



---

CIS-700  
Spring 2020

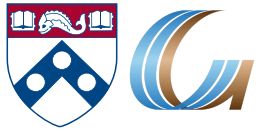
# Reasoning for Natural Language Understanding

Dan Roth

Computer and Information Science  
University of Pennsylvania

Introduction Part III: Knowledge Representation and Reasoning – Classical View

# This class



- Understand early and current work on Reasoning

- (Learn to) read critically, present, and discuss papers

- Understand some of the difficulties in NLU from the perspective of reasoning

- Conceptual and technical

- Try some new ideas

- How:

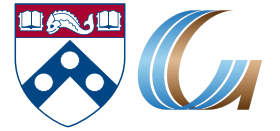
- Presenting/discussing papers
  - Probably: 2 presentations each; 4 discussants
- Writing a few critical reviews
- “Small” individual project (reproducing);
- Large project (pairs)
- Tentative details are on the web site.

- Today: discuss first project
  - Content + Timetable
- Today: release list of papers
  - Timetable

- Machine Learning
  - 519/419
  - 520
  - Other?
- NLP
  - Yoav Goldberg’s book
  - Jurafsky and Martin
  - Jacob Eisenstein
- Attendance is mandatory
- Participation is mandatory
- Time of class?
- **Expectations?**

- The classical view of reasoning:
  - Deriving conclusions from a corpus of explicitly stored information, as a mean to solve a range of problems.
- An ideal reasoning system will produce:
  - All-and-only the correct answer to every possible query
  - Produce answers that are as specific as possible
  - Be expressive enough to permit any possible fact to be stored and and query to be asked
  - Be efficient
- Probably impossible for many reasons (?)
- Most of the classical research focused on tradeoffs:
  - As correct systems become more expressive, they can become less efficient
- This was studied both in the context of logic- and of probability-based reasoning.
- Less effort was devoted to connecting things to applications where reasoning is needed
  - Representation (and Mapping) – are these realistic?
  - Formulation – is it satisfactory?

# Towards a Formulation



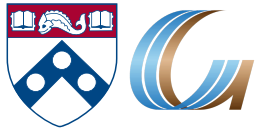
**Deduction:** Conclusion from given axioms (facts or observations)

<i>All humans are mortal.</i>	(axiom)
<i>Socrates is a human.</i>	(fact/ premise)
<i>Therefore, it follows that Socrates is mortal.</i>	(conclusion)

**Abduction:** Simple and mostly likely explanation, given observations

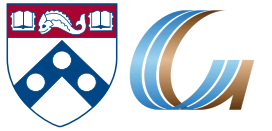
<i>All humans are mortal</i>	(theory)
<i>Socrates is mortal</i>	(observation)
<i>Therefore, Socrates must have been a human</i>	(diagnosis)

- Of these, abduction might be the most useful (?) in many situations.
- But, we need to formalize these.
- And, maybe think about the relations to Induction
- And, always ask, are these forms of reasoning sufficient?



- Propositions ( $p, q, \dots$ ). Connectives ( $\wedge, \neg, \dots$ ).
  - Implications:  $\varphi \Rightarrow x$ . Equivalences:  $\varphi \Leftrightarrow x$ .
- Reasoning *semantics* through entailment  $\models$ .
- Proof procedures  $\vdash$  to *compute* entailment.
  
- Given formulas in  $KB$  and an input  $O$ , *deduce* whether a result  $R$  is entailed ( $KB \cup O \models R$ ).
- Given formulas in  $KB$  and an input  $O$ , *abduce* an explanation  $E$  that entails  $O$  ( $KB \cup E \models O$ ).
  - The question of how to compute deduction (and abduction) is also an interesting question here.

