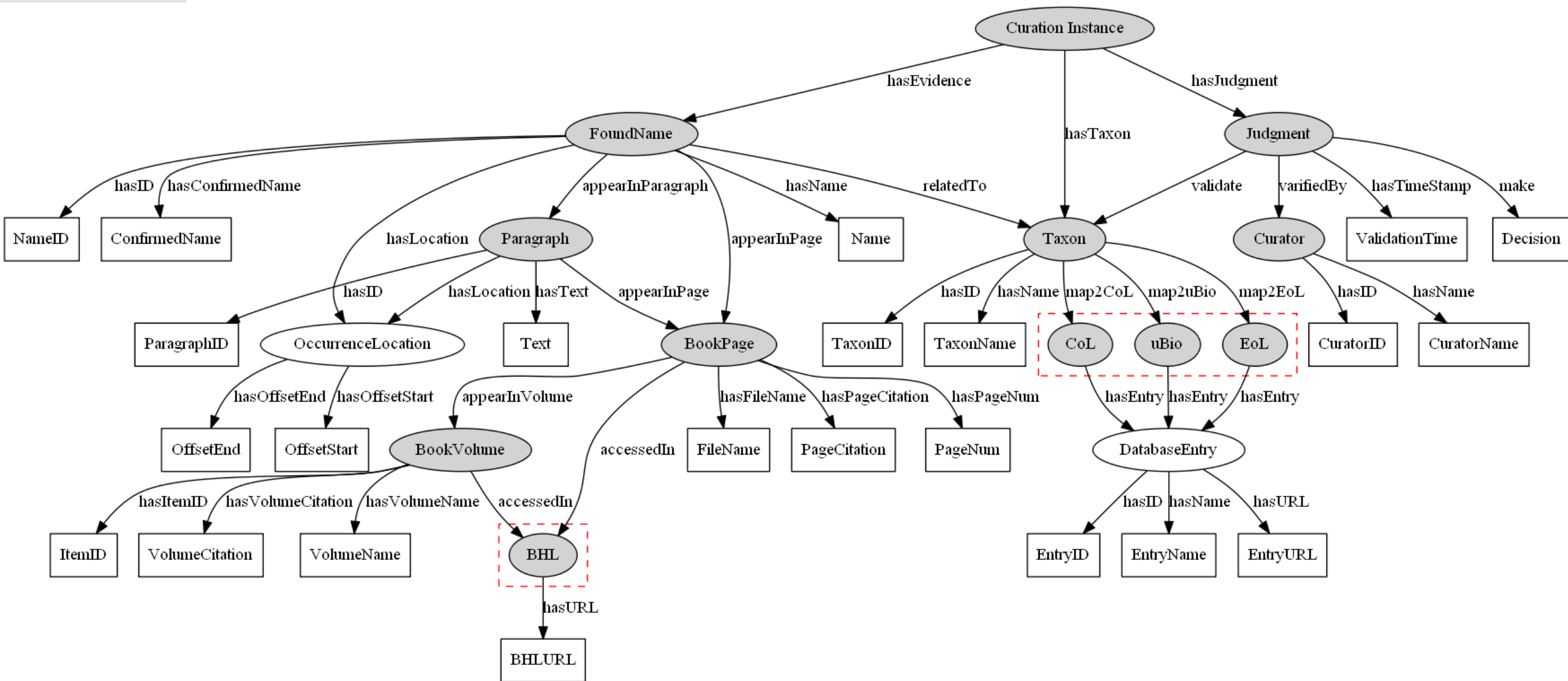
USB Cables

```
graph TD; A[Computer Accessories] --> C[USB Cables]; B[Cell Phone Accessories] --> C;
```

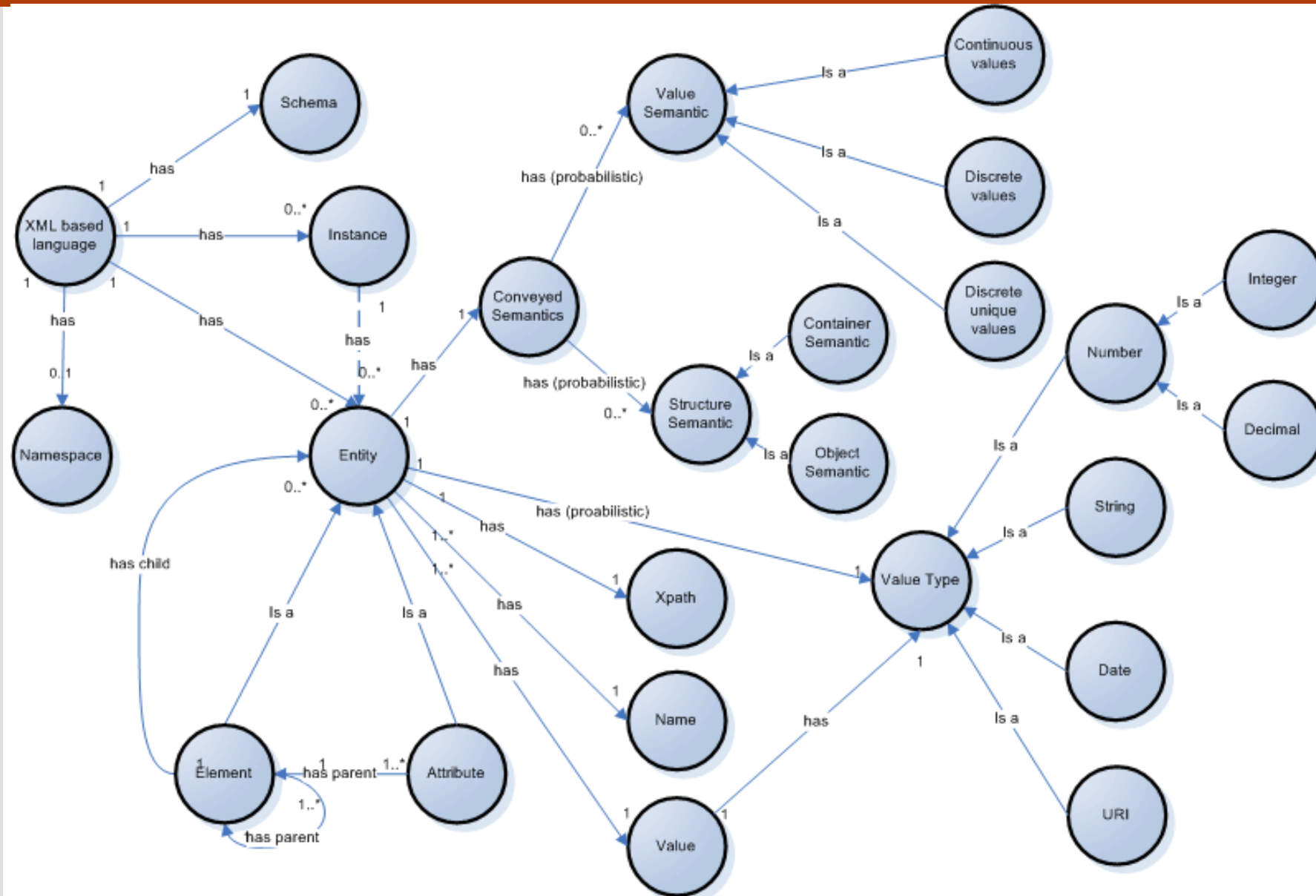
## Ontology

- Objects could be:
  - classes
  - instances of the class
  - class attributes
- More types of relations:
  - is-a
  - has-a
  - use-a
- the relationships aren't necessarily binary—for example, a co-worker

# Example of taxonomic name curation



# Example of XML entities ontology

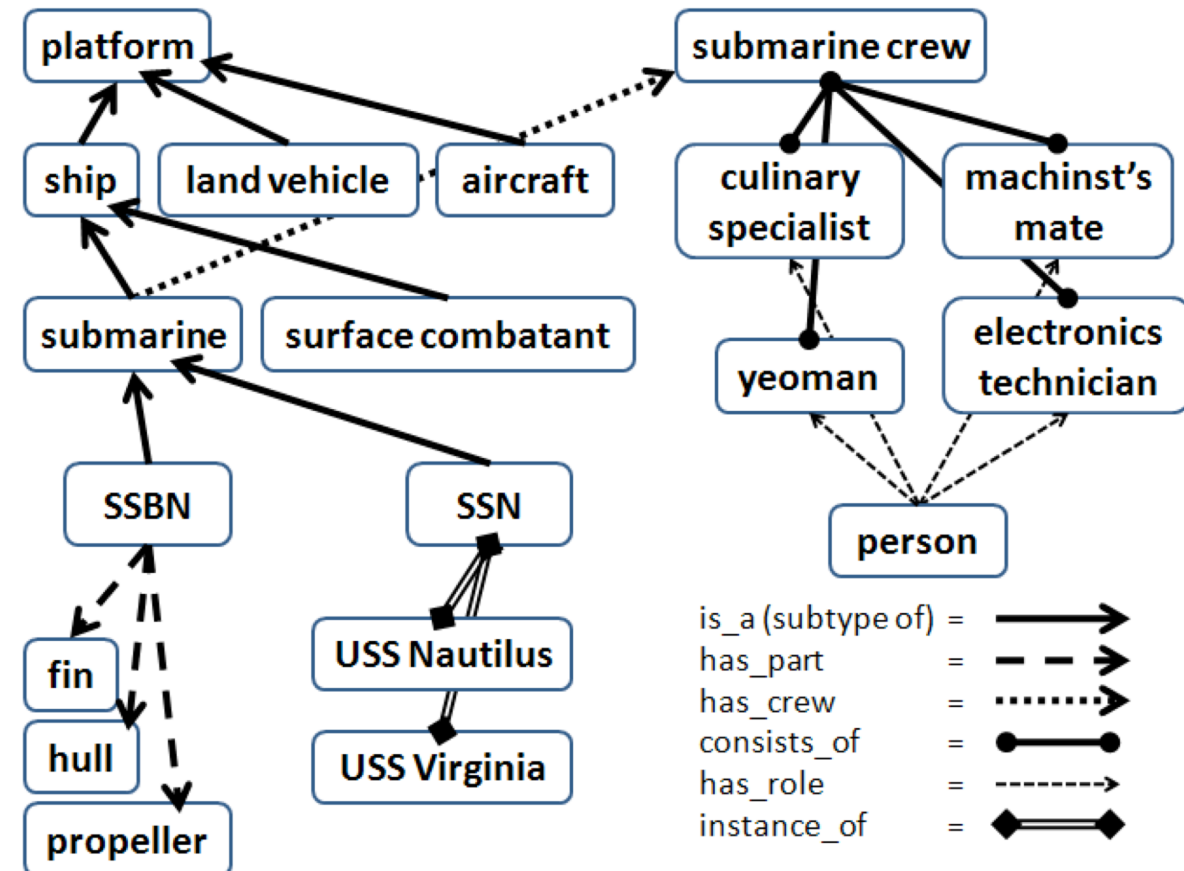


# What is ontology?

A structured, **taxonomic model** or representation of the **entities** and **relations** existing within a particular domain of reality.

