

# KNOWLEDGE REPRESENTATION

13

## 13 KNOWLEDGE REPRESENTATION\*

---

13.1 Knowledge

13.2 Ontology

13.3 Production systems

13.4 Qualitative physics

13.5 Structured descriptions

13.6 Frame and semantic networks

13.7 Semantic web

13.8 Knowledge graph

13.9 Knowledge embedding

13.10 Change

13.11 Explanation and diagnosis

13.12 Mental states

13.13 Commonsense

# KR and AI

---

## Knowledge Representation (KR)

- Symbols standing for things in the world (symbolism)  
e.g., encoding of propositions believed (by some agent)
- Embedding spaces for things in the world (connectionism)  
e.g., encoding of vectors of preserving the inherent structure of the knowledge

Much of AI involves building systems that are knowledge-based (symbolism)

reasoning over explicitly represented knowledge

- language understanding, planning, diagnosis, etc.

Much of ML (machine learning) involves building systems that are data-driven (connectionism)

Some AI, to a certain extent of mix

- game-playing, vision, speech, motor control, etc.

How much of intelligent behavior is knowledge-based??

# Knowledge

---

Knowledge: things in the world

- *Language*, e.g., FOL
- *Representation*, symbols standing for things  
e.g., declarative knowledge
- *Reasoning*, e.g., proofs and model checking

In philosophy, the study of knowledge is called **epistemology**

Plato defined knowledge as “justified true belie” (ongoing debate)

**Belief**: not necessarily true and/or held for appropriate reasons

# Declarative vs. procedural knowledge

---

Declarative (descriptive): “know-what”

knowledge expresses in declarative sentences

Procedural knowledge: “know-how”

E.g., knowledge about the computation of matrixes vs.  
procedure of computing matrixes (in a programming language)

Theorem proving (like resolution) is a general domain-independent method of reasoning

does not require the user to know how knowledge will be used

Want to communicate to theorem-proving procedure some guidance based on properties of the domain

- perhaps a specific method to use
- perhaps merely a method to avoid

# Knowledge base

---

Separation between the knowledge base and reasoning procedure should be maintained

**Knowledge base (KB)**: to store structured and unstructured inform

- needed to know facts about the world
- to distinguish from database
  - – not just tables with numbers and strings
- scaled up with Internet documents/hypertext/multimedia
  - – known as Web Content Management

A good KB should be expressive, concise, unambiguous, context-insensitive, effective, clear and correct

**Knowledge engineering** (expert systems, knowledge-based systems): the process of building a knowledge base

# Knowledge engineering vs. software engineering

The knowledge engineer or agent usually interviews the real experts or environments to become educated about the domain and to elicit required knowledge in a process called **knowledge acquisition**

## Knowledge engineering

1. Choosing a logic
2. Building a knowledge base
3. Implementing the proof theory
4. Inferring new facts

## Software engineering (Programming)

- Choosing a programming language
- Writing a program
- Choosing or writing a compiler
- Running a program

Should be less work

# Ontology

---

**Ontology:** a vocabulary for the domain knowledge

**Ontological engineering:** representing various ontology

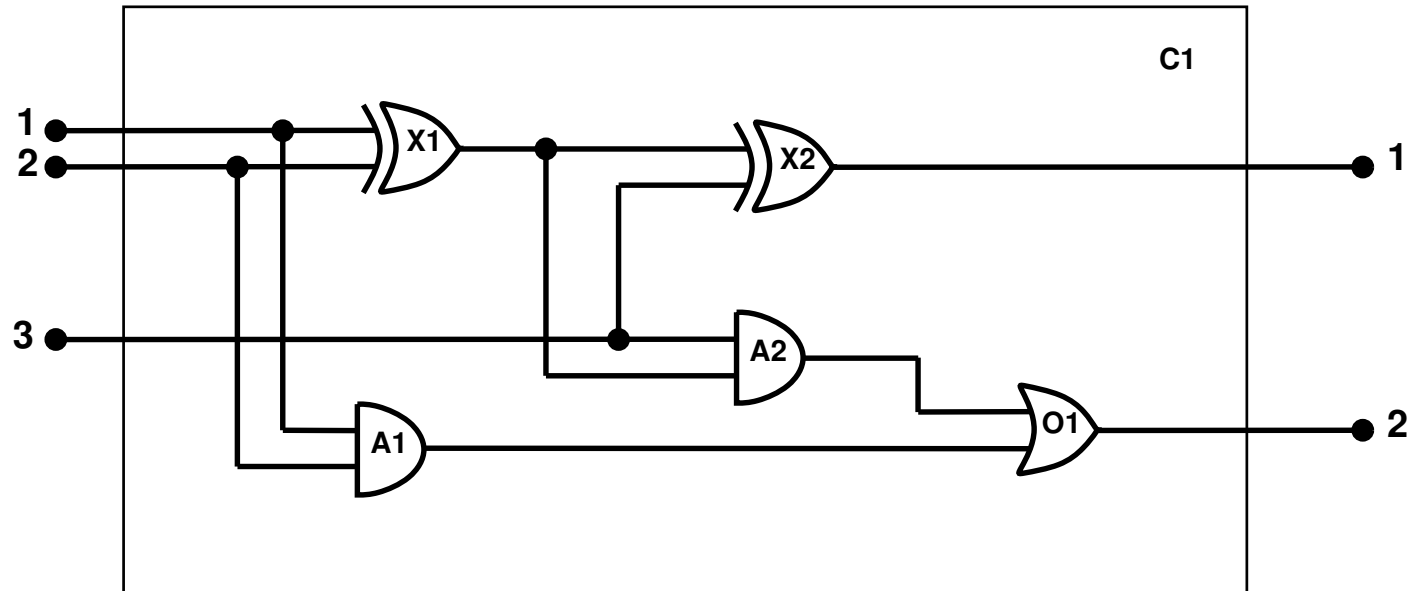
The five-step methodology

1. Decide what to talk about
2. Decide on a vocabulary of predicates, functions and constants
3. Encode general knowledge about the domain
4. Encode a description of the specific problem instance
5. Pose queries to the inference procedure and get answers