# Non-Monotonic Reasoning

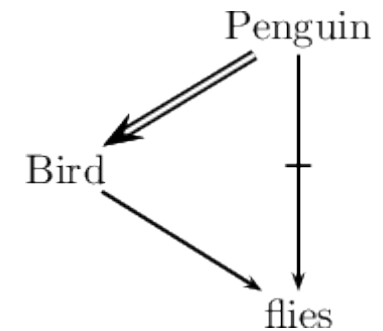


- Non-monotonicity typically viewed as property of *extending input O for fixed KB*, and having result R become “smaller”.

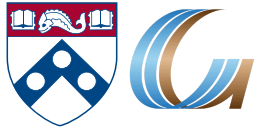
- Birds fly
- Tweedy is a bird; does Tweedy fly?
- Tweedy is a penguin

□ This is a problem to most formalisms

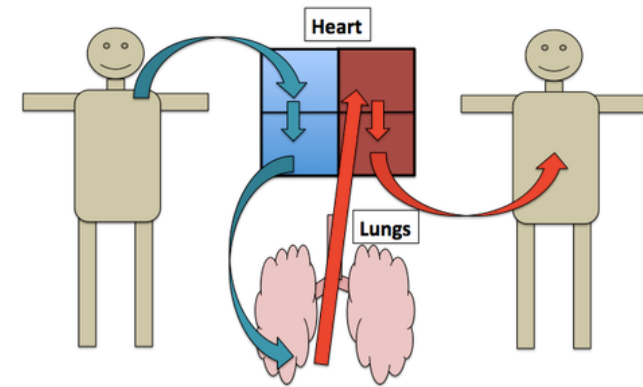
□ Involving learning in the process provides ways to address these difficulties.



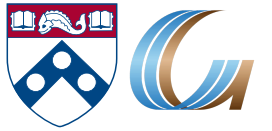
# Analogy



- The heart is a pump
- Is this an important reasoning setting?



# Quantitative Reasoning



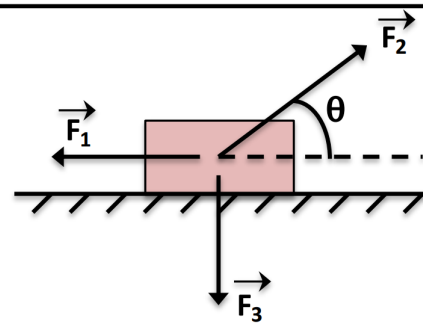
**Example:** The sum of two numbers is 111. One of the numbers is consecutive to the other number. Find the two numbers.

**Example:** Bill's father's uncle is twice as old as Bill's father. 2 years from now Bill's father will be 3 times as old as Bill. The sum of their ages is 92. Find Bill's age.

**Example:** The distance between New York to Los Angeles is 3000 miles. If the average speed of a jet plane is 600 miles per hour find the time it takes to travel from New York to Los Angeles by jet.

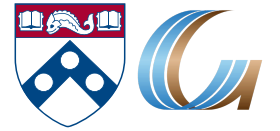
**Example:** Ram Emanuel's campaign contributions total that of all his competitors together.

The figure shows three forces applied to a trunk that moves leftward by 3.00 m over a frictionless floor. The force magnitudes are  $F_1 = 5.00\text{N}$ ,  $F_2 = 9.00\text{N}$ , and  $F_3 = 3.00\text{N}$ , and the indicated angle is  $\theta = 60.0^\circ$ . During the displacement, what is the net work done on the trunk by the three forces?