For examples:

- Gene ontology,
- lexical ontology (wordnet),
- Ontology for General Medical Science
- other domain ontology



# Ontology libraries & repositories

## LIBRARIES:

- A library system that offers various functions for managing, adapting and standardizing groups of ontologies.
- It should fulfill the needs for re-use of ontologies. In this sense, an ontology library system should be easily accessible and offer efficient support for reusing existing relevant ontologies and standardizing them based on upper-level ontologies and ontology representation languages.

## REPOSITORIES:

- A structured collection of ontologies (...) by using an Ontology Metadata Vocabulary.
- References and relations between ontologies and their modules build the semantic model of an ontology repository. Access to resources is realized through semantically-enabled interfaces applicable for humans and machines.
- Therefore, a repository provides a formal query language.

### Ontology libraries

[OBO Foundry](#)

WebProtégé

Romulus

DAML ontology library

Colore

VEST/AgroPortal Map of standards

[FAIRsharing](#)

DERI Vocabularies

OntologyDesignPatterns

SemanticWeb.org

W3C Good ontologies

TaxoBank

BARTOC

GFBio Terminology Service

agINFRA Linked Data Vocabularies

oeGOV

### Ontology repositories

[NCBO BioPortal\\*](#)

[Ontobee](#)

[EBI Ontology Lookup Service](#)

[AberOWL](#)

[CISMEF HeTOP](#)

[SIFR BioPortal\\*](#)

OKFN Linked Open Vocabularies

ONKI Ontology Library Service

MMI Ontology Registry and Repository\*

# Applications

## Vertical need

- For those uses who want to do very precise things, e.g.
  - reasoning,
  - using specific relationsusing only suitable ontologies (developed by the same communities and in the same format).
- For those users who may just use the repositories as libraries to find and download ontologies, and work in their own environment.

## Horizontal need

- For those who wants to work with a wide range of ontologies and vocabularies useful in their domain but developed by different communities, overlapping and in different formats.
- Such users greatly appreciate the unique endpoints (Web application and programmatic for REST and SPARQL queries) offered by the repositories under a simplified common model.

# Ontology: challenges and applications

- **Metadata & selection**
- **Multilingualism**
- **Ontology alignment**
- **Generic ontology-based services**
- **Annotations and Linked Data**
- **Scalability & interoperability**

# Leslie Valiant

- A fundamental question for AI is to characterize the computational **building blocks** that are necessary for cognition.
- A specific challenge is to build on the success of machine learning so as to cover broader issues in intelligence.
- This requires, in particular a reconciliation between two contradictory characteristics
  - The apparent logical nature of reasoning and
  - the statistical nature of learning.
- Professor Valiant has developed a formal system, called **robust logics**, that aims to achieve such a reconciliation.

- *Turing Award, 2010*
- *European Association for Theoretical Computer Science Award, 2008*
- *Knuth Prize, 1997*
- *Nevanlinna Prize, 1986*



JANUARY 7, 2013

Computer scientist  
Leslie Valiant named  
2012 ACM Fellow

# Ontology driven methods for Machine Learning and AI

A background image on the left side of the slide showing a laptop screen with a Windows-style interface, a notebook with colorful patterns, and a piece of paper with handwritten notes.

# Case studies

1. Robocup
2. Semantic Text processing and mining
3. Approximate reasoning

# ROBOCUP: robocup.org

[HOME](#)[ABOUT](#)[ORGANIZATION](#)[LEAGUES](#)[EVENTS](#)[NEWS](#)[RESEARCH](#)[GALLERY](#)[INFO](#)

## RoboCupSoccer - Simulation



This is one of the oldest leagues in RoboCupSoccer. The Simulation League focus on artificial intelligence and team strategy. Independently moving software players (agents) play soccer on a virtual field inside a computer. There are 2 subleagues: 2D and 3D.

## || LEAGUES

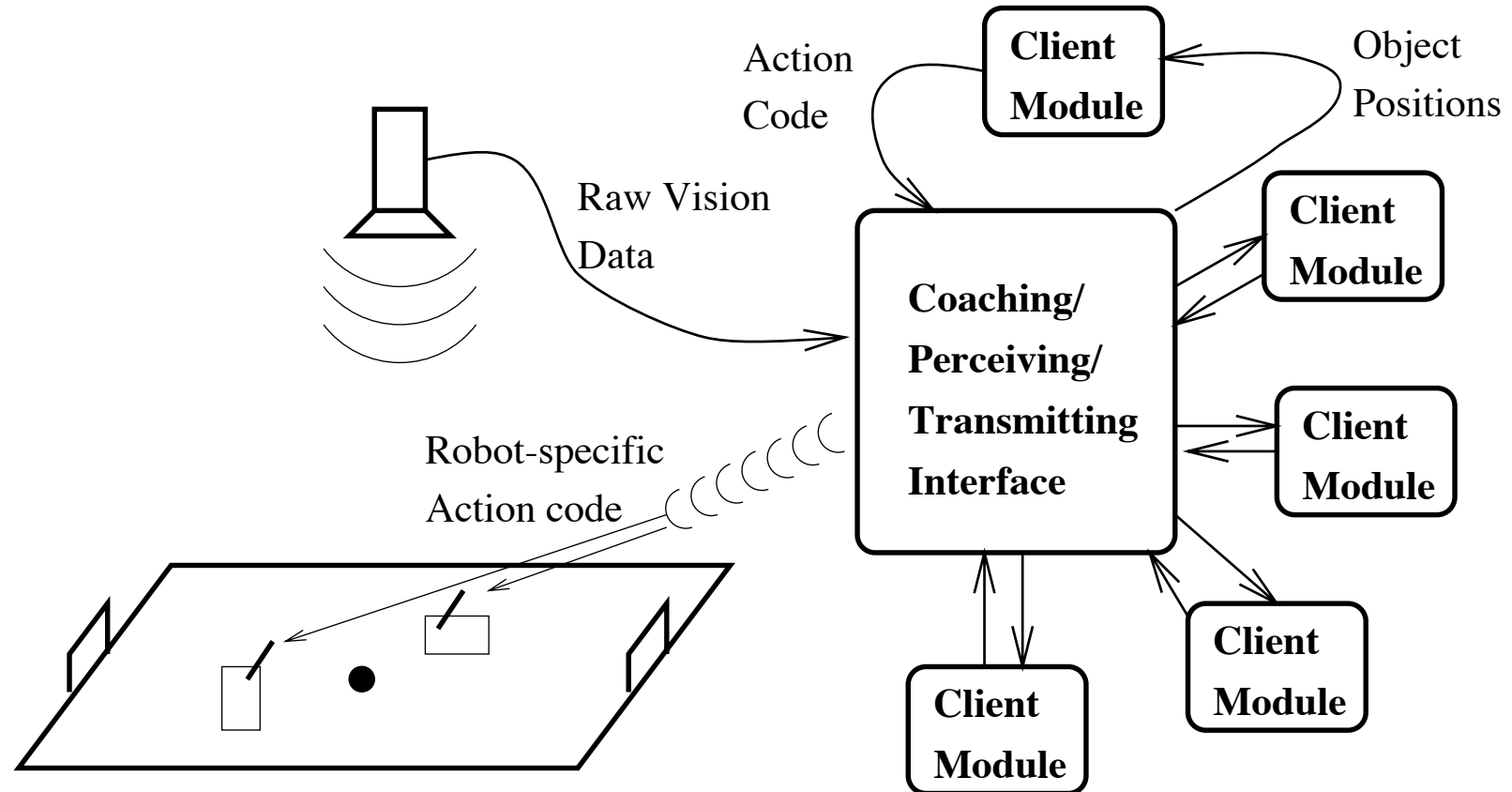
[Humanoid](#)[Standard Platform](#)[Middle Size](#)[Small Size](#)[Simulation](#)

## || SUB-LEAGUES

[Simulation 2D](#)[Simulation 3D](#)

# robotic soccer architecture

as a distributed deliberative and reactive system



# Robocup and simulated robotic soccer

## Noda's Soccer Server

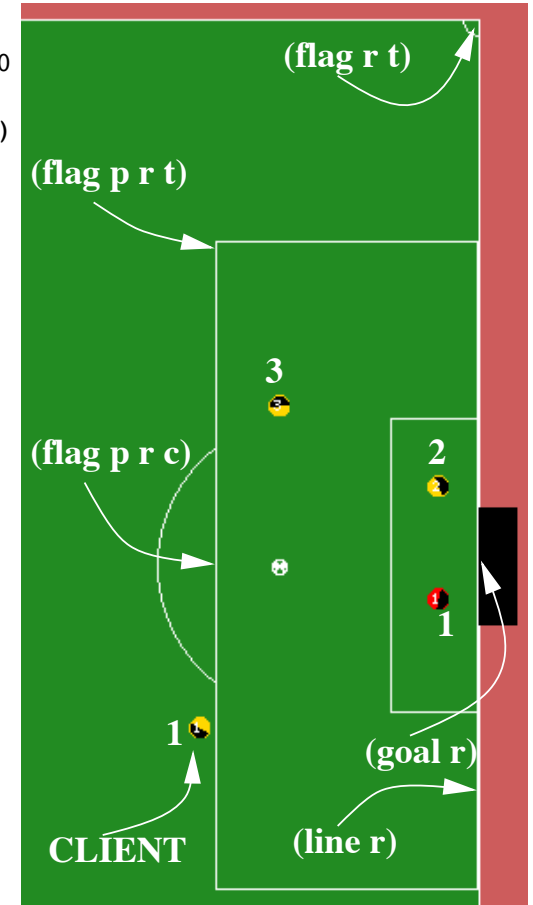
- the players' vision is limited ( $45^\circ$ );
- the players can communicate by posting to a blackboard that is visible to all players;
- all players are controlled by separate processes;
- each player has 10 teammates and 11 opponents;
- each player has limited stamina;
- actions and sensors are noisy; and
- play occurs in real time.

The simulator, acting as a server, provides a domain and supports users who wish to build their own agents (clients).