# The electronic circuits domain

---



# Ontological engineering

---

1. Decide what to talk about  
e.g., gates AND, OR, XOR and NOT
2. Decide on a vocabulary of predicates, functions and constants  
e.g.,  $Out(1, X_1)$
3. Encode general knowledge about the domain  
e.g.,  $\forall t_1 t_2 Connected(t_1, t_2) \Rightarrow Signal(t_1) = Signal(t_2)$
4. Encode a description of the specific problem instance  
e.g.,  $Type(X_1) = XOR$
5. Pose queries to the inference procedure and get answers  
e.g., what combinations of inputs would cause the first output of  $C_1$  (the sum bit) to be off? The answer ...

# General ontology

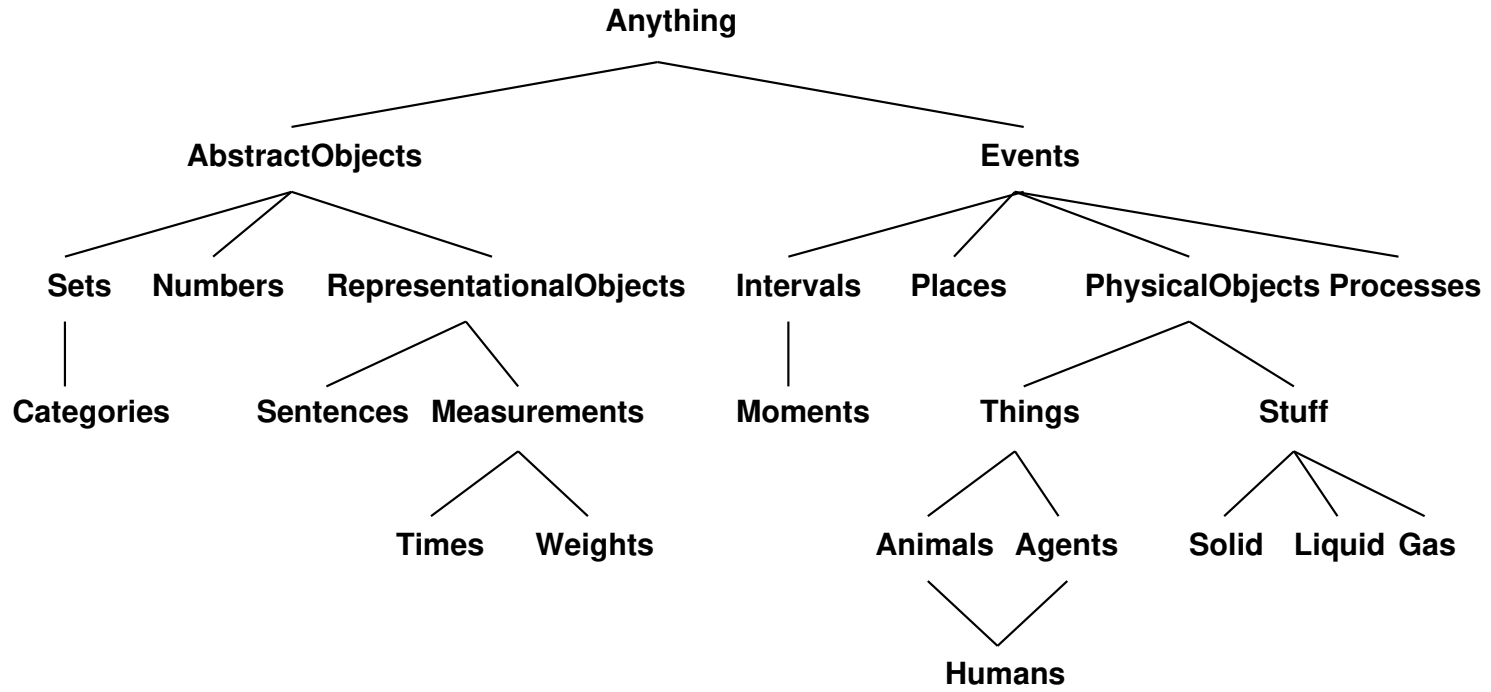
---

A general-purpose ontology has advantages over a special-purpose one

- Categories
- Measures
- Composite objects
- Time, Space, and Change
- Events and Processes
- Physical objects
- Substances
- Mental objects and belief

# The world ontology

---



Hard to build a real-world ontology

## Some KBs or ontologies

---

There are some routes for building very larger KBs or ontologies:

- CYC: creating the ontology and writing axioms from cyclopedia (1990)

- DBpedia: importing categories, attributes and values from Wikipedia (2007)

- TextRunner: building by reading a large corpus of Web pages (2008)

- OpenMind: building by volunteers who proposed facts and commonsense knowledge in English (2002)

- Knowledge Graph (KG, previous Freebase): building by Google and holding more than 70 billion facts (2012)

- Wikidata, Linking Open Data (LOD), YAGO etc.

Other data sources (also known Deep Web): MusicBrainz, DrugBank etc.