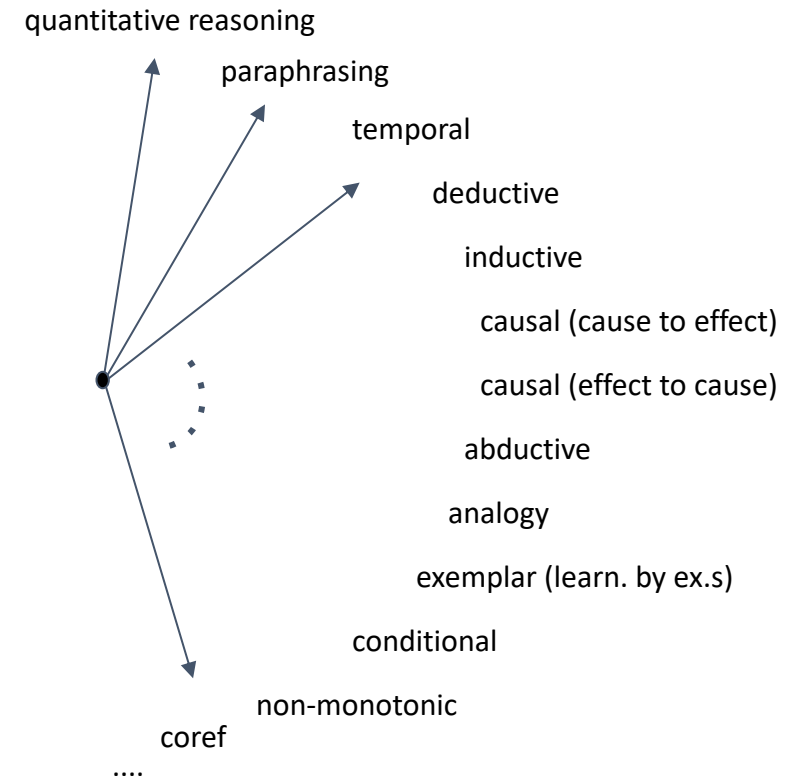




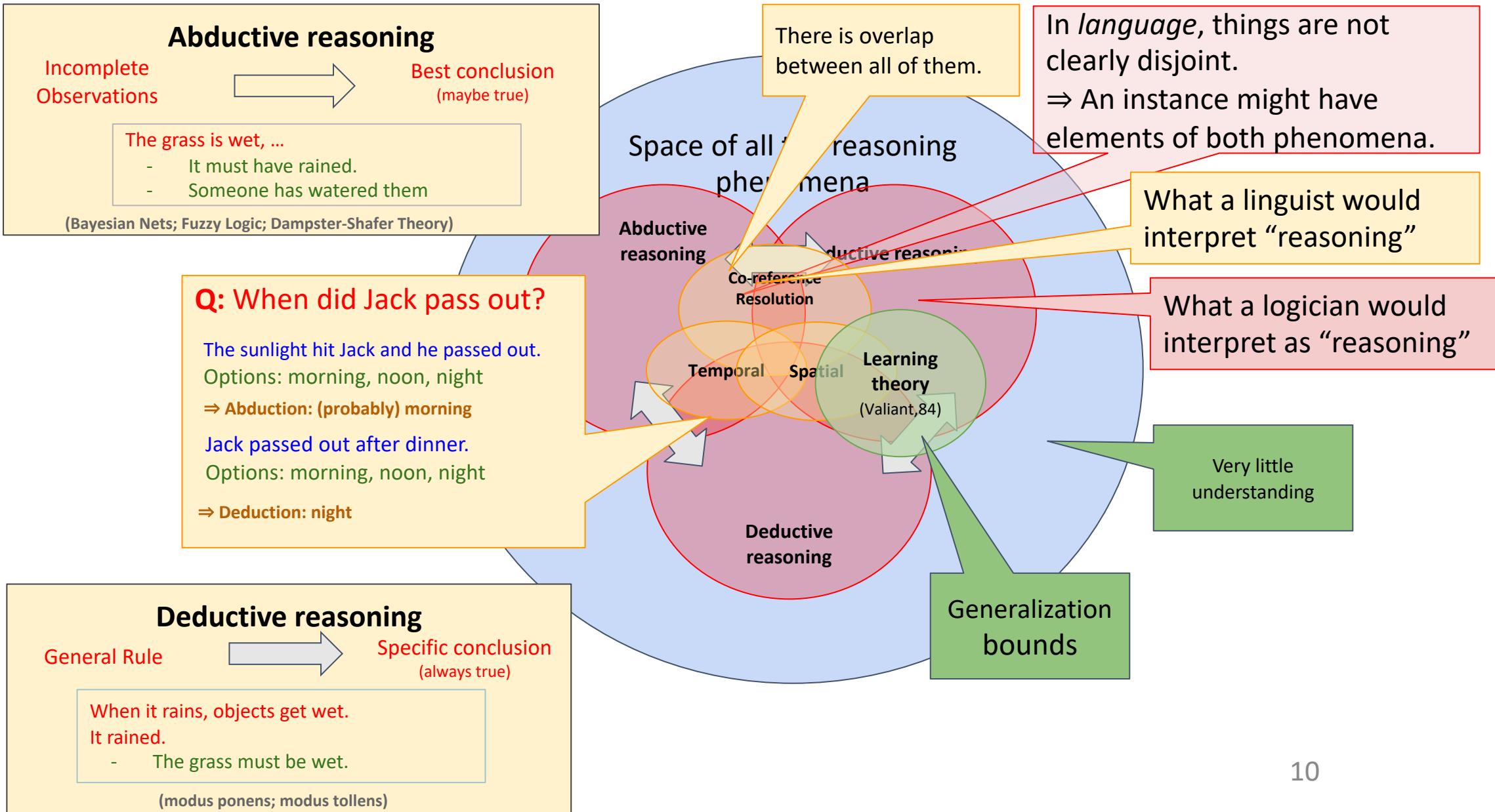
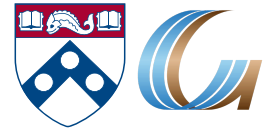
# The many faces of reasoning