# Example

```
(see 124 ((goal r) 20.1 34) ((flag r t) 47.5 -4) ((flag p r t) 30.3 -24) ((flag p r c) 10.1 -20)
((ball) 11 0) ((player usa 2) 21 19) ((player usa 3) 21 -11) ((player brazil 1) 17 35) ((line r) 40
**-> (dash 80)
(see 129 ((goal r) 16 43) ((flag r t) 42 -6) ((flag p r t) 25 -30) ((flag p r c) 5 -40) ((ball) 6 1)
((player usa 2) 16.3 24) ((player usa 3) 15.3 -17) ((line r) 32.8 -27))
**-> (turn 1)
**-> (dash 60)
(see 134 ((flag r t) 40 -9) ((flag p r t) 23.3 -35) ((ball) 3.7 2) ((player usa 2) 14.4 24)
((player usa 3) 13.3 -22) ((line r) 28.2 -30))
**-> (turn 2)
**-> (dash 30)
(hear 138 18 shoot the ball)
(see 139 ((flag r t) 38.1 -11) ((flag p r t) 22 -39) ((ball) 1.9 0) ((player usa 2) 12.8 27)
((player usa 3) 11.6 -27) ((line r) 25.5 -31))
**-> (say shooting now)
**-> (kick 53 51)
(hear 141 self shooting now)
(see 144 ((flag r t) 38.1 -11) ((flag p r t) 22 -39) ((ball) 8.1 42) ((player usa 2) 12.8 27)
((player usa 3) 11.6 -27) ((line r) 25.5 -31))
**-> (turn 42)
(see 149 ((goal r) 13.6 9) ((ball) 13.5 5 0) ((player usa 2) 12.8 -14) ((player brazil 1) 11 18)
((line r) 14 -73))
**-> (turn 5)
**-> (dash 81)
(hear 150 referee goal_l_1)
(hear 150 referee kick_off_r)
```





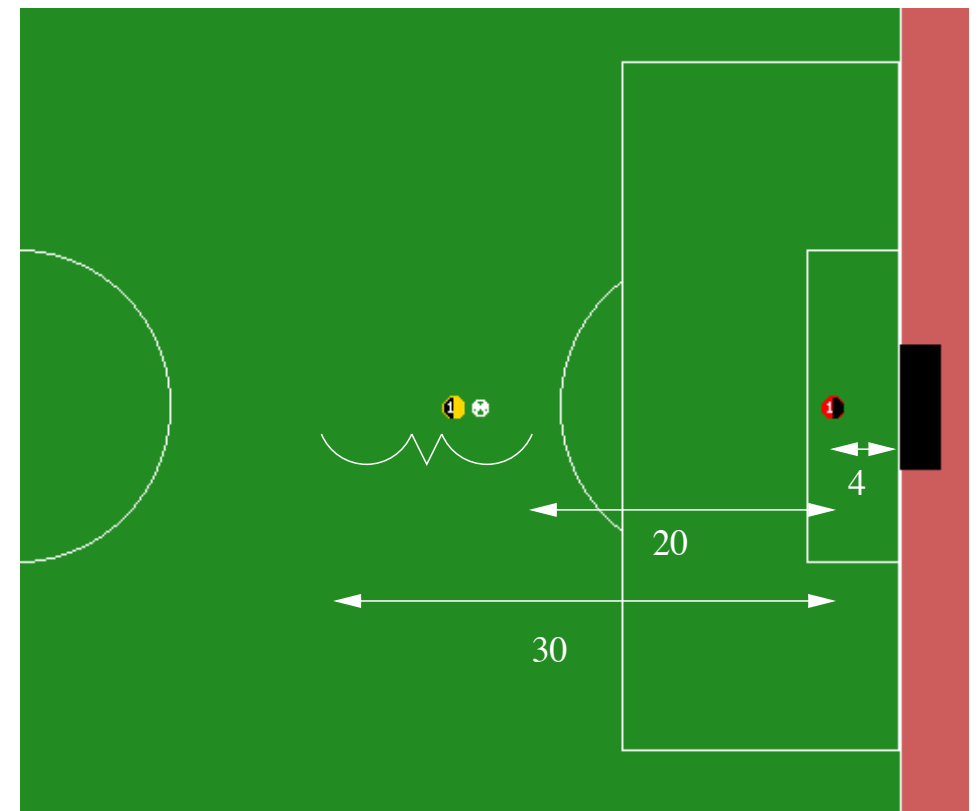
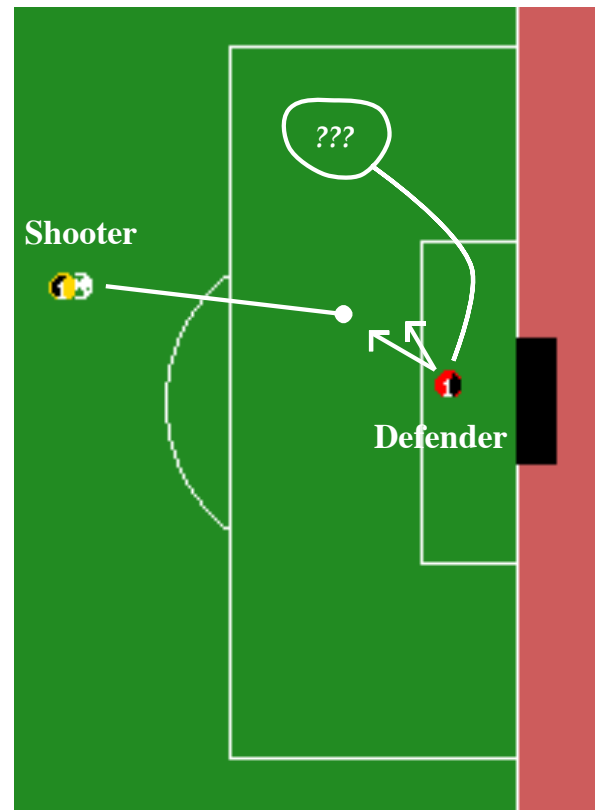
# simulated robotic soccer

- An example of MAS;
- Enough complexity to be realistic;
- Easy accessibility to researchers worldwide;
- Embodiment of most MAS issues: reactivity, modeling, cooperation, competition, role playing, resource management, communication, convention, commitment/decommitment strategies
- Straightforward evaluation
- Good multiagent ML opportunities.

# Learning a lower-level skill

## Intercepting a moving ball:

- Co-Learning for the shooter and the defender
- Using neural networks (NN)

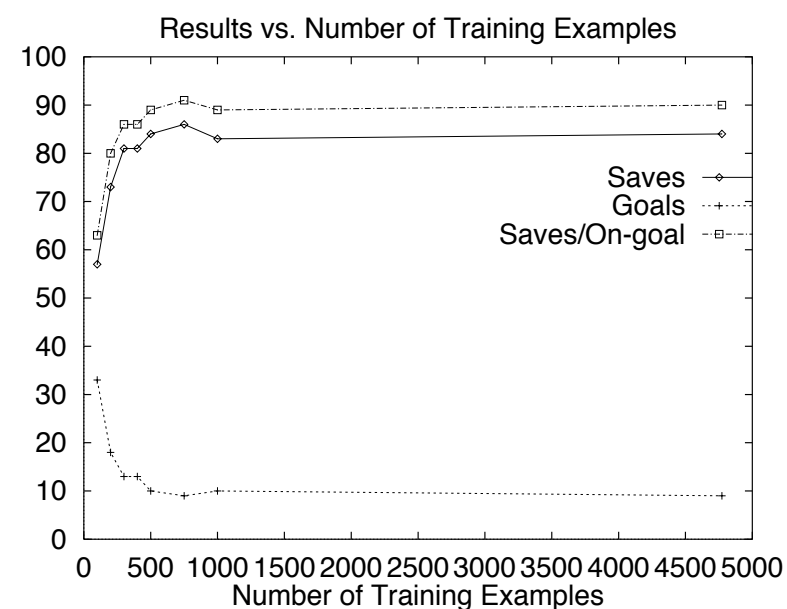


# Performance

## *A Layered Approach to Learning Client Behaviors in the RoboCup Soccer Server*

Peter Stone Manuela Veloso 

Training Examples	Saves		
	Saves(%)	Goals(%)	$\frac{\text{Saves}}{\text{Goals+Saves}}(\%)$
100	57	33	63
200	73	18	80
300	81	13	86
400	81	13	86
500	84	10	89
750	<b>86</b>	9	<b>91</b>
1000	83	10	89
4773	84	9	90



# Learning a Higher-level Decision: pass, dribble, shoot



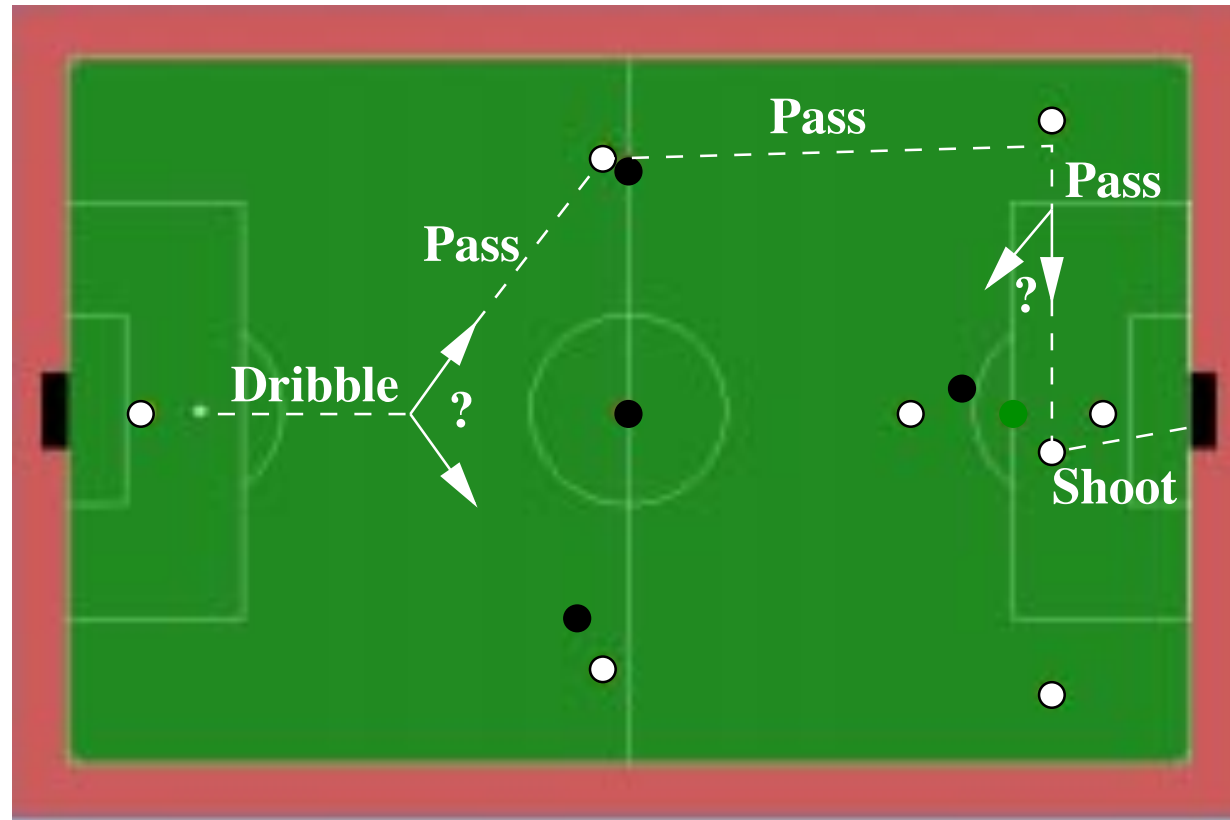
# Learning a Higher-level Decision: passing

174 attributes were used to construct a Decision Tree:

- Distance and Angle to the receiver (2);
- [?] Distance and Angle to other teammates (up to 9) sorted by angle from the receiver (18);
- [?] Distance and Angle to opponents (up to 11) sorted by angle from the receiver (22);
- Counts of teammates, opponents, and players within given distances and angles of the receiver (45);
- Distance and Angle from receiver to teammates (up to 10) sorted by distance (20);
- [?] Distance and Angle from receiver to opponents (up to 11) sorted by distance (22);
- [?] Counts of teammates, opponents, and players within given distances and angles of the passer from the receiver's perspective (45);

Result	Overall	Success Confidence:		
		.8–.9	.7–.8	.6–.7
(Number)	(5000)	(1050)	(3485)	(185)
SUCCESS (%)	65	<b>79</b>	63	58
FAILURE (%)	26	15	29	31
MISS (%)	8	5	8	10

# Team-level Strategies



○ Teammate

● Defender

## Next layers:

- More flexible and powerful approach would be to allow the dribbling player to learn:
  - when to continue dribbling,
  - when to pass, and
  - when to shoot.
- With these three possibilities as the action space and with appropriate predicates to discretize the state space, **TD-lambda and other reinforcement learning methods** will be applicable.
- By keeping track of whether an opponent or a teammate possesses the ball next, a player can propagate reinforcement values for each decision made while it possesses the ball.