# Production systems

---

Production (rule-based) systems:

- *working memory*: a knowledge base
- *rule memory*: a set of inference rules with form

$$p_1 \wedge \cdots \wedge p_n \Rightarrow act_1 \wedge \cdots \wedge act_m$$

where  $p_i$  are literals, and  $act_j$  are actions to take when the  $p_i$  are all satisfied – forward chaining

- *match phase*: in each cycle, the system computes the subset of rules whose left-hand side is satisfied by the current contents of the working memory

- *conflict resolution phase*: the system decides which of the rules should be executed

- *act phase*: in each cycle, the system executes the action(s) in the chosen rule(s)

# Production systems

---

Inefficient forward chaining unification match algorithm:

E.g., If there are  $w = 100$  elements in working memory and  $r = 200$  rules each with  $n = 5$  elements in the left-hand side, and solving a problem requires  $c = 1000$  cycles, then the naive match algorithm must perform  $wrnc = 10^8$  unifications

Rete algorithm of OPS5

E.g., rule memory

$$A(x) \wedge B(x) \wedge C(y) \Rightarrow addD(x)$$

$$A(x) \wedge B(x) \wedge D(y) \Rightarrow addE(x)$$

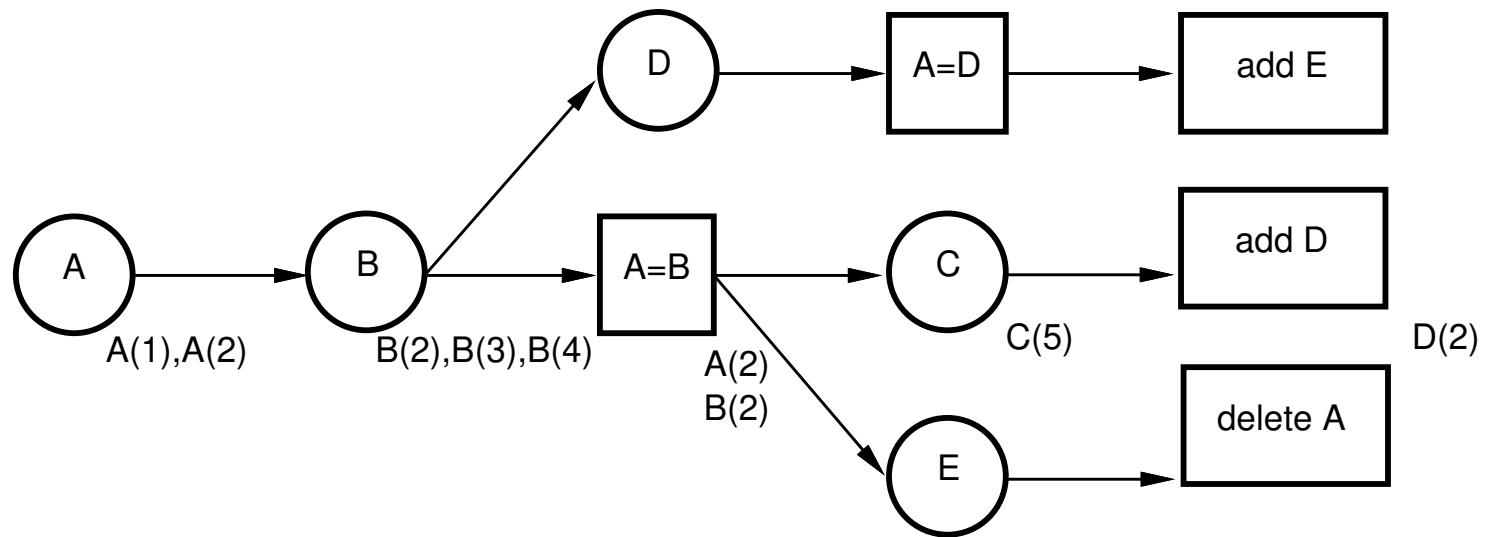
$$A(x) \wedge B(x) \wedge E(y) \Rightarrow addDeleteA(x)$$

and working memory

$$\{A(1), A(2), B(2), B(3), B(4), C(5)\}$$

# Production systems

---





# Production systems

---

Conflict resolution phase: some control strategy

- No duplication
- Recency
- Specificity
- Operation priority

OPS5

1. discard rule instances that have already been used
2. order remaining instances in terms of recency of working memory matching 1st condition (and then of 2nd condition, etc.)
3. if still no single rule, order rules by the number of conditions
4. select arbitrarily among those remaining

Production systems are essentially programming languages

# Expert systems

---

So-called **expert systems** are production systems

MYCIN (Stanford): aided physicians in treating bacterial infections  
– approximately 500 rules for recognizing about 100 causes of infection

E.g.

*IF*

the type of  $x$  is primary bacteremia

the suspected entry point of  $x$  is the gastrointestinal tract

the site of the culture of  $x$  is one of the sterile sites

*THEN*

there is evidence that  $x$  is bacteroides

– certainty factors: numbers from  $[0, 1]$  attached to conclusions to rank order

Recently, IBM Watson Health systems

# Qualitative physics

---

Qualitative physics (qualitative reasoning) concerns specifically with constructing a logical, nonnumeric theory of physical objects and processes

**Measure:** the values of the properties that we assign for objects

$$Price(tomato) = \$(0.3)$$

$$\forall d. d \in Days \Rightarrow Duration(d) = Hours(24)$$

# Composite object

---

Composite object: any object that has parts

E.g.,  $PartOf(taiwan, china)$

Schema (script): structure description

$\forall a Biped(a) \Rightarrow$

$\exists l_1 l_2 b Leg(l_1) \wedge Leg(l_2) \wedge Body(b) \wedge$

$Attached(l_1, b) \wedge Attached(l_2, b) \wedge$

$l_1 \neq l_2 \wedge \forall l_3 Leg(l_3) \wedge PartOf(l_3, a) \Rightarrow (l_3 = l_1 \vee l_3 = l_2)$

Various common knowledge of physics processes

say, water cycle

# Structured descriptions

---

In FOL, all categories and properties of objects are represented by atomic predicates

- correspond to simple nouns, e.g., *Person*
- seem to be more like noun phrases, e.g., *MarriedPerson*
- – have an internal structure and connections to other predicates  
e.g., a married person must be a person

These connections hold by definition, not by virtue of the facts we believe about the world

Need a way to break apart a predicate to see how it is formed from other predicates

# Tradeoff between expressiveness and tractability

Reasoning procedures required for more expressive languages  
may not work very well in practice

Tradeoff: expressiveness vs. tractability

E.g., Description Logics (see later)

