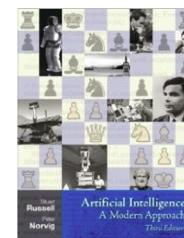


Logic and AI

- Logical Agents
- Propositional Logic
 - Syntax
 - Semantics
 - Resolution
- First Order Logic
 - Syntax
 - Semantics
 - Unification
 - Skolem
 - Resolution

Material from
Russell & Norvig,
chapters 7-9



Many slides based on
Russell & Norvig's slides
[Artificial Intelligence:
A Modern Approach](#)

Some based on Slides by
Vincent Conitzer and Sriraam Natarajan

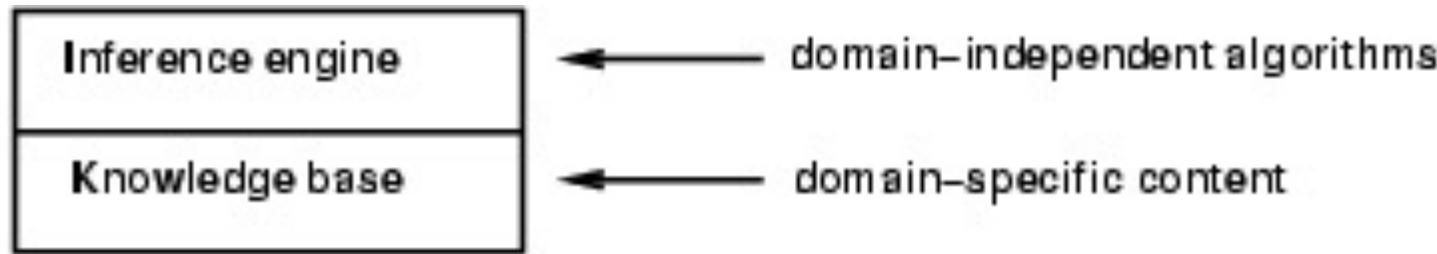
Logic and AI

- Would like our AI to have knowledge about the world, and logically draw conclusions from it
- Search algorithms generate successors and evaluate them, but do not “understand” much about the setting
- Example question: is it possible for a chess player to have 8 pawns and 2 queens?
 - Search algorithm could search through tons of states to see if this ever happens, but...

A “roommate” story

- You roommate comes home; he/she is completely wet
- You know the following things:
 - Your roommate is wet
 - If your roommate is wet, it is because of rain, sprinklers, or both
 - If your roommate is wet because of sprinklers, the sprinklers must be on
 - If your roommate is wet because of rain, your roommate must not be carrying the umbrella
 - The umbrella is not in the umbrella holder
 - If the umbrella is not in the umbrella holder, either you must be carrying the umbrella, or your roommate must be carrying the umbrella
 - You are not carrying the umbrella
- **Can you conclude that the sprinklers are on?**
- **Can AI conclude that the sprinklers are on?**

Let us separate knowledge from inference engines: Knowledge Bases



- Knowledge base = set of **sentences** in a **formal** language
- **Declarative** approach to building an agent (or other system):
 - Tell it what it needs to know
- Then it can Ask itself what to do - answers should follow from the KB
- Agents can be viewed at the **knowledge level** i.e., what they know, regardless of how implemented

A simple knowledge-based agent

```
function KB-AGENT(percept) returns an action
    static: KB, a knowledge base
        t, a counter, initially 0, indicating time

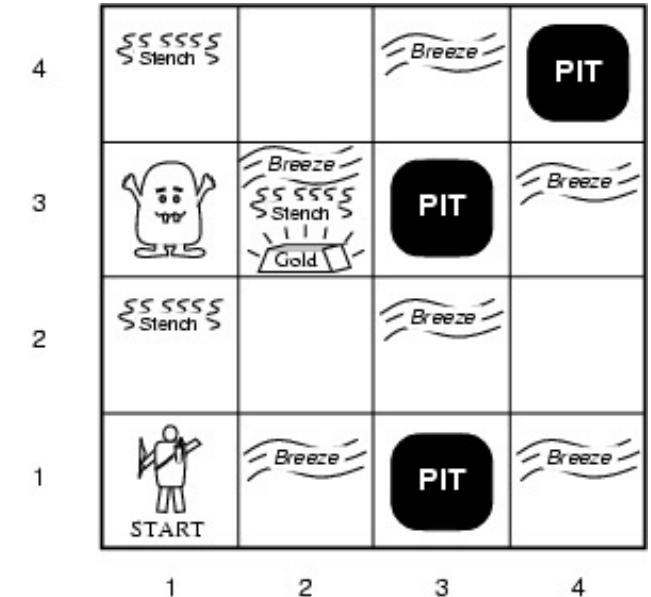
    TELL(KB, MAKE-PERCEPT-SENTENCE(percept, t))
    action  $\leftarrow$  ASK(KB, MAKE-ACTION-QUERY(t))
    TELL(KB, MAKE-ACTION-SENTENCE(action, t))
    t  $\leftarrow$  t + 1
    return action
```