A background image showing a portion of a laptop on the left side, displaying a Windows-style desktop with various icons. Below the laptop, there are some papers with handwritten notes and colorful, abstract drawings. The main part of the slide is a solid orange rectangle containing the title and a list of bullet points.

## Next layer

- Learning moving behavior to be a targeted receiver
- Learn to cooperate with the teammates, learn to thwart the opponents
- Last updates allow to have one more agent: the coach (trainer)

# ROBOCUP: summary

## The Dream

We proposed that the ultimate goal of the RoboCup Initiative to be stated as follows:

“ *By the middle of the 21st century, a team of fully autonomous humanoid robot soccer players shall win a soccer game, complying with the official rules of FIFA, against the winner of the most recent World Cup.* ”

- Challenging project that cover many issues in AI and Data Science
- P. Stone. *Layered Learning in Multiagent Systems : A Winning Approach to Robotic Soccer.*
- Why soccer (football)?

**Layered Learning in  
Multiagent Systems**

A Winning Approach to Robotic Soccer

**Peter Stone**

# Semantic text processing

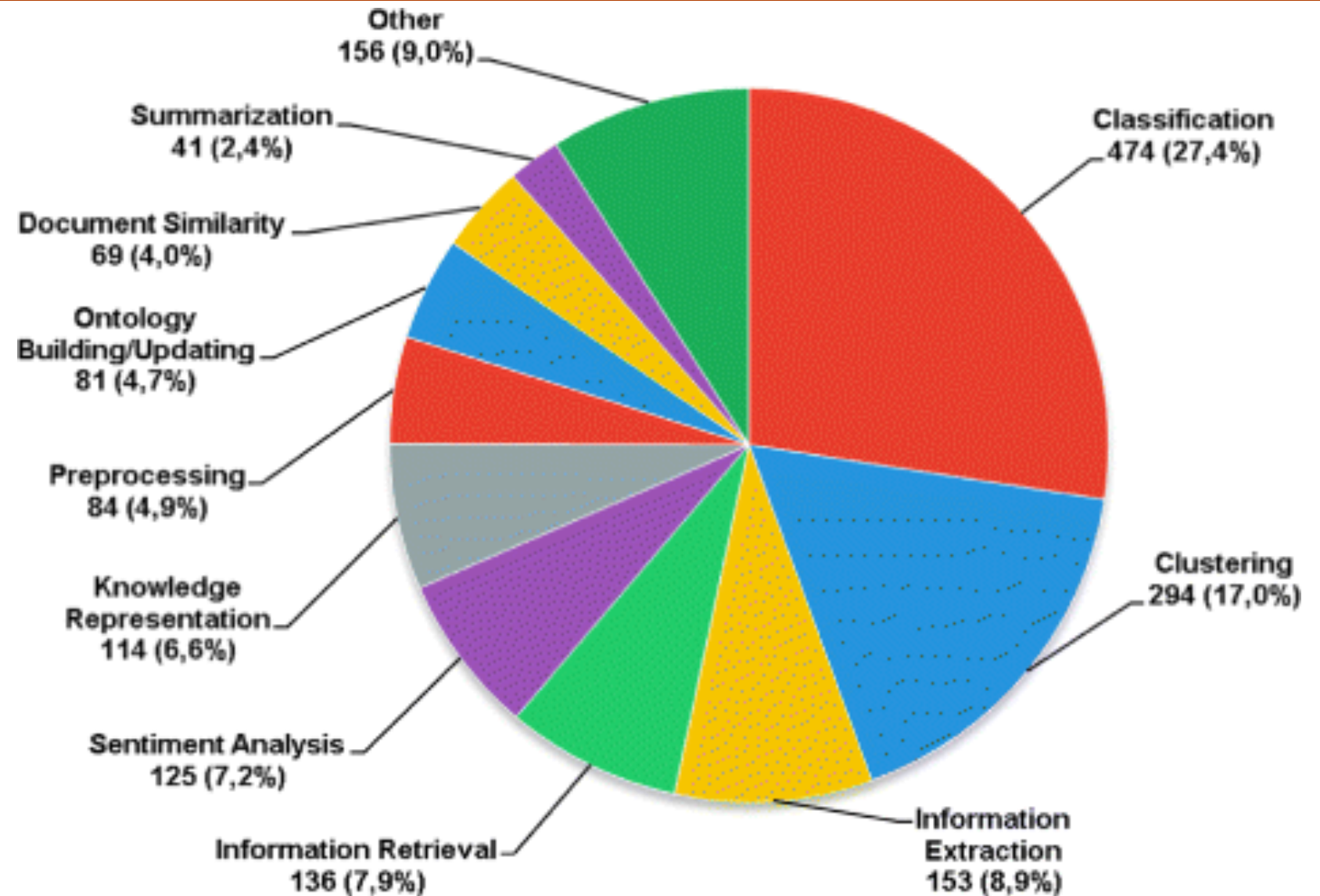


A background image showing a portion of a laptop on the left and a notebook with handwritten notes at the bottom left. The laptop screen displays a Windows-style interface with various application icons. The notebook has some illegible handwriting in blue ink.

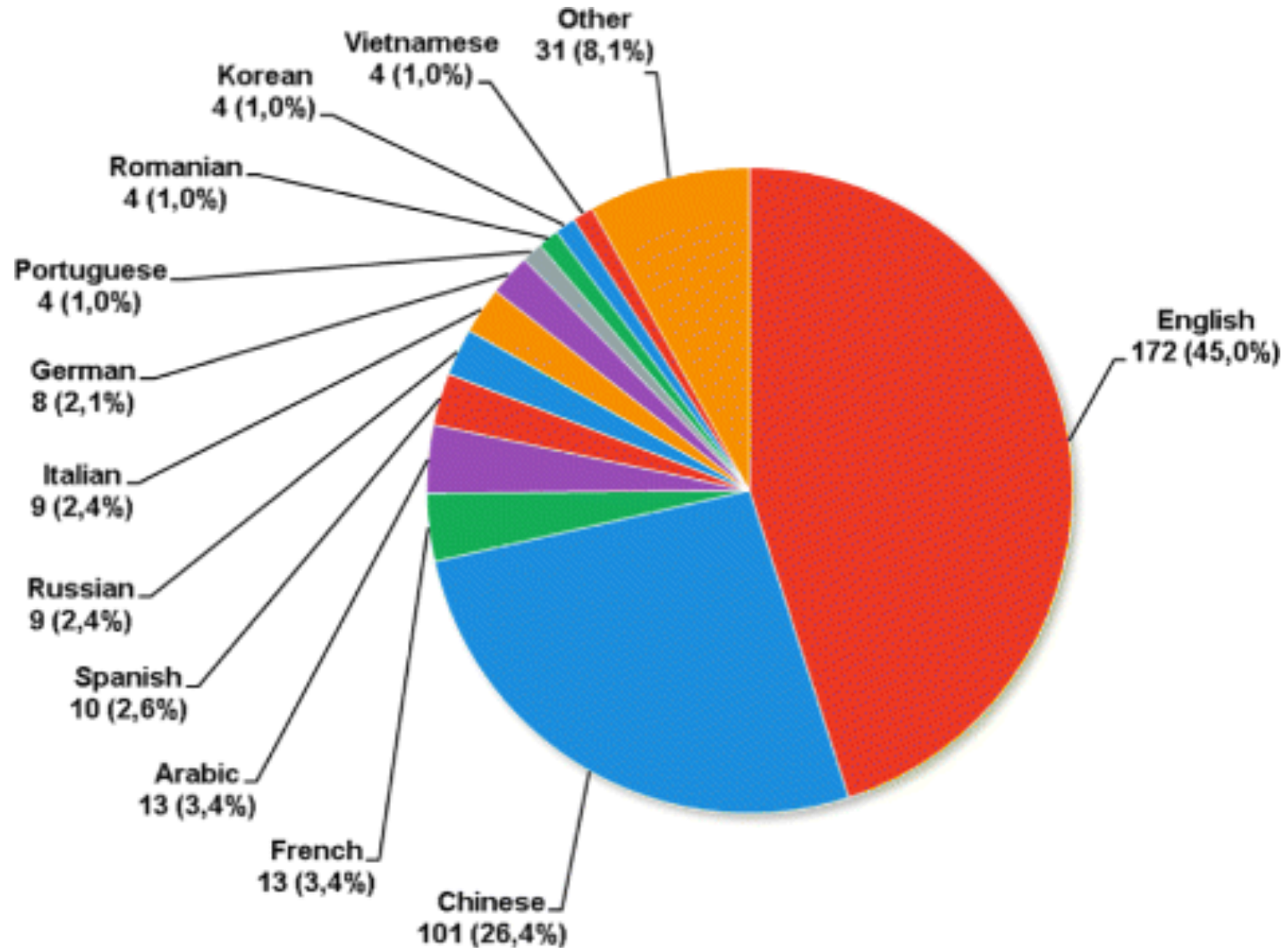
# What for?