– limited languages: between propositional language and first-order language with decidability

– vivid reasoning: easy to implement

# Categories

---

**Category:** include as members all objects having certain properties

E.g., An object (penguin) is a member of a category (birds)

*penguin*  $\in$  *birds*

Subclass relations organize categories into a **taxonomy (hierarchy)**

E.g., a category is a subclass of another category

*tomatoes*  $\in$  *fruit*

**Inheritance:** the individual inherits the property of the category from their membership

E.g.,  $Child(x, y) \wedge Familyname(john, y) \rightarrow Familyname(john, x)$

The problem: natural kind or inheritance with exception

E.g.,  $\forall x. x \in Typical(bird) \Rightarrow Flies(x)$

# Semantic networks

---

Inheritance is the result of (path-based) transitivity reasoning over paths in a network

- inheritance in trees
- inheritance in DAGs

- if-then reasoning in graphical form
- “does  $a$  inherit from  $b$ ?” is the same as “is  $b$  in the transitive closure of  $:IS-A$  (or subsumption) from  $a$ ?”



# Inheritance reasoning

---

## Inheritance with exceptions

$$\forall rxb. Holds(r, x, b) \Leftrightarrow \\ Val(r, x, b) \vee (\exists px \in p \wedge Rel(r, p, b) \wedge \neg InterveningRel(x, p, r))$$

$$\forall xpr. InterveningRel(x, p, r) \Leftrightarrow \\ \exists i Intervening(x, i, p) \wedge \exists b' Rel(r, i, b')$$

$$\forall aip. Intervening(x, i, p) \Leftrightarrow (x \in i) \wedge (i \subset p)$$

## Multiple inheritance

- credulous accounts choose arbitrarily
- skeptical accounts are more conservative

# Description logics

---

**Description logics (DLs):** a family of logics with notations designed to describe the definitions and properties of categories

**Subsumption:** checking if one category is a subset of another based on their definitions

**Classification:** checking if an object belongs to a category

DLs focus on the tractability of inference and serve as theoretical foundation for ontology

## Concepts, roles, constants

---

In DL, there are sentences that will be true or false (as in FOL)

In addition, there are three sorts of expressions that act like nouns and noun phrases in natural language

- **concepts** are like category nouns, e.g., *GraduateStudent*
- **roles** are like relational nouns, e.g., *:AreaOfStudy*  
(note “:” at start)
- **constants** are like proper nouns, e.g., *johnSmith*

These correspond to unary predicates, binary predicates and constants respectively in FOL

However, unlike in FOL, concepts need not be atomic and can have semantic relationships with each other

# A description logic: language\*

---

Symbols of *DL*:

- atomic concepts, e.g., *GraduateStudent*
- roles: (all are atomic), e.g., *AreaOfStudy*
- constants

Four types of logical symbols:

- punctuation:  $[, ], (, )$
- positive integers:  $1, 2, 3, \dots$
- concept-forming operators: *ALL, EXISTS, FILLS, AND*
- connectives:  $\sqsubseteq, \doteq, \text{and } \rightarrow$

## A description logic: syntax\*

---

The set of concepts is the least set satisfying

- Every atomic concept is a concept
- If  $r$  is a role and  $d$  is a concept, then  $[ALL\ r\ d]$  is a concept
- If  $r$  is a role and  $n$  is an integer, then  $[EXISTS\ n\ r]$  is a concept
- If  $r$  is a role and  $c$  is a constant, then  $[FILLS\ r\ c]$  is a concept
- If  $d_1, \dots, d_k$  are concepts, then so is  $[AND\ d_1, \dots, d_k]$

Three types of sentences

- If  $d$  and  $e$  are concepts, then  $(d \sqsubseteq e)$  is a sentence
- if  $d$  and  $e$  are concepts, then  $(d \doteq e)$  is a sentence
- If  $d$  is a concept and  $c$  is a constant, then  $(c \rightarrow d)$  is a sentence

## A description logic: semantics\*

---

Constants stand for individuals, concepts for sets of individuals, and roles for binary relations

The meaning of a complex concept is derived from the meaning of its parts the same way a noun phrases is

- $[EXISTS\ n\ r]$  describes those individuals that stand in relation  $r$  to at least  $n$  other individuals
- $[FILLS\ r\ c]$  describes those individuals that stand in the relation  $r$  to the individual denoted by  $c$
- $[ALL\ r\ d]$  describes those individuals that stand in relation  $r$  only to individuals that are described by  $d$
- $[AND\ d_1, \dots, d_k]$  describes those individuals that are described by all of the  $d_i$