The agent must be able to:

- Represent states, actions, etc.
- Incorporate new percepts
- Update internal representations of the world
- Deduce hidden properties of the world
- Deduce appropriate actions

Wumpus World PEAS description

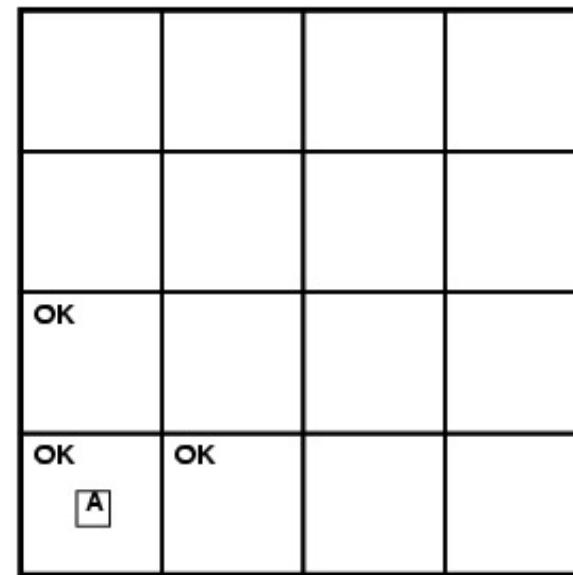
- **Performance measure**
 - gold +1000, death -1000
 - -1 per step, -10 for using the arrow
- **Environment**
 - Squares adjacent to wumpus are smelly
 - Squares adjacent to pit are breezy
 - Glitter iff gold is in the same square
 - Shooting kills wumpus if you are facing it
 - Shooting uses up the only arrow
 - Grabbing picks up gold if in same square
 - Releasing drops the gold in same square
- **Sensors:** Stench, Breeze, Glitter, Bump, Scream
- **Actuators:** Left turn, Right turn, Forward, Grab, Release, Shoot



Wumpus world characterization

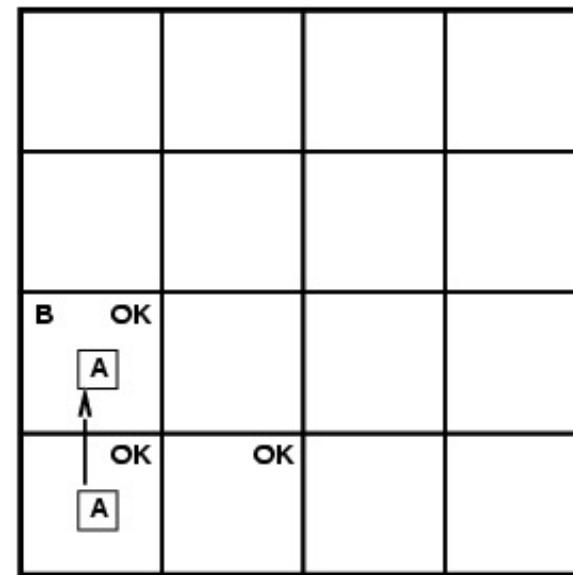
- Fully Observable No – only local perception
- Deterministic Yes – outcomes exactly specified
- Episodic No – sequential at the level of actions
- Static Yes – Wumpus and Pits do not move
- Discrete Yes
- Single-agent? Yes – Wumpus is essentially a natural feature