- Extract relevant and useful information from large bodies of unstructured data
- Find an answer to a question without having to ask a human
- Discover the meaning of colloquial speech in online posts
- Uncover specific meanings of words used in foreign languages mixed with our own

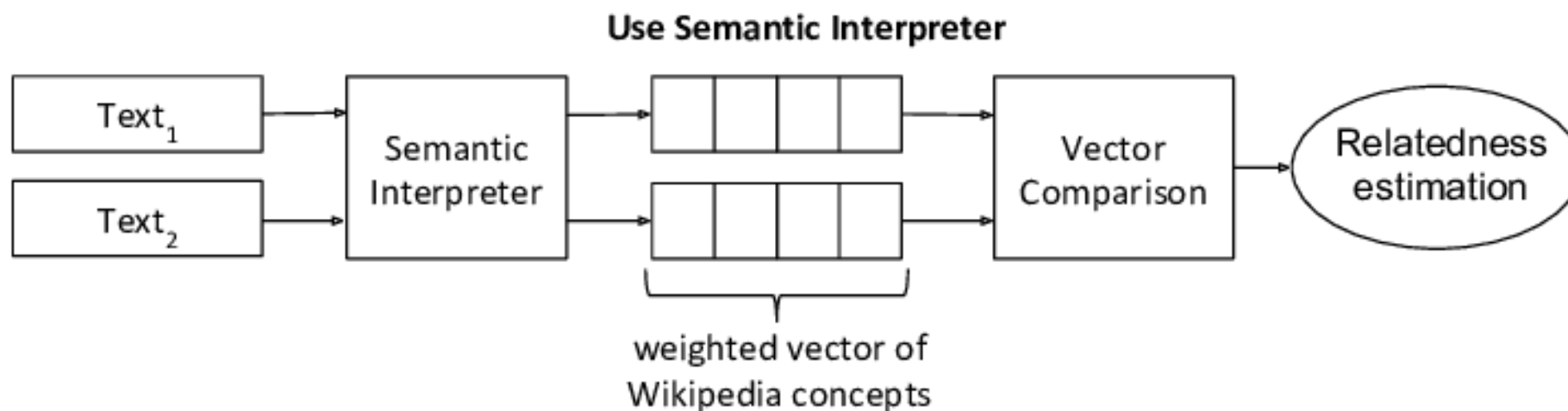
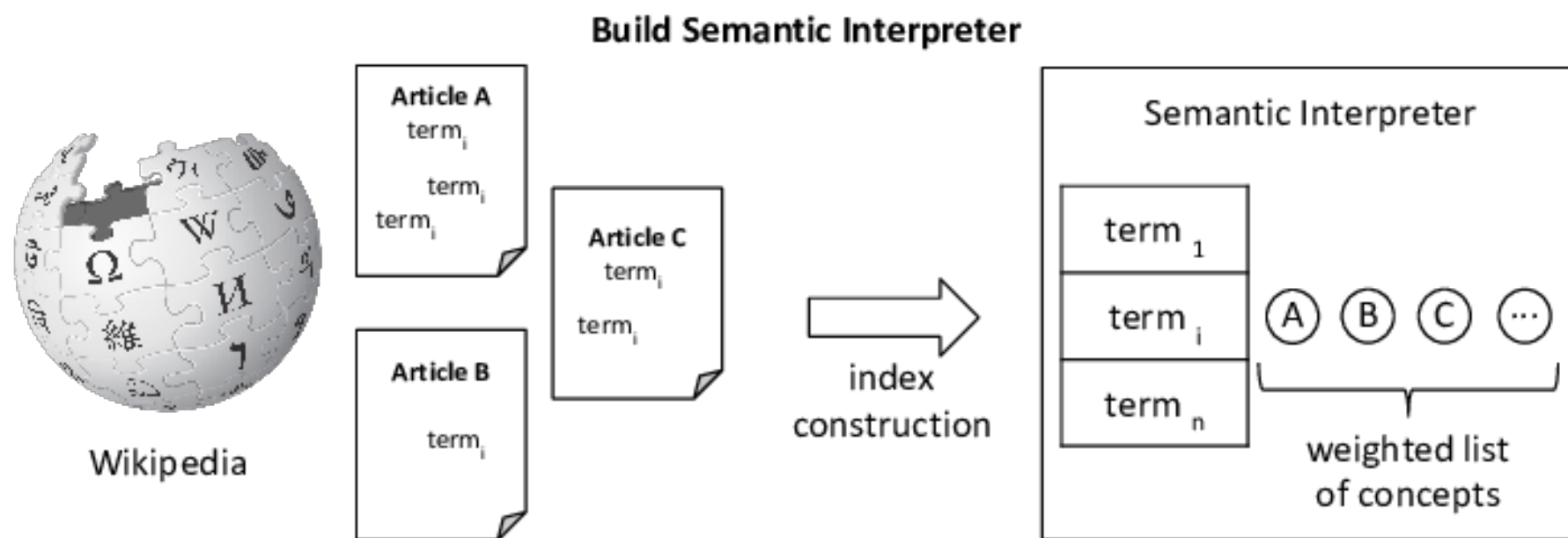
# Semantic text mining tasks:



# What are the natural languages being considered when working with text semantics?



# Explicit Semantic Analysis (ESA)



# Semantic interpreter

	<b>term<sub>0</sub></b>	<b>term<sub>1</sub></b>	...	<b>term<sub>M</sub></b>
<b>doc<sub>0</sub></b>	$w_{00}$	$w_{01}$	...	$w_{0M}$
<b>doc<sub>1</sub></b>	$w_{10}$	...	...	...
...	...	...	$w_{ij}$	...
<b>doc<sub>N</sub></b>	$w_{N0}$	...	...	$w_{NM}$

Representation of system data

	<b>concept<sub>0</sub></b>	<b>concept<sub>1</sub></b>	...	<b>concept<sub>K</sub></b>
<b>term<sub>0</sub></b>	$c_{00}$	$c_{01}$	...	$c_{0K}$
<b>term<sub>1</sub></b>	$c_{10}$	...	...	...
...	...	...	$c_{jk}$	...
<b>term<sub>M</sub></b>	$c_{M0}$	...	...	$c_{MK}$

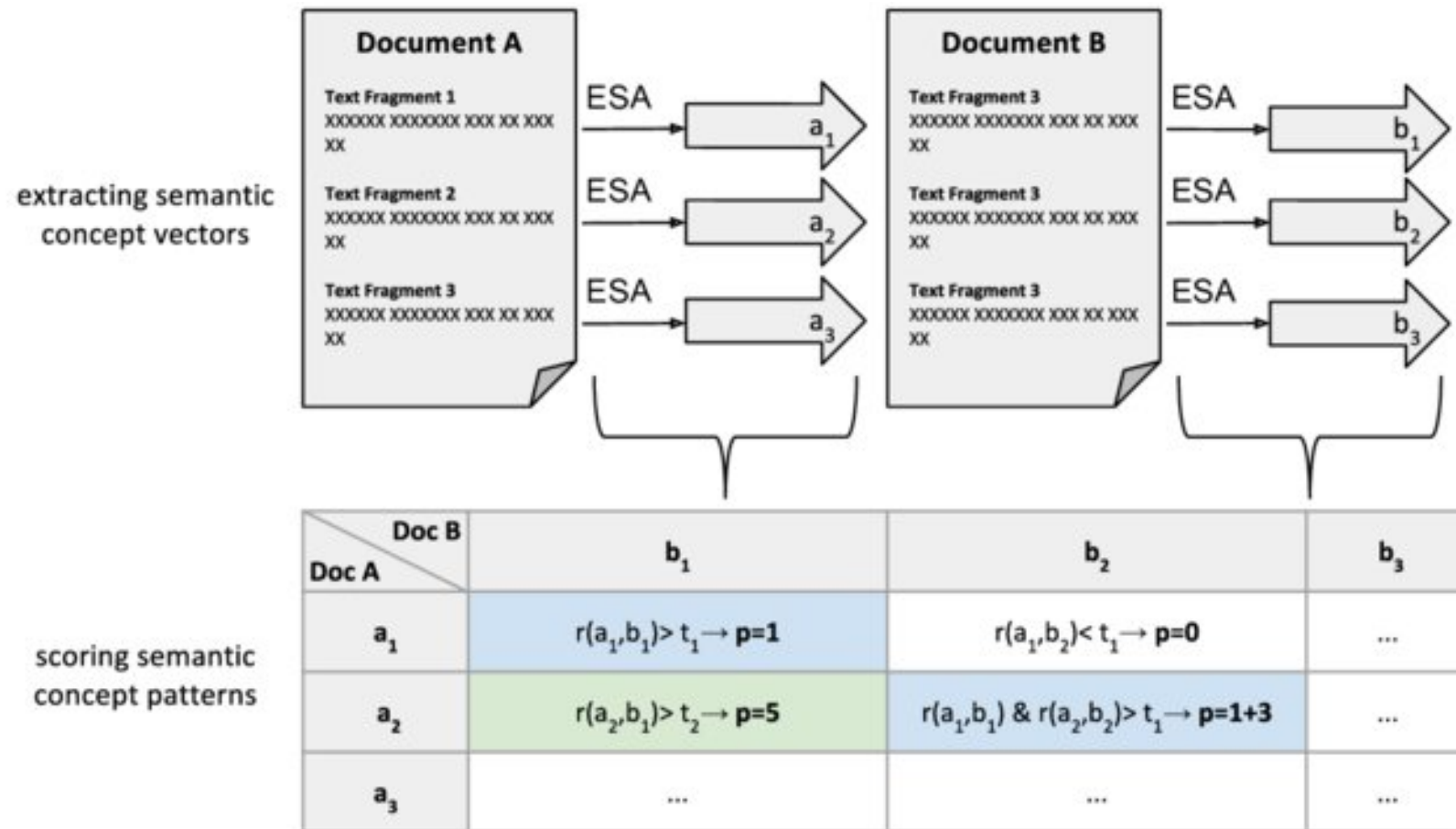
Representation of knowledge base

$$u_{ik} = \sum_{t_j \in T} w_{ij} \times c_{jk}$$

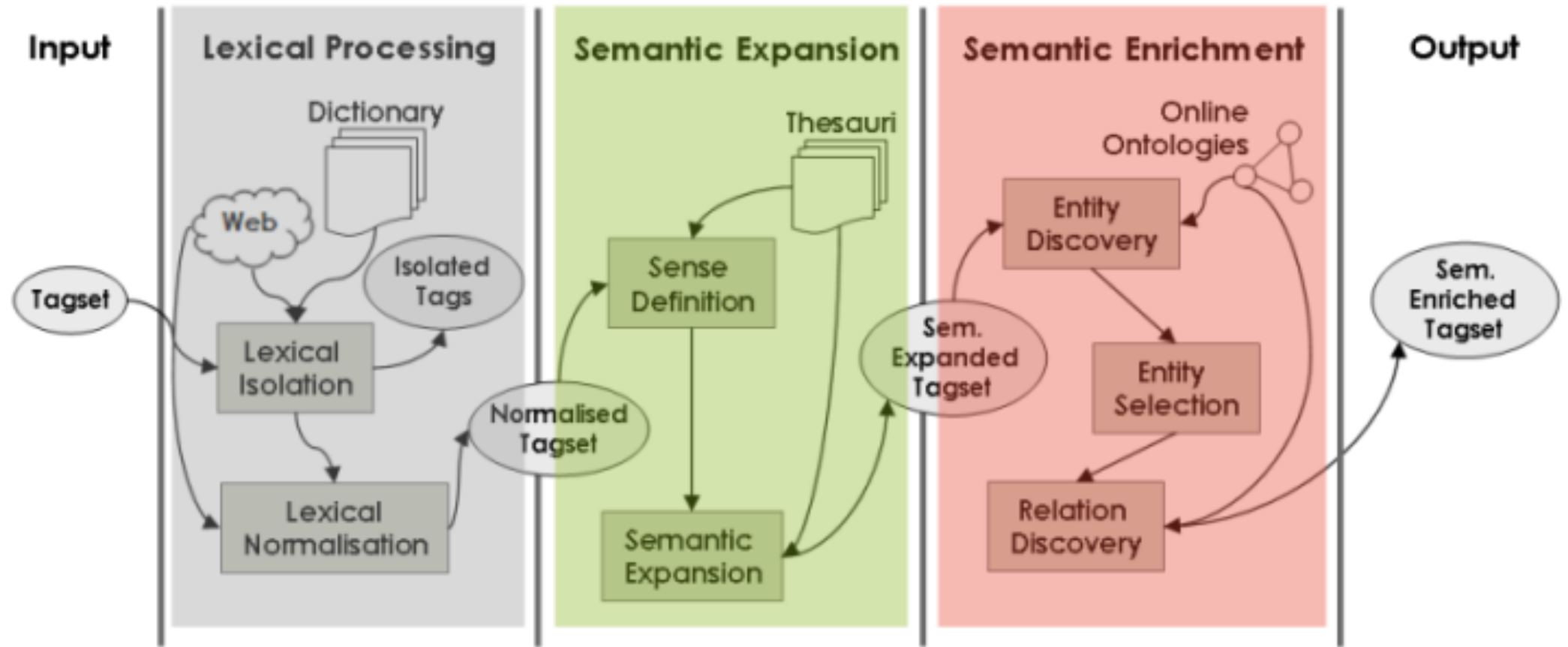
	<b>concept<sub>0</sub></b>	<b>concept<sub>1</sub></b>	...	<b>concept<sub>K</sub></b>
<b>doc<sub>0</sub></b>	$u_{00}$	$u_{01}$	...	$u_{0K}$
<b>doc<sub>1</sub></b>	$u_{10}$	...	...	...
...	...	...	$u_{ik}$	...
<b>doc<sub>N</sub></b>	$u_{N0}$	...	...	$u_{NK}$