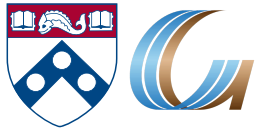


- Reasoning is often studied in a very narrow sense.
  - But probably has many forms
  - Realistic examples typically span multiple reasoning aspects.



# The many faces of reasoning





- **Idea:** represent all your knowledge in First Order Logic (KB).
- **Given a query  $\alpha$ , determine whether it holds in the KB: (KB implies  $\alpha$ )**
- For efficiency reasons:
  - FOL (too complex to compute with)  $\rightarrow$  Propositional Logic

Facts:

- Joe is married to Sue
- Bill has a brother with no children.
- Henry's friends are Bill's cousins.

(Declarative) Knowledge:

- *Ancestor* is the transitive closure of *parent*.
- *Brother* is *sibling* restricted to males
- *Favourite-cousin* is a special type of *cousin*.

Representation:

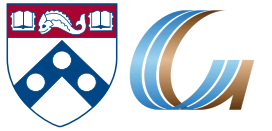
$\forall x \text{ Friend}(\text{henry}, x) \equiv \text{Cousin}(\text{bill}, x)$

- **Problem I:** complexity of inference.
  - Key solution: relax expressivity.
- **(but of, course, there were many other problems – incomplete knowledge, uncertainty)**
  - E.g., what if the knowledge is not given, but rather learned?

- **Given a query  $\alpha$ , determine whether it holds in the KB: (KB implies  $\alpha$ )**
  - Assume that KB is a collection of propositional rules:  $\mathbf{p} \rightarrow \mathbf{q}$  ; this is equivalent to:  $\neg \mathbf{p} \vee \mathbf{q} \equiv \mathbf{T}$  (a tautology)
    - $\mathbf{p}$  itself can be a conjunction of propositions;
    - $\mathbf{q}$  can be a disjunction of propositions (if it a conjunction, we'll split to multiple rules.)
  - Then the KB is a conjunction of disjunctions: **a CNF**
  - Answering  $\text{KB} \models \alpha$  is equivalent to solving satisfiability for  $\text{KB} \wedge \neg \alpha$ 
    - Determining that  $\text{KB} \wedge \neg \alpha$  has no satisfying assignments.
    - There is a lot of algorithmic proof theory to develop, under some conditions, efficient algorithms for  $\text{KB} \models \alpha$ 
      - E.g., if all the rule in KB are Horn rules (monotone antecedent, a single head proposition) there is an efficient algorithm.