Exploring a wumpus world



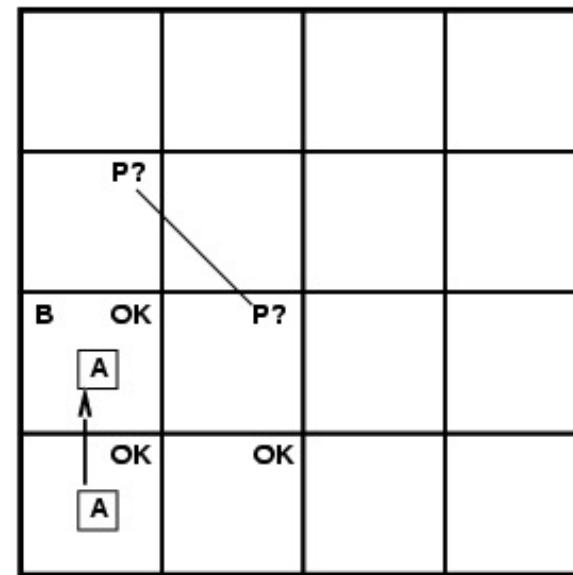
Sensors: Stench, Breeze, Glitter, Bump, Scream

Exploring a wumpus world



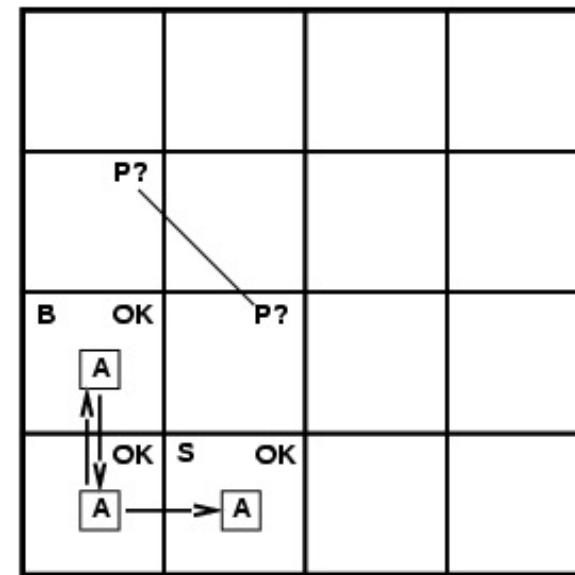
Sensors: Stench, Breeze, Glitter, Bump, Scream

Exploring a wumpus world



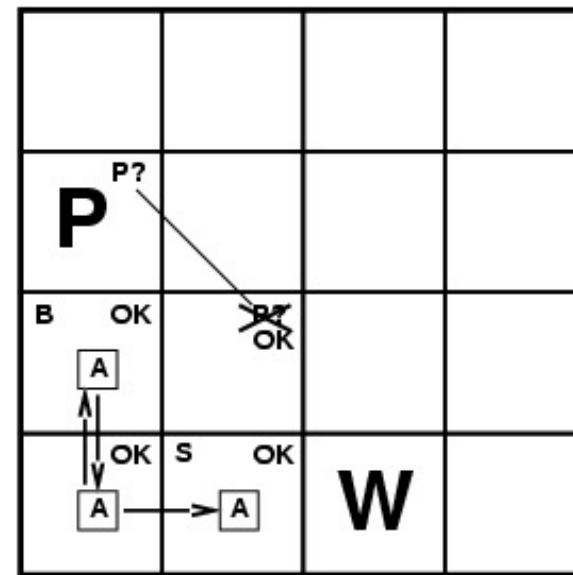
Sensors: Stench, Breeze, Glitter, Bump, Scream