Formal semantics can be defined

## A description logic: example\*

---

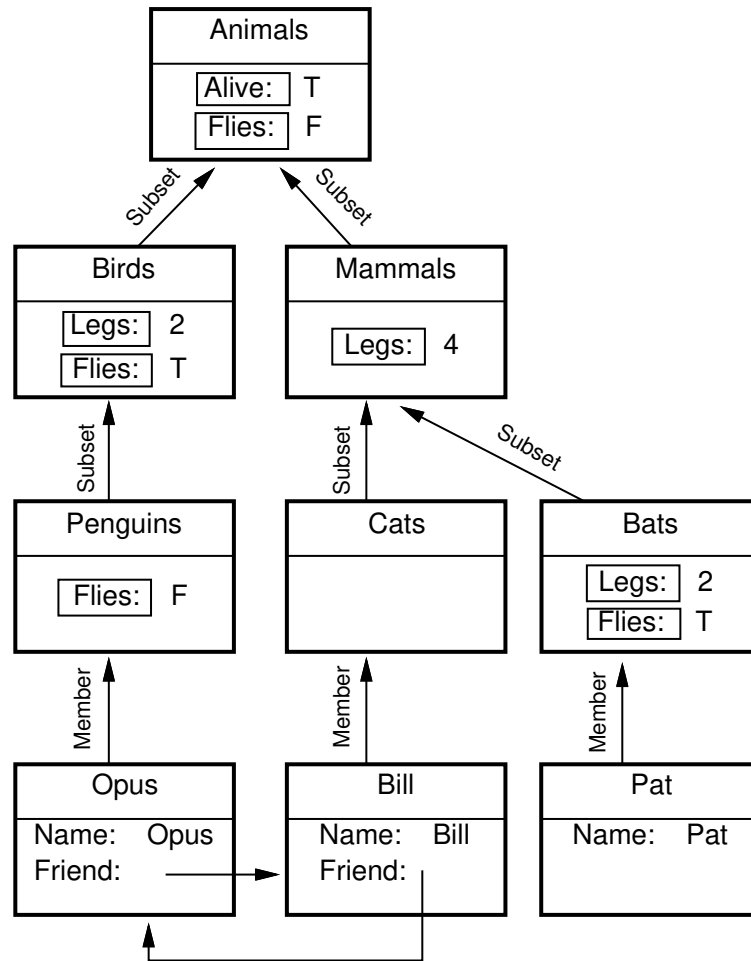
*[AND Company*  
    *[EXISTS 7 : Director]*  
    *[ALL : Manager [AND Woman*  
        *[FILLS : Degree PhD]]]*  
    *[FILLS : MinSalary \$24.00/hour]]]*

“a company with at least 7 directors, whose managers are all women with PhDs, and whose min salary is \$24/hr”

A DL knowledge base is a set of DL sentences serving mainly to

- give names to definitions
- give names to partial definitions
- assert properties of individuals

# Frame



(a) A frame-based knowledge base

Rel(Alive,Animals,T)  
Rel(Flies,Animals,F)

Birds  $\subset$  Animals  
Mammals  $\subset$  Animals

Rel(Flies,Birds,T)  
Rel(Legs,Birds,2)  
Rel(Legs,Mammals,4)

Penguins  $\subset$  Birds  
Cats  $\subset$  Mammals  
Bats  $\subset$  Mammals  
Rel(Flies,Penguins,F)  
Rel(Legs,Bats,2)  
Rel(Flies,Bats,T)

Opus  $\in$  Penguins  
Bill  $\in$  Cats  
Pat  $\in$  Bats  
Name(Opus,"Opus")  
Name(Bill,"Bill")  
Friend(Opus,Bill)  
Friend(Bill,Opus)  
Name(Pat,"Pat")

(b) Translation into first-order logic



# Frame vs. logic

Link Type	Semantics	Example
$A \xrightarrow{\text{Subset}} B$	$A \subset B$	$Cats \subset Mammals$
$A \xrightarrow{\text{Member}} B$	$A \in B$	$Bill \in Cats$
$A \xrightarrow{R} B$	$R(A, B)$	$Bill \xrightarrow{\text{Age}} 12$
$A \xrightarrow{\boxed{R}} B$	$\forall x x \in A \Rightarrow R(x, B)$	$Birds \xrightarrow{\boxed{\text{Legs}}} 2$
$A \xrightarrow{\boxed{\boxed{R}}} B$	$\forall x \exists y x \in A \Rightarrow y \in B \wedge R(x, y)$	$Birds \xrightarrow{\boxed{\boxed{\text{Parent}}}} Birds$

Frame and semantic networks can be formalized by FOL, and hence can be thought of applications of FOL

# Procedural knowledge

---

Organizing procedural knowledge  
knowing facts by executing code

Object-oriented (OO) representation

– with enough procedures/sentences in a KB, need to organize them

– in terms of **objects**

– – clustering procedures for determining properties, identifying parts, interacting with parts, as well as constraints between parts, all of objects

to make some things easier to find

# Basic frame language

---

Frames are object structures

- individual frames: represent a single object, e.g., *person*
- generic frames: represent categories of objects, e.g., *students*

An individual frame is a named list of buckets called **slots**. What goes in the bucket is called a **filler** of the slot

(*frame-name*  
    < *slot-name1* *filler1* >  
    < *slot-name2* *filler2* > ...)

where frame names and slot names are atomic, and fillers are either numbers, strings or the names of other individual frames

Notation: **attribute-value pair** (AVP)

individual frames: *birds*

slot names: *:Fly* (note “:” at start)

generic frames: *Animals*

## IS-A and inheritance

---

Individual frames have a special slot: *INSTANCE-OF* (*IS-A*)  
whose filler is the name of a generic frame

(*birds*

<: *INSTANCE-OF* *Animals* >

<: *Fly* *T* > ...)

Similarly, generic frames have a slot: *IS-A*  
whose filler is the name of another generic frame

Slots in generic frames can have **attached procedures**

- computing a filler (when no slot filler is given)
- propagating constraints (when a slot filler is given)

# IS-A and inheritance

---

Specialization relationships imply that procedures and fillers of a more general frame are applicable to more specific frame  $\Rightarrow$  inheritance

Basic (local) reasoning with frames

- user instantiates a frame, i.e., declares that an object or situation exists

- slot fillers are inherited where possible

- causing more frames to be instantiated and slots to be filled

Object-oriented programming

specifying problems with frames is a style of **programming** rather than declarative object-oriented modeling of the world

# Semantic Web

---

**Semantic Web** is the next generation of Web led by W3C (World Wide Web Consortium, <http://www.w3c.org>) that makes the Web pages understandable to machine

- RDF (Resource Description Framework) as underlying meta-data representational language is a language of categories (ontology)
- The “semantics” of Web (data and facts) is realized by ontology, OWL (Web Ontology Language) is an ontology representation language
- Various Knowledge Graphs are implementations of KB
- Description logics are theoretical foundation of ontology and the standards of ontology language
- The spirit of the semantic Web came from AI and can be viewed as an application of AI (so-called Internet + AI)

# Resource description framework

---

**RDF:** W3C specifications as a metadata data model

- based on description logics
- syntax in **XML** (eXtensible Markup Language)

**Triples:** *subject-predicate-object*, or

**(h,r,t)** (*head-relation-tail*), entity-relationship (**ER** model)

- *subject, object* (entities): the (web) resources
- *predicate*: the relationship between the subject and the object

**RDF graph:** a collection of RDF triples represents a labeled, directed multi-graph

**SPARQL:** query language for RDF graphs

- an SQL-like language

# Web ontology language

---

OWL: a family of ontology representation languages

- based on description logics
- RDF/OWL syntax in **XML** (eXtensible Markup Language)

OWL extends RDF Schema

- Class equivalent  
Property  
sameIndividualAs  
...
- RDFS  
subClassOf  
resource  
ID  
...