New representation of system data

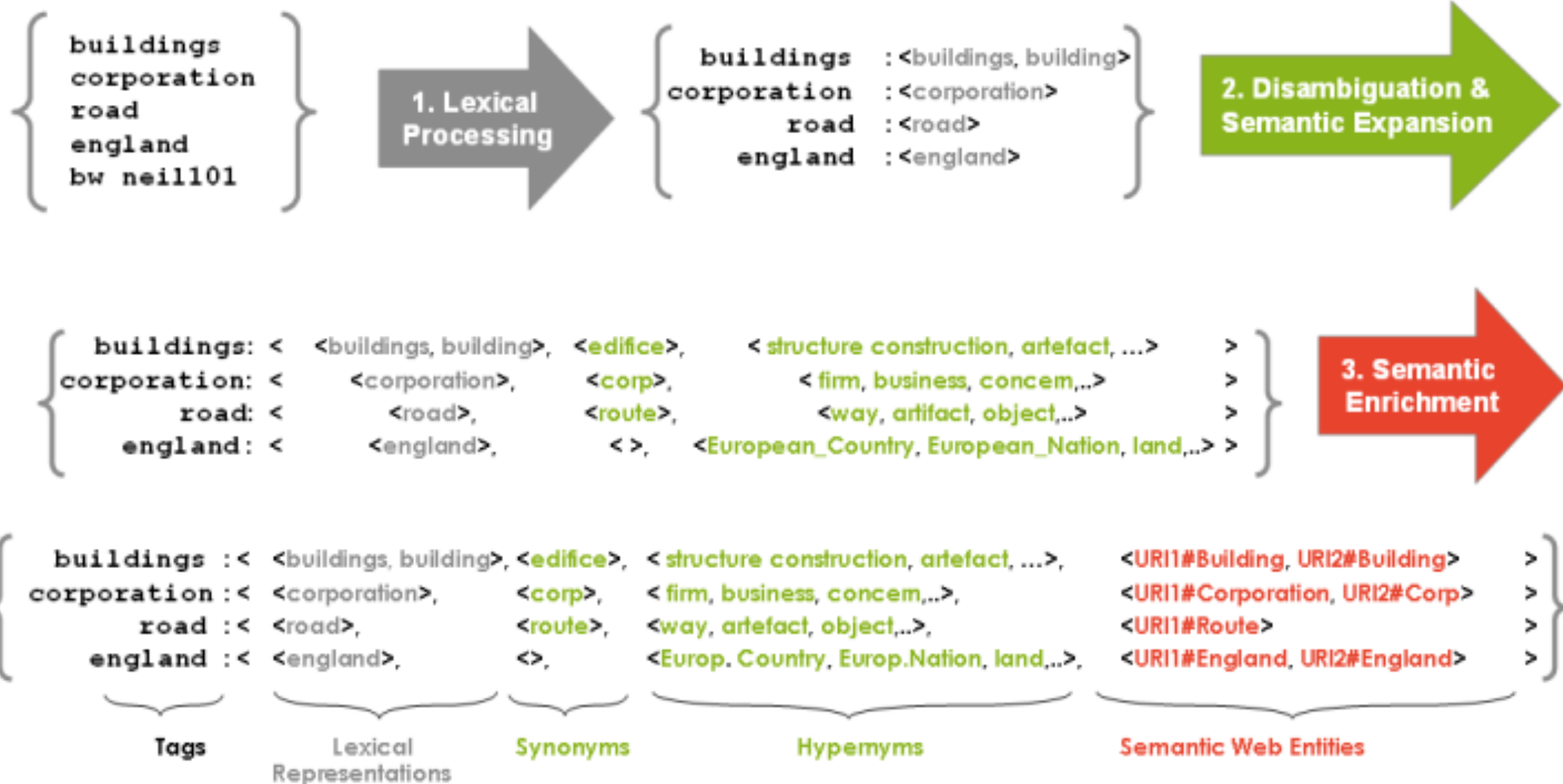
# Example: Semantic similarity



# Semantic enrichment



# Example



# Improved ESA

	term <sub>0</sub>	term <sub>1</sub>	...	term <sub>M</sub>
doc <sub>0</sub>	w <sub>00</sub>	w <sub>01</sub>	...	w <sub>0M</sub>
doc <sub>1</sub>	w <sub>10</sub>	...	...	...
...	...	...	w <sub>ij</sub>	...
doc <sub>N</sub>	w <sub>N0</sub>	...	...	w <sub>NM</sub>

Representation of system data

	concept <sub>0</sub>	concept <sub>1</sub>	...	concept <sub>K</sub>
term <sub>0</sub>	c <sub>00</sub>	c <sub>01</sub>	...	c <sub>0K</sub>
term <sub>1</sub>	c <sub>10</sub>	...	...	...
...	...	...	c <sub>jk</sub>	...
term <sub>M</sub>	c <sub>M0</sub>	...	...	c <sub>MK</sub>

Representation of knowledge base

$$u_{ik} = \sum_{t_j \in T} w_{ij} \times c_{jk}$$

	concept <sub>0</sub>	concept <sub>1</sub>	...	concept <sub>K</sub>
doc <sub>0</sub>	u <sub>00</sub>	u <sub>01</sub>	...	u <sub>0K</sub>
doc <sub>1</sub>	u <sub>10</sub>	...	...	...
...	...	...	u <sub>ik</sub>	...
doc <sub>N</sub>	u <sub>N0</sub>	...	...	u <sub>NK</sub>

New representation of system data



# JRS'2012 Data mining competition

Large biomedical document repositories, such as MEDLINE, hire experts to index their resources with MeSH terms.

- MeSH contains over 26,000 main *headings*.
- Headings can be used in a context of 83 qualifiers (*subheadings*).
- Medical doctors use MeSH *heading/subheading* pairs to search for information.
- 670,943 articles were indexed (semi-)manually in 2007.



Experts need support in their work.

- Over 1 million articles are expected in 2015...

# Challenges

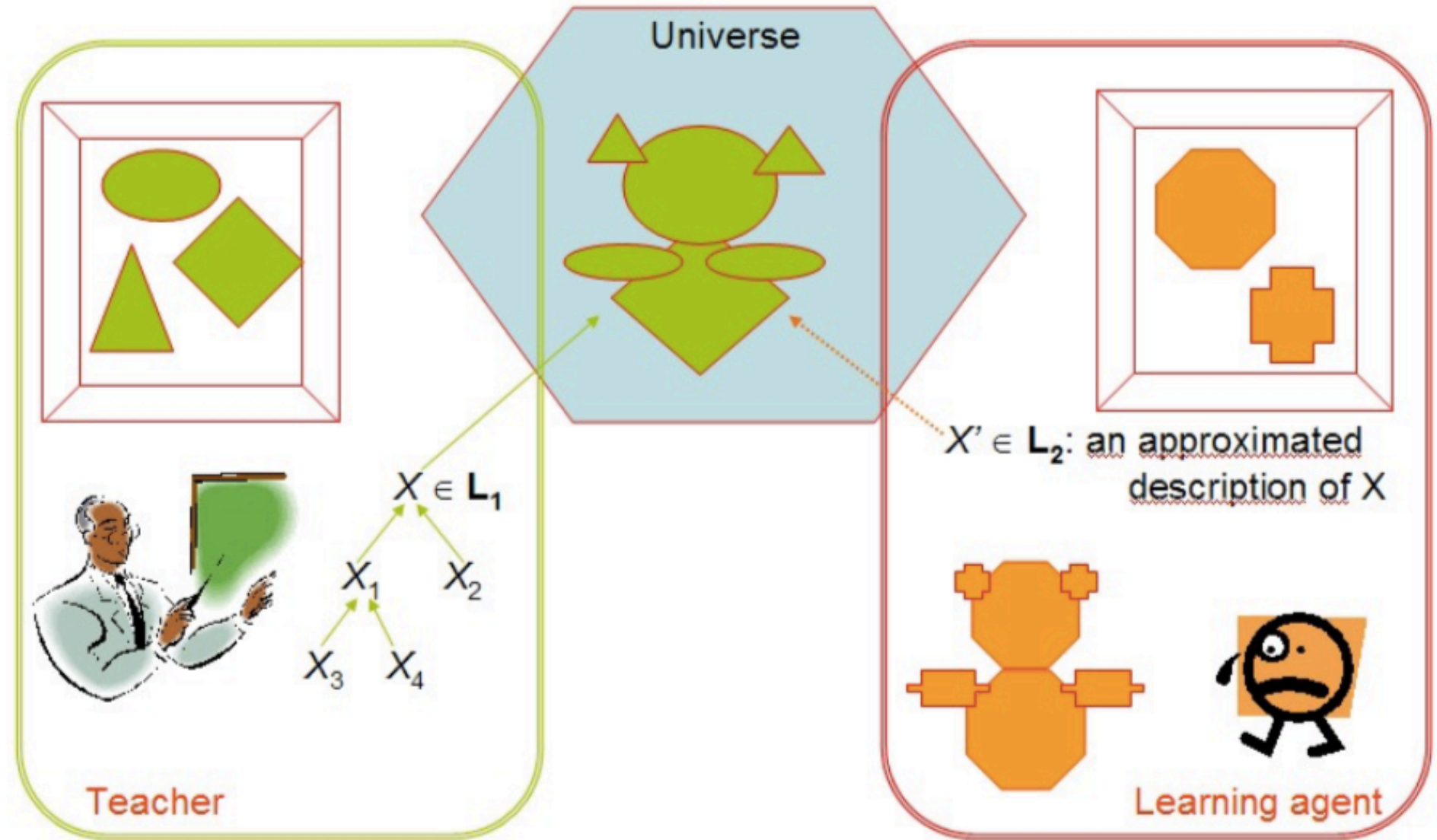
- **Scalability:** deeper semantic analysis vs time and space complexity
- Text representation model.
- Semantic analysis < text understanding
- Example: Named Entity Recognition (NER) problem:

- One of the major problems in NER is ambiguous names: e.g. one protein name may refer to multiple gene products
- Example: using sense-tagged corpora and unified medical language system (UMLS) to resolve ambiguous terms.