Exploring a wumpus world



Sensors: Stench, Breeze, Glitter, Bump, Scream

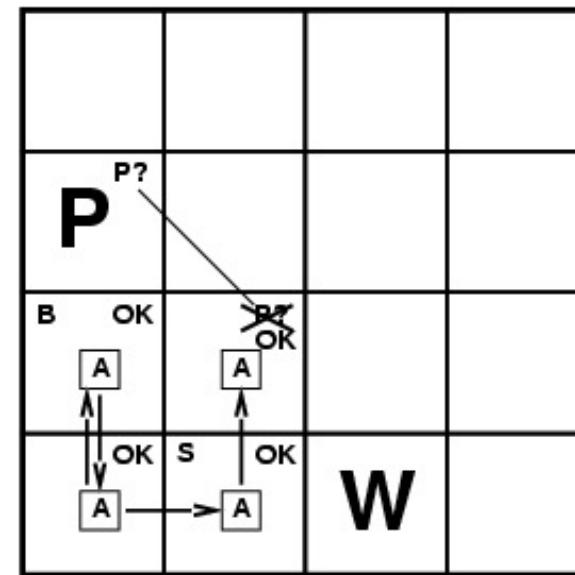
Exploring a wumpus world



Reasoning
turns our
observation
into knowledge

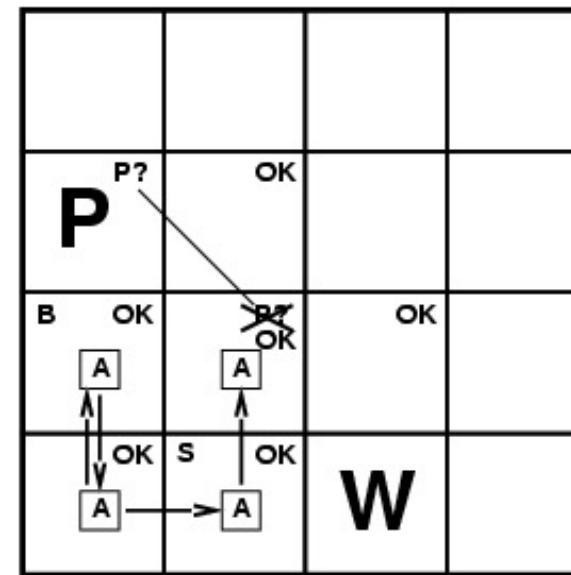
Sensors: Stench, Breeze, Glitter, Bump, Scream

Exploring a wumpus world



Sensors: Stench, Breeze, Glitter, Bump, Scream

Exploring a wumpus world



Sensors: Stench, Breeze, Glitter, Bump, Scream

Knowledge base for the roommate story

- RoommateWet
- RoommateWet => (RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers)
- RoommateWetBecauseOfSprinklers => SprinklersOn
- RoommateWetBecauseOfRain =>
NOT(RoommateCarryingUmbrella)
- UmbrellaGone
- UmbrellaGone => (YouCarryingUmbrella OR RoommateCarryingUmbrella)
- NOT(YouCarryingUmbrella)

Syntax of Propositional Logic)

- What do well-formed sentences in the knowledge base look like?
- A BNF grammar:
- $\text{Symbol} \rightarrow P, Q, R, \dots, \text{RoommateWet}, \dots$
- $\text{Sentence} \rightarrow \text{True} \mid \text{False} \mid \text{Symbol} \mid \text{NOT}(\text{Sentence}) \mid (\text{Sentence AND Sentence}) \mid (\text{Sentence OR Sentence}) \mid (\text{Sentence} \Rightarrow \text{Sentence})$
- We will drop parentheses sometimes, but formally they really should always be there

Semantics

- **Interpretation** specifies which of the proposition symbols are true and which are false
- Given a interpretation, I should be able to tell you whether a sentence is true or false. **Truth table** defines semantics of operators:

a	b	NOT(a)	a AND b	a OR b	a => b
false	false	true	false	false	true
false	true	true	false	true	true
true	false	false	false	true	false
true	true	false	true	true	true

- Given an interpretation, can compute truth of sentence recursively with these.
- A **model** (of a set of sentences) is an interpretation in which all the sentences are true.

Caveats

- TwolsAnEvenNumber OR
ThreelsAnOddNumber
is true (not exclusive OR)
- TwolsAnOddNumber =>
ThreelsAnEvenNumber
is true (if the left side is false it's always true)

All of this is assuming those symbols are assigned their natural values...

Tautologies

- A sentence is a **tautology** if it is true for any setting of its propositional symbols

P	Q	P OR Q	NOT(P) AND NOT(Q)	(P OR Q) OR (NOT(P)) AND NOT(Q))
false	false	false	true	true
false	true	true	false	true
true	false	true	false	true
true	true	true	false	true

- $(P \text{ OR } Q) \text{ OR } (\text{NOT}(P) \text{ AND } \text{NOT}(Q))$ is a tautology

Is this a tautology?