## Example

---

Define the terms “Camera” and “SLR”, state that SLRs are a type of Camera

```
< owl : Classrdf : ID = “Camera” / >
```

```
< owl : Classrdf : ID = “SLR” >
```

```
< rdfs : subclassOfrdf : resource = “#Camera” / >
```

```
< /owl : Class >
```

# Knowledge graph

---

Knowledge graph (KG, or *linked data*, multi-relational data): representing entities and relations

Triples  $(h,r,t)$ : *head-relation-tail*

– based on RDF and description logics

E.g. KG services to enhance search engine's results with information gathered from a variety of sources

– Google, knowledge panel

# Knowledge embedding

---

Knowledge embedding representation (KER): representing knowledge into low-dimensional (continuous) vector spaces (so-called **embedding**), so as to simplify computations preserving the inherent structure of the knowledge

E.g. fit any dataset with a (continuous, differentiable) scalar function with a single real-valued parameter

$$f_{\alpha}(x) = \sin^2(2^{x\tau} \arcsin \sqrt{\alpha})$$

any dataset can be viewed as a list of numerical values  $\mathcal{X} = [x_0, \dots, x_n]$  describing the data content regardless of the underlying modality (time-series, images, sound etc.)

– say, animal shapes obtained with the different values of  $\alpha$

Ref. Boué L, Real numbers, data science and chaos: How to fit any dataset with a single parameter, arXiv, 2019.

# Knowledge graph embedding

---

**KG embedding (KGE)**: embedding entities and relations of a KG into low-dimensional vector spaces, so as to simplify computations preserving the inherent structure of the KG

Machine learning

- Representing entities and relations.  
Entities are represented as vectors (deterministic points in the vector space)
- Defining a **scoring function**  
Defined each fact to measure its plausibility (facts observed in the KG have higher scores).
- Learning entity and relation representations (i.e., embeddings)  
Optimization problem that maximizes the total plausibility of observed facts

## Distance-based embedding

---

The relationships are represented as **translations** in the embedding space: if  $(h, r, t)$  holds, then the embedding of  $t$  should be close to the embedding of  $h$  plus some vector that depends on  $l$

– i.e., the functional relation induced by the  $l$ -labeled edges corresponds to a translation of the embeddings

I.e., we want that  $\vec{h} + \vec{l} \approx \vec{t}$  when  $(h, r, t)$  holds, while  $\vec{h} + \vec{l}$  should be far away from  $\vec{t}$  otherwise

There are various ways to embed knowledge into mathematical structures

# Change

---

Time:

E.g.,  $At(\text{evening}, \text{sleep})$

Event:

E.g.,  $\text{worldWarII}, \text{SubEvent}(\text{battleOfBritain}, \text{worldWarII})$

An event that includes as subevents all events occurring in a given time period is called **interval**

Space:

E.g.,  $In(\text{beijing}, \text{china})$

$\forall x.l. \text{Location}(x) = l \Leftrightarrow$

$\text{At}(x, l) \wedge \forall l_1 \text{At}(x, l_1) \Rightarrow In(l, l_1)$

Process: liquid event

E.g.,  $T(\text{working}(\text{teacher}), \text{todayLessonHours})$

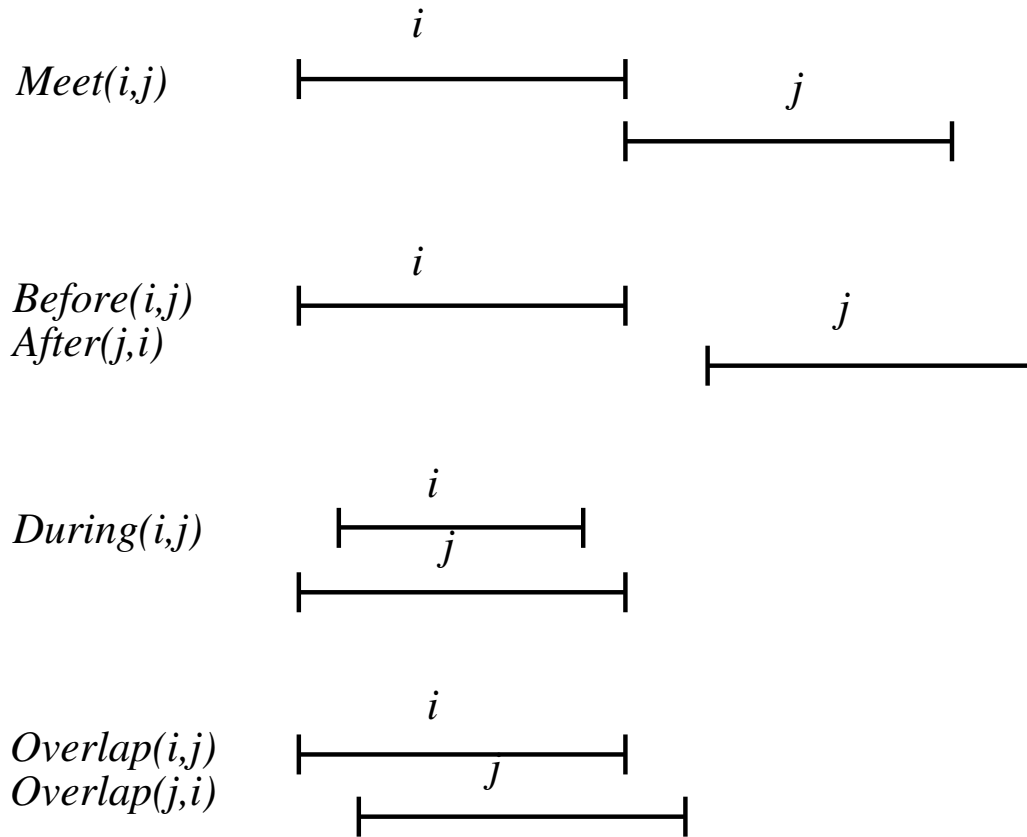
$T(c, i)$  means that some event of type  $c$  occurred over exactly the interval  $i$