Machine-learning techniques have been applied to sense-tagged corpora, in which senses (or concepts) of ambiguous terms have been most manually annotated

>>> quite an expansive manual work

# Concept approximation



# Example nursery data set

- Creator: Vladislav Rajkovic et al. (13 experts)
- Donors: Marko Bohanec (marko.bohanec@ijs.si)  
Blaz Zupan (blaz.zupan@ijs.si)
- Date: June, 1997
- Number of Instances: 12960 (instances completely cover the attribute space)
- Number of Attributes: 8

## Attributes

NURSERY	not_recom, recommend, very_recom, priority, spec_prior
. EMPLOY	<i>Employment of parents and child's nursery</i>
. . parents	usual, pretentious, great_pret
. . has_nurs	proper, less_proper, improper, critical, very_crit
. STRUCT_FINAN	<i>Family structure and financial standings</i>
. . STRUCTURE	<i>Family structure</i>
. . . form	complete, completed, incomplete, foster
. . . children	1, 2, 3, more
. . housing	convenient, less_conv, critical
. . finance	convenient, inconv
. SOC_HEALTH	<i>Social and health picture of the family</i>
. . social	non-prob, slightly_prob, problematic
. . health	recommended, priority, not_recom

# Layered learning

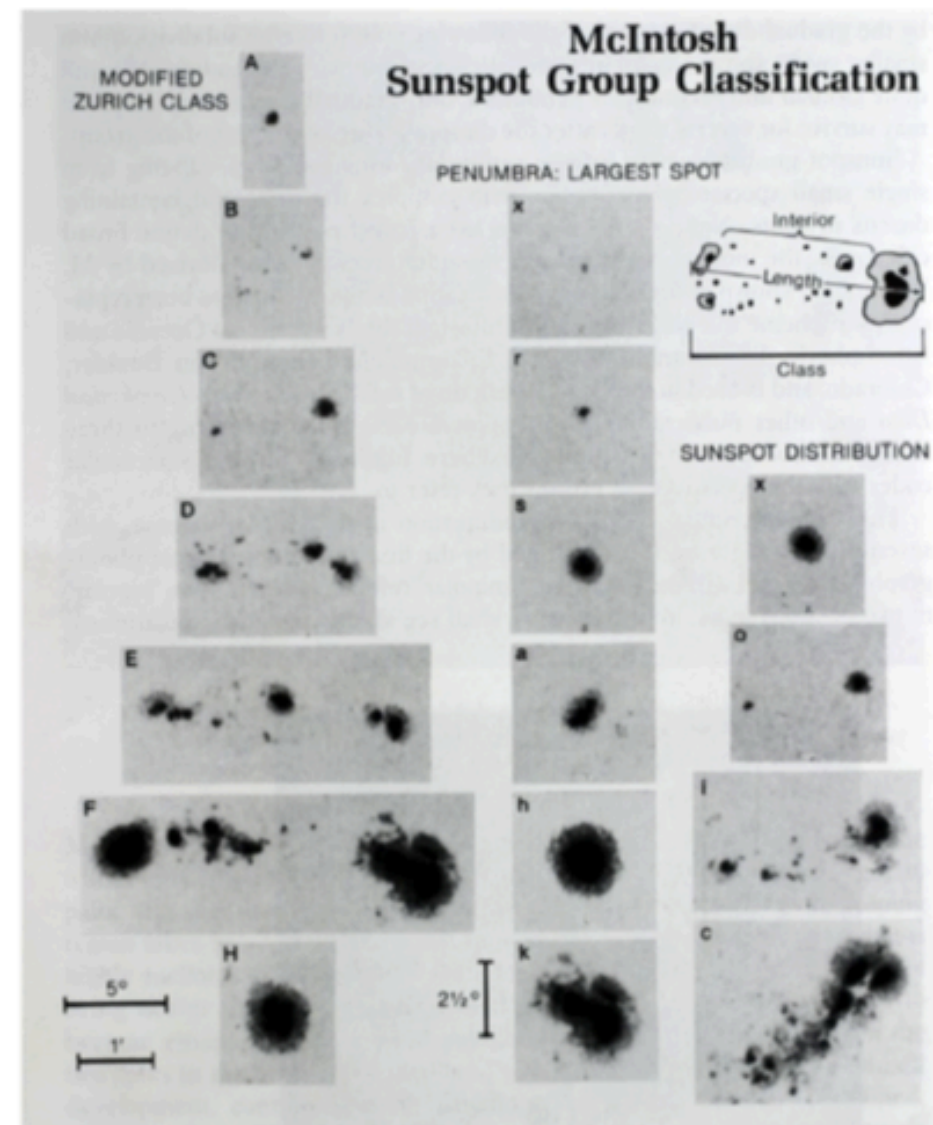
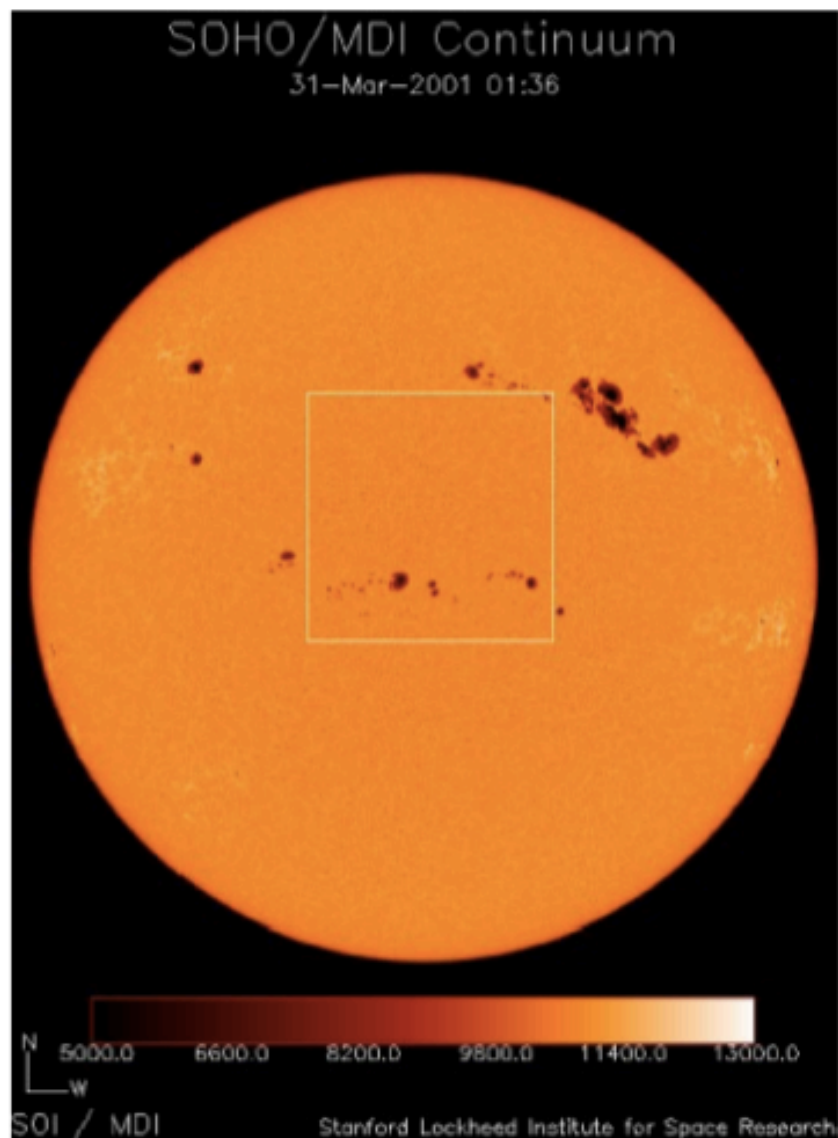
## Method:

- 1 Use clustering algorithm to approximate intermediate concepts;
- 2 Use rule based algorithm (RSES system) to approximate the target concept;

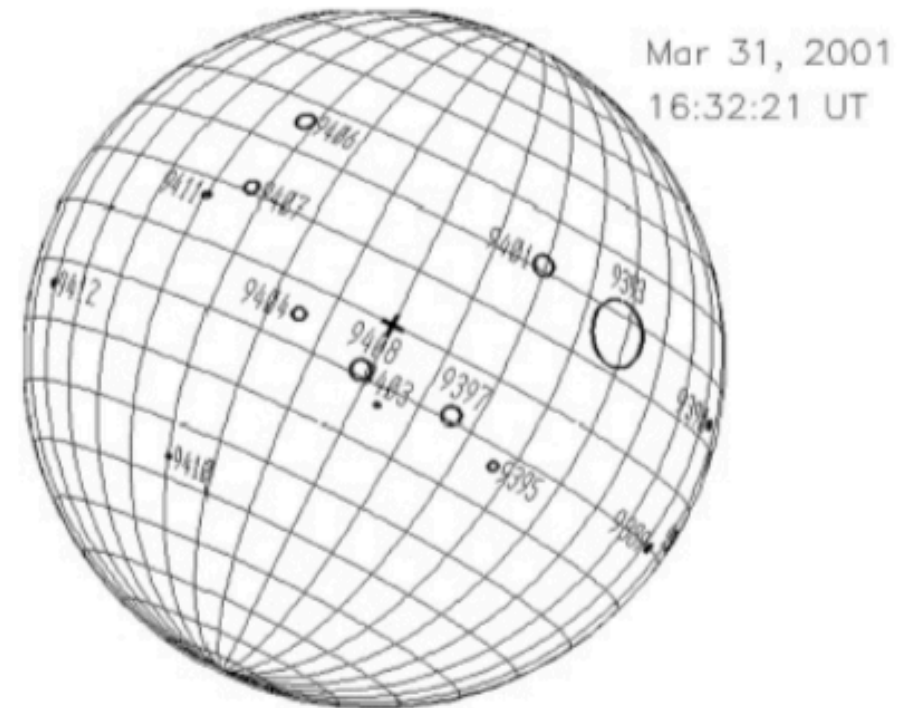
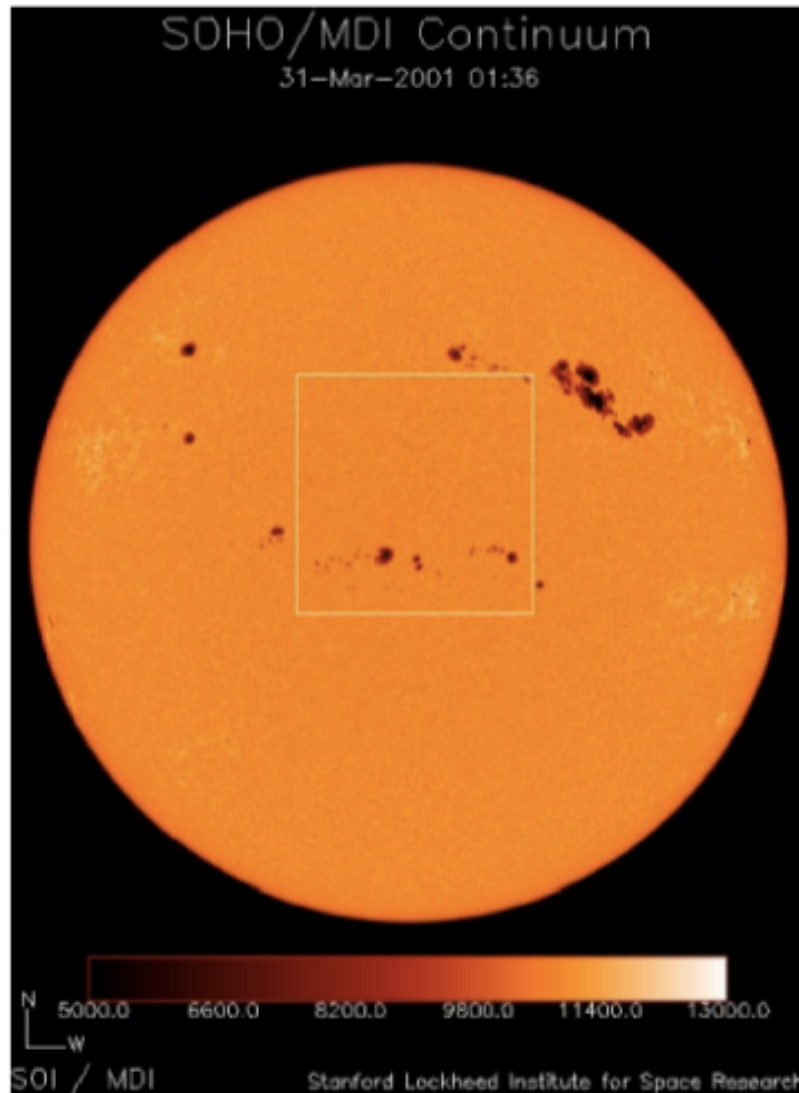
## Results: (60% – training, 40% – testing )