- $(P \Rightarrow Q) \text{ OR } (Q \Rightarrow P)$

Logical equivalences

- Two sentences are **logically equivalent** if they have the same truth value for every setting of their propositional variables

P	Q	P OR Q	NOT(NOT(P)) AND NOT(Q))
false	false	false	false
false	true	true	true
true	false	true	true
true	true	true	true

- P OR Q and NOT(NOT(P)) AND NOT(Q)) are logically equivalent
- Tautology = logically equivalent to True**

Famous logical equivalences

they can be used for rewriting and simplifying rules

- $(a \text{ OR } b) \equiv (b \text{ OR } a)$ *commutativity*
- $(a \text{ AND } b) \equiv (b \text{ AND } a)$ *commutativity*
- $((a \text{ AND } b) \text{ AND } c) \equiv (a \text{ AND } (b \text{ AND } c))$ *associativity*
- $((a \text{ OR } b) \text{ OR } c) \equiv (a \text{ OR } (b \text{ OR } c))$ *associativity*
- $\text{NOT}(\text{NOT}(a)) \equiv a$ *double-negation elimination*
- $(a \Rightarrow b) \equiv (\text{NOT}(b) \Rightarrow \text{NOT}(a))$ *contraposition*
- $(a \Rightarrow b) \equiv (\text{NOT}(a) \text{ OR } b)$ *implication elimination*
- $\text{NOT}(a \text{ AND } b) \equiv (\text{NOT}(a) \text{ OR } \text{NOT}(b))$ *De Morgan*
- $\text{NOT}(a \text{ OR } b) \equiv (\text{NOT}(a) \text{ AND } \text{NOT}(b))$ *De Morgan*
- $(a \text{ AND } (b \text{ OR } c)) \equiv ((a \text{ AND } b) \text{ OR } (a \text{ AND } c))$ *distributivity*
- $(a \text{ OR } (b \text{ AND } c)) \equiv ((a \text{ OR } b) \text{ AND } (a \text{ OR } c))$ *distributivity*

Is this a tautology?

- $(P \Rightarrow Q) \text{ OR } (Q \Rightarrow P)$
- $(\text{not}(P) \text{ OR } Q) \text{ OR } (\text{not}(Q) \text{ or } P)$
- $\text{not}(P) \text{ OR } Q \text{ OR } \text{not}(Q) \text{ or } P$
- $(\text{not}(P) \text{ OR } P) \text{ OR } (\text{not}(Q) \text{ or } Q)$
- $(\text{true}) \text{ OR } (\text{true})$
- true

Wumpus world sentences

- Let $P_{i,j}$ be true if there is a pit in $[i, j]$.
- Let $B_{i,j}$ be true if there is a breeze in $[i, j]$.
 - $\neg P_{1,1}$
 - $\neg B_{1,1}$
 - $B_{2,1}$
- "Pits cause breezes in adjacent squares"
 - $B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})$
 - $B_{2,1} \Leftrightarrow (P_{1,1} \vee P_{2,2} \vee P_{3,1})$

Wumpus world sentences

- Let $P_{i,j}$ be true if there is a pit in $[i, j]$.
- Let $B_{i,j}$ be true if there is a breeze in $[i, j]$.
 - $\neg P_{1,1}$
 - $\neg B_{1,1}$
 - $B_{2,1}$
- "Pits cause breezes in adjacent squares"
 - $B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})$
 - $B_{2,1} \Leftrightarrow (P_{1,1} \vee P_{2,2} \vee P_{3,1})$

Great, we know how to code our knowledge. We know what is a model, and we have seen basic options to rewrite our knowledge without changing its semantics.

But how do we infer that something is true?

Inference / Entailment

- We have a knowledge base of things that we know are true
 - RoommateWetBecauseOfSprinklers
 - RoommateWetBecauseOfSprinklers => SprinklersOn
- Can we conclude that SprinklersOn?
- We say SprinklersOn is **entailed** by the knowledge base if, for every setting (models) of the propositional variables for which the knowledge base is true, SprinklersOn is also true

RWBOS	SprinklersOn	Knowledge base
false	false	false
false	true	false
true	false	false
true	true	true