# Reasoning about change

---

## Time interval

E.g.,  $\forall ij. Meet(i, j) \Leftrightarrow Time(End(i)) = Time(Start(j))$



# Reasoning about change

---

## Action

E.g.,  $\forall xyi_0.T(\text{engaged}(x, y), i_0) \Rightarrow$   
 $\exists i_1(\text{Meet}(i_0, i_1) \vee \text{After}(i_1, i_0)) \wedge$   
 $T(\text{Marry}(x, y) \vee \text{BreakEngagement}(x, y), i_1)$

## Fluent: something that changes across situations

E.g.,  $\text{President}(USA)$   
 $T(\text{democrat}(\text{president}(USA)), AD2003)$

## Context

E.g.,  $\text{President}(USA, 45th) = \text{DonaldTrump}$



# Explanation and diagnosis

---

## Reasoning

- deduction, such as if  $\alpha$ ,  $\alpha \Rightarrow \beta$ , then  $\beta$
- induction, such as  $\alpha$ ,  $\beta$ , then  $\alpha \Rightarrow \beta$
- abduction, such as  $\beta$ ,  $\alpha \Rightarrow \beta$ , then  $\alpha$

## Abductive reasoning

- given  $\alpha \Rightarrow \beta$ , from  $\beta$ , abduce  $\alpha$
- $\alpha$  is sufficient for  $\beta$  or  
one way for  $\beta$  to be true is for  $\alpha$  to be true

Can be used for causal reasoning: (*cause*  $\Rightarrow$  *effect*)

# Diagnosis

---

One simple diagnosis uses abductive reasoning

KB has facts about symptoms and diseases

including:  $(Disease \wedge Hedges \Rightarrow Symptoms)$

Goal: find disease(s) that best explain observed symptoms

Observe: we typically do not have knowledge of the form

$(Symptom \wedge \dots \Rightarrow Disease)$

so the reasoning is not deductive

Non-uniqueness: multiple equally good explanations

# Explanation\*

---

Given  $KB$ , and  $\beta$  to be explained, we want an  $\alpha$  s.t.

1.  $\alpha$  is sufficient to account for  $\beta$

$$KB \cup \{\alpha\} \models \beta$$

2.  $\alpha$  is not ruled out by  $KB$

$$KB \not\models \neg\alpha, \text{ i.e., } KB \cup \{\alpha\} \text{ is consistent}$$

otherwise,  $(p \wedge \neg p)$  would count as an explanation

3.  $\alpha$  is as simple as possible

4.  $\alpha$  is in the appropriate vocabulary

Call such  $\alpha$  an **explanation** of  $\beta$  w.r.t.  $KB$

The simplest explanation is the negation of a clause with a minimal set of literal

# Mental states

---

Propositional attitudes (modalities): e.g., know, believe, want, expect, etc.

Multi-agents: e.g., an agent reasons about the mental processes of the other agents

Formalizing reasoning about mental states

- syntactic theory

- possible worlds (modal logic)

Modal operators:  $B, K$

$B(a, \psi)$  or  $B_a(\psi)$ : agent  $a$  believes that sentence  $\psi$  is true

$K(a, \psi)$  or  $K_a(\psi)$ : agent  $a$  knows that sentence  $\psi$  is true

$B(A, \psi)$ ,  $A = \{a_1, \dots, a_n\}$ : every agent of  $A$  believes that sentence  $\psi$  is true

# A formal theory of belief\*

---

Extending first-order language  $L$ :

*Belief* formulas:  $Believes(Agent, fluent)$

*Strings*:  $Flies(Clark)$  represented as  $[F, l, i, e, s, (, C, l, a, r, k, ), ]$

– referential opaque: an equal term cannot be substituted for the one (mental object) in the scope of belief, e.g., “*Clark*”  $\neq$  “*Superman*”

*Den* function: mapping a string to the object that it denotes

*Name* function: mapping an object to a string that is the name of a constant that denotes the object

E.g.,

$Den("Clark") = ManOfSteel \wedge Den("Superman") = ManOfSteel$

$Name(ManOfSteel) = K_{11}$

# Belief theory\*

---

*Inference rules*, e.g., Modus Ponens

$\forall apq. \text{LogicalAgent}(a) \wedge \text{Believes}(a, p) \wedge \text{Believes}(a, \text{Concat}(p, " \Rightarrow ", q)) \Rightarrow \text{Believes}(a, q)$

where *Concat* is a function on strings that concatenates their elements together, abbreviate  $\text{Concat}(p, " \Rightarrow ", q)$  as " $p \Rightarrow q$ "

E.g., belief rules: if a logical agent believes something, then it believes that it believes it

$\forall ap. \text{LogicalAgent}(a) \wedge \text{Believes}(a, p) \Rightarrow \text{Believes}(a, " \text{Believes}(\text{Name}(a), p) ")$

# Knowledge and belief

---

*Logical omniscience:*

$Believes(a, \phi), Believes(a, \phi \Rightarrow \psi) \models Believes(a, \psi)$