- **Given a query  $\alpha$ , determine whether it holds in the KB: (KB implies  $\alpha$ )**
  
- But, exact reasoning could be too hard.
  
- And, what if KB is only approximate?
  - Model theory may makes more sense here.
    - $KB \models \alpha$  means that all the assignments that satisfy KB also satisfy  $\alpha$ .
    - Of course, there are too many assignments...
  - PAC semantics: what if you “sample” KB.
    - See Learning to Reason,(Khardon & Roth 96); an approach that is independent of the size of KB
    - This algorithm is complete, but not sound.
      - If  $KB \models \alpha$  it never errs. Otherwise, it may not find a counter examples.
  - It is also possible, under some conditions, to develop **exact** Learning to Reason
    - Under some assumptions on the type of queries, it is possible to find a polynomial size set of examples in KB such that is sufficient to test the query on these.

# Many Other Representations



- Limited Forms of FOL

- Relational Databases: `Course(csc248)`   `Dept(csc248,ComputerScience)`   `Enrollment(csc248,42)`  
`Course(mat100)`   `Dept(mat100,Mathematics)`

- And the hope is that you can address questions such as:

*How many courses are offered by the Computer Science Department?*

- Many other representations were developed, some along with inference systems.

- Logic Programs (Prolog): a collection of Horn sentences

$\forall x_1 \dots x_n [P_1 \wedge \dots \wedge P_m \supset P_{m+1}]$    where  $m \geq 0$  and each  $P_i$  is atomic

- For example:

```
parent(bill,mary).  
parent(bill,sam).  
mother(X,Y) :- parent(X,Y), female(Y)  
female(mary).
```

- Now I can infer who is the Mother of Bill (if I execute the program)

- Knowledge:

- $\text{Actor}(a) \Rightarrow \neg \text{Director}(a)$
- $\text{Director}(a) \Rightarrow \neg \text{WorkedFor}(a,b)$
- $\text{InMovie}(m,a) \wedge \text{WorkedFor}(a,b) \Rightarrow \text{InMovie}(m,b)$

- Input:

- $\text{Actor}(\text{Brando}), \text{Actor}(\text{Cruise}), \text{Director}(\text{Coppola}),$
- $\text{WorkedFor}(\text{Brando}, \text{Coppola}), \text{etc.}$

- Query:

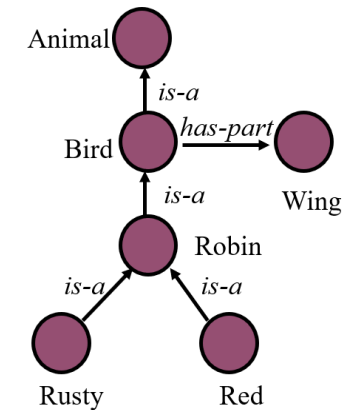
- is  $(\text{InMovie}(\text{GodFather}, \text{Brando}))$  ?
- is (what is the probability that:  $\text{Pr}(\text{InMovie}(\text{GodFather}, \text{Brando})) = ?$

- Abductive version:

- What is the most likely table for  $\text{InMovie}$ ?