- **SprinklersOn is entailed!**

Wumpus Truth tables for inference

$B_{1,1}$	$B_{2,1}$	$P_{1,1}$	$P_{1,2}$	$P_{2,1}$	$P_{2,2}$	$P_{3,1}$	KB	α_1
<i>false</i>	<i>true</i>							
<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	<i>true</i>
:	:	:	:	:	:	:	:	:
<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>
<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	<u><i>true</i></u>	<u><i>true</i></u>
<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	<u><i>true</i></u>	<u><i>true</i></u>
<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>true</i>	<u><i>true</i></u>	<u><i>true</i></u>
<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>
:	:	:	:	:	:	:	:	:
<i>true</i>	<i>false</i>	<i>false</i>						

Simple algorithm for inference

- Want to find out if sentence a is entailed by knowledge base...
- *Go through the possible settings of the propositional variables,*
 - *If knowledge base is true and a is false, return false*
- *Return true*
- Not very efficient: $2^{\#\text{propositional variables}}$ settings

Inconsistent knowledge bases

- Suppose we were careless in how we specified our knowledge base:
 - $\text{PetOfRoommateIsABird} \Rightarrow \text{PetOfRoommateCanFly}$
 - $\text{PetOfRoommateIsAPenguin} \Rightarrow \text{PetOfRoommateIsABird}$
 - $\text{PetOfRoommateIsAPenguin} \Rightarrow \text{NOT}(\text{PetOfRoommateCanFly})$
 - $\text{PetOfRoommateIsAPenguin}$
- It entails both $\text{PetOfRoommateCanFly}$ and $\text{NOT}(\text{PetOfRoommateCanFly})$
- Therefore, technically, this knowledge base implies anything: The Moon Is Made Of Cheese

So, make sure that your (logical) knowledge base is consistent!

The Moon is Made of Cheese

- PetOfRoommateCanFly AND NOT(PetOfRoommateCanFly)
 - PetOfRoommateCanFly, NOT(PetOfRoommateCanFly)
 - Now, “a true statement OR anything else” is always true, hence
 - PetOfRoommateCanFly OR MoonMadeOfCheese
 - Therefore, MoonMadeOfCheese has to be true.
-
- Please note that you can really put anything there!

So, we justify the Aristotelian claim that
“there cannot be contradictions”
(The Principle of Non-Contradiction)

Reasoning patterns

- Obtain new sentences directly from some other sentences in knowledge base according to reasoning patterns
- If we have sentences a and $a \Rightarrow b$, we can correctly conclude the new sentence b
 - This is called modus ponens
 - If we have $a \text{ AND } b$, we can correctly conclude a
 - All of the logical equivalences from before also give reasoning patterns

Formal proof that the sprinklers are on

- 1) RoommateWet
- 2) RoommateWet => (RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers)
- 3) RoommateWetBecauseOfSprinklers => SprinklersOn
- 4) RoommateWetBecauseOfRain => NOT(RoommateCarryingUmbrella)
- 5) UmbrellaGone
- 6) UmbrellaGone => (YouCarryingUmbrella OR RoommateCarryingUmbrella)
- 7) NOT(YouCarryingUmbrella) Knowledge Base
- 8) YouCarryingUmbrella OR RoommateCarryingUmbrella (*modus ponens on 5 and 6*)
- 9) NOT(YouCarryingUmbrella) => RoommateCarryingUmbrella (*equivalent to 8*)
- 10) RoommateCarryingUmbrella (*modus ponens on 7 and 9*)
- 11) NOT(NOT(RoommateCarryingUmbrella)) (*equivalent to 10*)
- 12) NOT(NOT(RoommateCarryingUmbrella)) => NOT(RoommateWetBecauseOfRain) (*equivalent to 4 by contraposition*)
- 13) NOT(RoommateWetBecauseOfRain) (*modus ponens on 11 and 12*)
- 14) RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers (*modus ponens on 1 and 2*)
- 15) NOT(RoommateWetBecauseOfRain) => RoommateWetBecauseOfSprinklers (*equivalent to 14*)
- 16) RoommateWetBecauseOfSprinklers (*modus ponens on 13 and 15*)
- 17) SprinklersOn (*modus ponens on 16 and 3*)