	original attributes only	using intermediate concepts
Accuracy	83.4	99.9%
Coverage	85.3%	100%
Nr of rules	634	42 (for the target concept) 92 (for intermediate concepts)

# Sunspot recognition and classification



# Sunspot recognition and classification



Joint USAF/NOAA Solar Region Summary (MAR 30, 2001 24:00:00 UT)

NUMR	LOCATI	LO	AREA	Mcl	LL	NN	MAG	TYPE
9387	N08W92	216	0060	Hex	02	01	Alpha	
9389	S10W63	187	0050	Bxo	08	08	Beta	
9390	N13W65	189	0050	Hex	02	01	Alpha	
9393	N17W30	154	2240	Fko	19	63	Beta-Gamma-Delta	
9395	S13W17	141	0050	Hex	02	02	Alpha	
9396	S08W65	209	0140	Dao	10	05	Beta	
9397	S09W06	130	0180	Eoo	15	23	Beta-Gamma	
9401	N21W11	135	0230	Eki	13	37	Beta-Gamma	
9403	S13E06	118	0010	Bxo	05	02	Beta	
9404	S05E23	101	0080	Cao	04	07	Beta	
9406	N26E41	083	0170	Hex	03	01	Alpha	

# Differential Calculus to Function Approximation

- **ill-defined data**: limited number of objects and large number of attributes;
- prediction of a **real decision variable** based on nominal attributes;
- the need for the knowledge about the **real mechanisms behind the data**;

No.	Combination	B-1	1-4	4-6	6-E	PB	PE	Binding affinity
1	A2B2C2D2a2b2	1	1	1	1	1	1	4.52526247
2	A1B2C1D1a2b2	-1	1	-1	-1	1	1	4.818066119
3	A1B2C2D1a2b2	-1	1	1	-1	1	1	5.036009902
...	...	...	...	...	...	...	...	
...	...	...	...	...	...	...	...	
39	A1B1C1D1a1b1	-1	-1	-1	-1	-1	-1	8.963821581
40	A1B1C1D1a2b1	-1	-1	-1	-1	1	-1	8.998482244

# Discrete Differential Calculus

## Input

### 1. A decision table

$\$$	$a_1$	$a_2$	...	$dec$
$u_1$	1	-1	...	4.23
$u_2$	1	1	...	4.31
...	...	...	...	...
$u_n$	-1	1	...	8.92

### 2. Domain knowledge

## First level

- Create comparing table

	$a_1$	$a_2$	...	change
$u_1, u_2$	$1 \rightarrow 1$	$-1 \rightarrow 1$	...	$\nearrow$
$u_1, u_3$	...	...	...	$\searrow$
...	...	...	...	...

- Learn the preference relation, i.e., decision rules of form

$$a_2 : -1 \rightarrow 1 \wedge a_6 = 1 \dots \implies change = \searrow$$

## Second level

- Ranking prediction;
- Decision value prediction;
- Experiment design.

# Semantic evaluation of clustering algorithm.



Which partition is better?

# External evaluation methods

Doc.	Soft Cluster			Expert Tag				
	$C_1$	$C_2$	$C_3$	Cosmonaut	astronaut	moon	car	truck
$d_1$	1			1		1	1	
$d_2$	1				1	1		
$d_3$	1	1		1				
$d_4$		1	1				1	1
$d_5$		1	1				1	
$d_6$			1					1

## Rand index:

Pairs of documents		Same cluster?	
		Yes	No
Same expert tag?	Yes	$a$	$b$
	No	$c$	$d$

$$Rand\ Index = \frac{a + d}{a + b + c + d}.$$

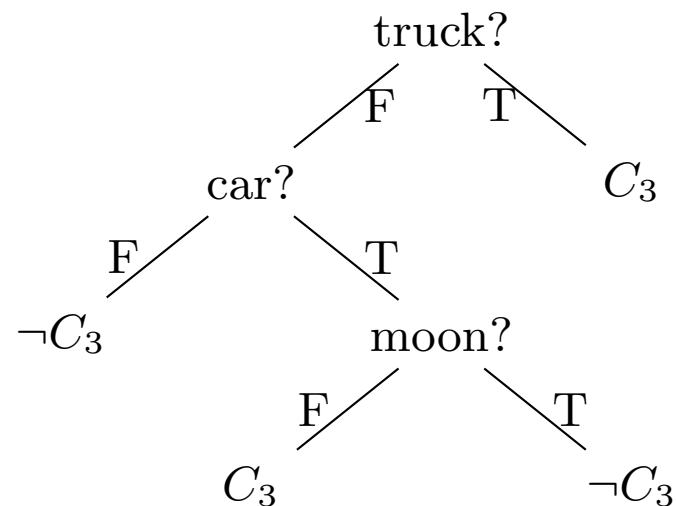
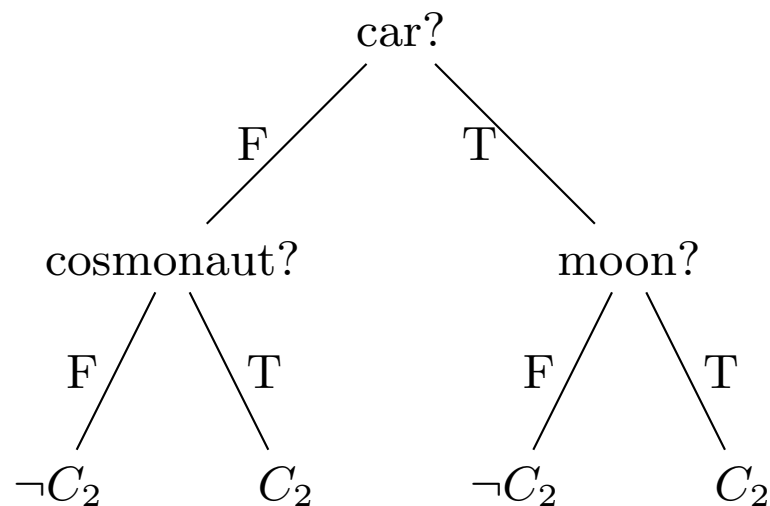
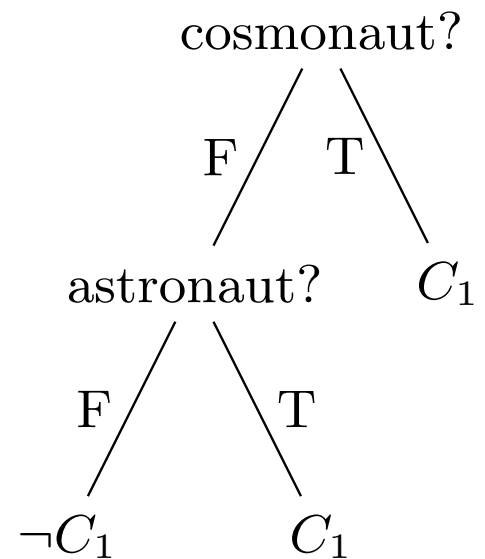
## MMI (Maximum Mutual Information):

	$C_1$	$C_2$	$C_3$
Cosmonaut	0.139	0.083	0
astronaut	0.083	0	0
moon	0.139	0	0
car	0.056	0.125	0.125
truck	0	0.042	0.208

$$MMI(X, Y) = \sum_x \sum_y p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right).$$

# Semantic Explorative Evaluation

Doc.	Expert Tag					decision $C_1$
	Cosm.	astron.	moon	car	truck	
$d_1$	1		1	1		1
$d_2$		1	1			1
$d_3$	1					1
$d_4$				1	1	
$d_5$				1		
$d_6$					1	



# Learning Ontology

A background image showing a portion of a laptop on the left and a notebook with handwritten notes at the bottom left. The laptop screen displays a Windows-style interface with various icons. The notebook has blue ink handwriting. The main content area is a solid brown bar with white text.

# Ontology learning

- Ontology Learning from Text:
- Linked Data Mining
- Concept Learning in Description Logics and OWL
- Crowdsourcing



# Challenges in ontology learning

- **Heterogeneity:** neither the integration of methods nor the homogenization of data has attracted high attention of ML community
- **Uncertainty:** Low-quality or unstructured data can lead to results that are less likely to be correct.
- **Reasoning:** ontology learning approaches are not capable of generating consistent (and coherent) ontologies
- **Scalability:** Extracting knowledge from the growing amounts of data on the web – un-structured, textual data on the one hand and structured data such as databases, linked data or ontologies on the other hand – requires scalable and efficient approaches
- **Quality:** Formal correctness, completeness and consistency are only a few of many possible criteria for judging the quality of an ontology
- **Interactivity:** The lesser the extent to which humans are involved in a semi-automatic ontology generation process, the lower the quality we can expect.

A vertical strip on the left side of the slide shows a close-up of a laptop screen displaying a Windows-style desktop with various icons, a keyboard, and a notebook with handwritten notes.

# CONCLUSIONS

- No free lunch theorem =>  
a need of knowledge modeling and involving in the learning process
- Layered learning = decomposition + synthesis of results
- Lack of a “back propagation” mechanism
- ML techniques are efficient in Knowledge Acquisition