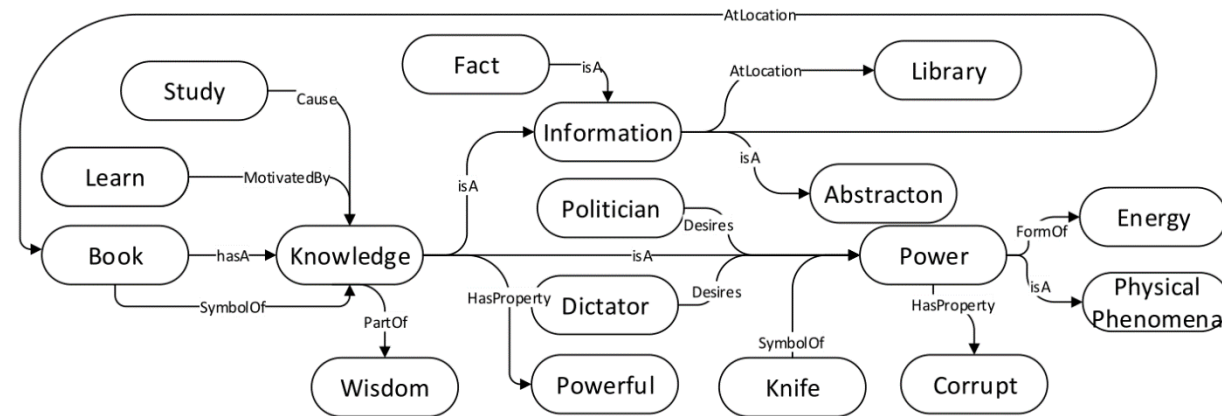
- Semantic Networks: allows the use of more expressive predicates, and more “intuitive inference”.
  - People talked about inference as a form of “spreading activation”
    - A graph of labeled nodes and labeled, directed arcs
    - Arcs define relationships that hold between objects denoted by the nodes.

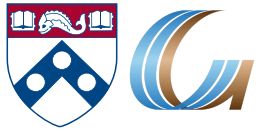
Link Type	Semantics	Example
$A \xrightarrow{\text{Subset}} B$	$A \subset B$	$Cats \subset Mammals$
$A \xrightarrow{\text{Member}} B$	$A \in B$	$Bill \in Cats$
$A \xrightarrow{R} B$	$R(A, B)$	$Bill \xrightarrow{\text{Age}} 12$
$A \xrightarrow{\boxed{R}} B$	$\forall x, x \in A \Rightarrow R(x, B)$	$Bird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





- This led to two directions:
- **(1) Concept nets:**
  - Based on Open Mind Common Sense (OMCS)
  - Intended to serve as a large commonsense knowledge base
  - Built on contributions of many people across the Web.





- (2) Formalization efforts:

- These networks were formalized in terms of **Description Logics**, and then elaborated into **Frame Description Forms**.
- **Frames** were used to describe types and their attributes: values, Restrictions, attached procedures (how an attribute should be used).

(Student  
with a dept is computer-science and  
with  $\geq 3$  enrolled-course is a  
(Graduate-Course  
with a dept is a Engineering-Department))

- Eventually, this led to theories of **Frames** (Minsky), and **Scripts** (Schank)
  - There are beginning to be influential again, where people think more about **Events**

- More generally, these languages had expressive grammars:

$$\begin{aligned} \langle type \rangle ::= & \langle atom \rangle \\ & | (\text{AND } \langle type_1 \rangle \dots \langle type_n \rangle) \\ & | (\text{ALL } \langle attribute \rangle \langle type \rangle) \\ & | (\text{SOME } \langle attribute \rangle) \end{aligned}$$
$$\begin{aligned} \langle attribute \rangle ::= & \langle atom \rangle \\ & | (\text{RESTR } \langle attribute \rangle \langle type \rangle) \end{aligned}$$

- Example: The set of all people the all their male friends are doctors with some specialty.

$$(\text{AND person } (\text{ALL } (\text{RESTR friend male}) (\text{AND doctor } (\text{SOME specialty}))))).$$

- And it came with inference algorithms – **subsumption**, and was extremely influential – all systems built in the 80-ith and later, were built on these languages.
  - It was also influential in areas such as Feature Extraction for machine learning, and theories of grammar.