Reasoning about penguins

- 1) PetOfRoommateIsABird => PetOfRoommateCanFly
- 2) PetOfRoommateIsAPenguin => PetOfRoommateIsABird
- 3) PetOfRoommateIsAPenguin => NOT(PetOfRoommateCanFly)
- 4) PetOfRoommateIsAPenguin
- 5) PetOfRoommateIsABird (*modus ponens on 4 and 2*)
- 6) PetOfRoommateCanFly (*modus ponens on 5 and 1*)
- 7) NOT(PetOfRoommateCanFly) (*modus ponens on 4 and 3*)
- 8) NOT(PetOfRoommateCanFly) => FALSE (*equivalent to 6*)
- 9) FALSE (*modus ponens on 7 and 8*)
- 10) FALSE => TheMoonIsMadeOfCheese (*tautology*)
- 11) TheMoonIsMadeOfCheese (*modus ponens on 9 and 10*)

Getting more systematic

(Think of it as making your life easier when implementing this for any knowledge base on a computer)

- Any knowledge base can be written as a single formula in **conjunctive normal form (CNF)**
 - CNF formula: (... OR ... OR ...) AND (... OR ...) AND ...
 - ... can be a symbol x , or $\text{NOT}(x)$ (these are called **literals**)
 - Multiple facts in knowledge base are effectively ANDed together

**RoommateWet => (RoommateWetBecauseOfRain OR
RoommateWetBecauseOfSprinklers)**

becomes

**(NOT(RoommateWet) OR RoommateWetBecauseOfRain
OR RoommateWetBecauseOfSprinklers)**

Converting story problem to conjunctive normal form

- RoommateWet
 - RoommateWet
- RoommateWet => (RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers)
 - NOT(RoommateWet) OR RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers
- RoommateWetBecauseOfSprinklers => SprinklersOn
 - NOT(RoommateWetBecauseOfSprinklers) OR SprinklersOn
- RoommateWetBecauseOfRain => NOT(RoommateCarryingUmbrella)
 - NOT(RoommateWetBecauseOfRain) OR NOT(RoommateCarryingUmbrella)
- UmbrellaGone
 - UmbrellaGone
- UmbrellaGone => (YouCarryingUmbrella OR RoommateCarryingUmbrella)
 - NOT(UmbrellaGone) OR YouCarryingUmbrella OR RoommateCarryingUmbrella
- NOT(YouCarryingUmbrella)
 - NOT(YouCarryingUmbrella)

Now that we have a normal form, we can implement general reasoning patterns for this normal form. One of them is

Unit resolution

If we have

- $I_1 \text{ OR } I_2 \text{ OR } \dots \text{ OR } I_k$

and

- $\text{NOT}(I_i)$

we can conclude

- $I_1 \text{ OR } I_2 \text{ OR } \dots \text{ OR } I_{i-1} \text{ OR } I_{i+1} \text{ OR } \dots \text{ OR } I_k$
- **Basically modus ponens**

Applying resolution to story problem

- 1) RoommateWet
- 2) NOT(RoommateWet) OR RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers
- 3) NOT(RoommateWetBecauseOfSprinklers) OR SprinklersOn
- 4) NOT(RoommateWetBecauseOfRain) OR NOT(RoommateCarryingUmbrella)
- 5) UmbrellaGone
- 6) NOT(UmbrellaGone) OR YouCarryingUmbrella OR RoommateCarryingUmbrella
- 7) NOT(YouCarryingUmbrella) Knowledge Base
- 8) NOT(UmbrellaGone) OR RoommateCarryingUmbrella (6,7)
- 9) RoommateCarryingUmbrella (5,8)
- 10) NOT(RoommateWetBecauseOfRain) (4,9)
- 11) NOT(RoommateWet) OR RoommateWetBecauseOfSprinklers (2,10)
- 12) RoommateWetBecauseOfSprinklers (1,11)
- 13) SprinklersOn (3,12)

Unfortunately, unit resolution is not enough!

Limitations of unit resolution

- $P \text{ OR } Q$
- $\text{NOT}(P) \text{ OR } Q$
- Can we conclude Q ?