– So we need a limited rational agent

*Belief and knowledge:* knowledge is justified true belief

$\forall ap. Knows(a, p) \Leftrightarrow Believes(a, p) \wedge T(Den(p) \wedge T(Den(KB(a)) \Rightarrow Den(p)))$

*Belief and Time:*  $Believes(agent, string, interval)$

*Knowledge and action:* knowledge producing actions

# Commonsense

---

KB

$$\forall x \text{Bird}(x) \Rightarrow \text{Flies}(x)$$

$$\text{Bird}(\text{Tweety})$$

$KB \vdash \text{Flies}(\text{Tweety})??$

With exceptions

$$\forall x \text{Bird}(x) \wedge x \neq \text{Penguin} \wedge \dots \Rightarrow \text{Flies}(x)$$

$$\forall x \text{Bird}(x) \wedge \neg \text{Abnormal}(x) \Rightarrow \text{Flies}(x)$$



# Commonsense reasoning: nonmonotonicity

---

Monotonicity of FOL

if  $KB \vdash P$  then  $(KB \wedge S) \vdash P$

i.e., if  $P$  follows from KB, then it still follows when KB is augmented by  $TELL(KB, S)$

*Nonmonotonicity*:  $KB \subset KB', \exists P, KB \vdash P$  but  $KB' \not\vdash P$

**Nonmonotonic logics** are the formalization of reasoning with incomplete knowledge

A solution to the frame problem and related problems

# Commonsense reasoning: paraconsistency

---

Triviality of FOL

$$\alpha \wedge \neg\alpha \vdash \beta$$

i.e., everything follows from a single contradiction

Paraconsistency:  $\{\alpha, \neg\alpha\} \subset KB, \exists\beta, KB \not\vdash \beta$

Paraconsistent logics are the formalization of reasoning with inconsistent knowledge

# Reasoning with incomplete knowledge

---

## Closed World Assumption (CWA)

Let  $KB$  be a (finite) set of sentences (belief set),  $T(KB)$  theory of  $KB$  is  $T(KB) = \{\phi \mid KB \models \phi\}$

The CWA of  $KB$ , written as  $CWA(KB) = KB \cup KB_{asm}$ , defined as follows

1.  $\phi \in T(KB)$  iff  $KB \models \phi$ ,  $\phi$  is a sentence
2.  $\neg p \in KB_{asm}$  iff  $p \notin T(KB)$ ,  $p$  is a ground atom
3.  $\phi \in CWA(KB)$  iff  $\{KB \cup KB_{asm}\} \models \phi$

# Reasoning with incomplete knowledge

---

CWA

$$KB = \{p(A), p(A) \Rightarrow q(A), p(B)\}$$

$$T(KB) \not\models q(B), T(KB) \not\models \neg q(B)$$

$$CWA(KB) \models \neg q(B)$$

Problems

$$KB = \{p(A) \vee p(B)\}$$

$$CWA(KB) \models \neg p(A) \wedge \neg p(B)$$

# Reasoning with incomplete knowledge\*

---

Predicate Completion (COMP) (Negation-as-failure in Prolog)

$$KB = \{p(A)\} \Leftrightarrow \forall x.x = A \Rightarrow P(x)$$

I.e., “if” half of a *definition* for  $P$

$$\forall x.P(x) \Rightarrow x = A$$

I.e., the *completion* formula for  $P$

The completion of  $P$  in  $KB$ , written as  $PC(KB; P)$ , defined as follows

$$COMP(KB; P) \equiv KB \wedge (\forall x.P(x) \Rightarrow x = A)$$

$$\forall x.P(x) \Leftrightarrow x = A$$

# Reasoning with incomplete knowledge\*

---

## Circumscription (CIRC)

Idea (Occam principle): the only objects satisfying the property  $P$  are those that must, given  $KB$

Preferential semantics: minimality and minimal entailment  $\models_m$

Let  $M_1, M_2$  be two models.  $M_1$  is less (preferential) than  $M_2$ , written as  $M_1 \prec_P M_2$ , if

1.  $|M_1| = |M_2|$
2.  $|M_1|_P \subset |M_2|_P$

# Reasoning with incomplete knowledge\*

---

## CIRC

Let  $M$  be a model of  $KB$ .  $M$  is said minimal (preferential) iff there is no other models  $M'$  of  $KB$  such that  $M' \prec_P M$

Define  $KB \models_m \psi$  iff  $\psi$  is true in all minimal models of  $KB$

# Reasoning with incomplete knowledge\*

---

Example KB

$$\forall x \text{Bird}(x) \wedge \neg \text{ab}(x) \Rightarrow \text{Flies}(x)$$

$$\text{Bird}(\text{Tweety})$$

$$\text{Penguin}(\text{Tweety}) \Rightarrow \neg \text{Flies}(\text{Tweety})$$

Set  $P = \{\text{ab}, \text{Penguin}, \text{Bird}\}$ . Then

$$KB \models_m \text{Flies}(\text{Tweety})$$



# Reasoning with incomplete knowledge\*

---

Default logic (DL)

*Default rule*: an inference rule (at meta level)

$$\frac{\alpha(x):\beta(x)}{\gamma(x)}$$

E.g.,  $\frac{Bird(x):Flies(x)}{Flies(x)}$

Default theory  $KB = (W, D)$ :  $D$  is a set of default rules,  $W$  is a set of sentences

*Extension of KB??*

# Reasoning with inconsistent knowledge

---

Maximal consistent subsets (MCS)

$MCS(KB)$  is the set of maximal consistent subset of  $KB$

Reasoning with incomplete and inconsistent knowledge??

In what extent commonsense reasoning can be formalized??

# The commonsense problem

---

Recall: What is common sense??

Commonsense is not explained, but  
rely on our routines of behavior that we have learned over time  
act in situations that are sufficiently unlike the routines we have  
seen before

Common sense is critical to human-level intelligence and AI  
 $AI \approx$  Commonsense