- Knowledge:

- $\text{Actor}(a) \Rightarrow \neg \text{Director}(a)$
- $\text{Director}(a) \Rightarrow \neg \text{WorkedFor}(a,b)$
- $\text{InMovie}(m,a) \wedge \text{WorkedFor}(a,b) \Rightarrow \text{InMovie}(m,b)$

- Input:

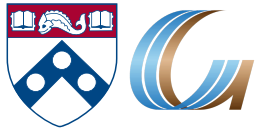
- $\text{Actor}(\text{Brando}), \text{Actor}(\text{Cruise}), \text{Director}(\text{Coppola}),$
- $\text{WorkedFor}(\text{Brando}, \text{Coppola}), \text{etc.}$

- Query:

- is  $(\text{InMovie}(\text{GodFather}, \text{Brando}))$  ?
- is (what is the probability that:  $\text{Pr}(\text{InMovie}(\text{GodFather}, \text{Brando})) = ?$

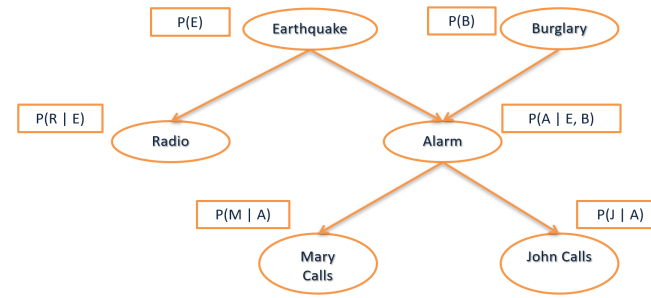
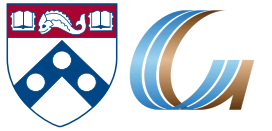
- Abductive version:

- What is the most likely table for  $\text{InMovie}$ ?



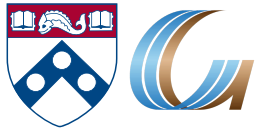
- In parallel to the progress on the logical representations, people argued that we need to deal with uncertainty, and need to move to **probabilistic representations**.
- Progress here proceeded in two camps
  - (Propositional) representation of distributions
    - **Bayesian Networks** (Pearl 1988)
  - Probabilistic extensions of the FOL/Prolog representations. (Haddawy 1993)
    - **Problog**
    - Markov Logic Network
- Two important comments:
  - The latter direction is presented today as fusing probabilities with declarative (logical) knowledge. This, in fact, was studied much earlier (in the 60—ies), but without practical implementations.
  - Fusing Probabilities with Declarative information is different from fusing Learning with Declarative Information. In fact, none of the bullets above came with a native approach for **learning**.
  - Fusing learning with declarative knowledge came later in the context of Structured Learning, e.g., ILP Formulations, Roth & Yih 2004, and following works.

# Bayes Nets



- Nodes are random variables
- Edges represent causal influences
- Each node is associated with a conditional Probability distribution
- **Computational Problems (Inference):**
  - **Computing the probability of an event:**
  - Given structure and parameters
  - Given an observation  $E$ , what is the probability of assignment  $Y$ ?
  - $P(R=\text{off}, A=\text{off} \mid E=e) = ?$  ( $E, Y$  are sets of instantiated variables)
- **Most likely explanation (Maximum A Posteriori assignment, MAP, MPE)**
  - Given structure and parameters
  - Given an observation  $E$ , what is the most likely assignment to  $Y$ ?
  - $\text{Argmax}_Y P(Y=y \mid E=e)$  (Say,  $Y = (R, A)$ )
  - ( $E, Y$  are sets of instantiated variables)

# Probabilistic Relational Representations



- Representation of distributions over relations, as opposed to propositional variables.
- Ability to build programs that do not only encode complex interactions between variables but also express inherent uncertainty.
- **Inference:** Becoming much harder. For the most part, done by **propositionalizing** relational representations (that is, substitution of all domain variables, and blowing up the representations to get a propositional BN).
- But, there are other ways, e.g., **lifted inference**.

```
0.3::stress(X) :- person(X).
0.2::influences(X,Y) :- person(X), person(Y).

smokes(X) :- stress(X).
smokes(X) :- friend(X,Y), influences(Y,X), smokes(Y).

0.4::asthma(X) :- smokes(X).

person(angelika).
person(joris).
person(jonas).
person(dimitar).

friend(joris,jonas).
friend(joris,angelika).
friend(joris,dimitar).
friend(angelika,jonas).
```