(General) resolution

if we have

- $I_1 \text{ OR } I_2 \text{ OR } \dots \text{ OR } I_k$
and
- $m_1 \text{ OR } m_2 \text{ OR } \dots \text{ OR } m_n$
where for some i, j , $I_i = \text{NOT}(m_j)$

we can conclude

- $I_1 \text{ OR } I_2 \text{ OR } \dots \text{ } I_{i-1} \text{ OR } I_{i+1} \text{ OR } \dots \text{ OR } I_k \text{ OR } m_1 \text{ OR } m_2 \text{ OR } \dots \text{ OR } m_{j-1} \text{ OR } m_{j+1} \text{ OR } \dots \text{ OR } m_n$
- Same literal may appear multiple times; remove those

Applying resolution to story problem (more clumsily)

- 1) RoommateWet
- 2) NOT(RoommateWet) OR RoommateWetBecauseOfRain OR RoommateWetBecauseOfSprinklers
- 3) NOT(RoommateWetBecauseOfSprinklers) OR SprinklersOn
- 4) NOT(RoommateWetBecauseOfRain) OR NOT(RoommateCarryingUmbrella)
- 5) UmbrellaGone
- 6) NOT(UmbrellaGone) OR YouCarryingUmbrella OR RoommateCarryingUmbrella
- 7) NOT(YouCarryingUmbrella) Knowledge Base

- 8) NOT(RoommateWet) OR RoommateWetBecauseOfRain OR SprinklersOn (2,3)
- 9) NOT(RoommateCarryingUmbrella) OR NOT(RoommateWet) OR SprinklersOn (4,8)
- 10) NOT(UmbrellaGone) OR YouCarryingUmbrella OR NOT(RoommateWet) OR SprinklersOn (6,9)
- 11) YouCarryingUmbrella OR NOT(RoommateWet) OR SprinklersOn (5,10)
- 12) NOT(RoommateWet) OR SprinklersOn (7,11)
- 13) SprinklersOn (1,12)

How to use this for a systematic inference?

- General strategy: if we want to see if sentence a is entailed, add $\text{NOT}(a)$ to the knowledge base and see if it becomes inconsistent (we can derive a contradiction)
- CNF formula for modified knowledge base is satisfiable if and only if sentence a is not entailed
 - Satisfiable = there exists a model that makes the modified knowledge base true = modified knowledge base is consistent

Resolution algorithm

- Given formula in conjunctive normal form,
repeat:
 - Find two clauses with complementary literals,
 - Apply resolution,
 - Add resulting clause (if not already there)
 - If the **empty clause results**, **formula is not satisfiable**
 - Must have been obtained from P and $\text{NOT}(P)$

Example

Our knowledge base:

- 1) RoommateWetBecauseOfSprinklers
- 2) NOT(RoommateWetBecauseOfSprinklers) OR SprinklersOn

Can we infer SprinklersOn?

- We add:
 - 3) NOT(SprinklersOn)
- From 2) and 3), get
 - 4) NOT(RoommateWetBecauseOfSprinklers)
- From 4) and 1), get empty clause

Special case: Horn clauses

- Horn clauses are implications with only positive literals
$$x_1 \text{ AND } x_2 \text{ AND } x_4 \Rightarrow x_3 \text{ AND } x_6$$
$$\text{TRUE} \Rightarrow x_1$$
- Try to figure out whether some x_j is entailed
- Simply follow the implications (modus ponens) as far as you can, see if you can reach x_j
- x_j is entailed if and only if it can be reached (can set everything that is not reached to false)
- Can implement this more efficiently by maintaining, for each implication, a count of how many of the left-hand side variables have been reached

Inference-based agents in the wumpus world

A wumpus-world agent using propositional logic:

$$\neg P_{1,1}$$

$$\neg W_{1,1}$$

$$B_{x,y} \Leftrightarrow (P_{x,y+1} \vee P_{x,y-1} \vee P_{x+1,y} \vee P_{x-1,y})$$

$$S_{x,y} \Leftrightarrow (W_{x,y+1} \vee W_{x,y-1} \vee W_{x+1,y} \vee W_{x-1,y})$$

$$W_{1,1} \vee W_{1,2} \vee \dots \vee W_{4,4}$$

$$\neg W_{1,1} \vee \neg W_{1,2}$$

$$\neg W_{1,1} \vee \neg W_{1,3}$$

...

⇒ 64 distinct proposition symbols, 155 sentences

Limitations of propositional logic

- Some English statements are hard to model in propositional logic:
“If your roommate is wet because of rain, your roommate must not be carrying any umbrella”
- Pathetic attempt at modeling this:
RoommateWetBecauseOfRain =>
(NOT(RoommateCarryingUmbrella0) AND
NOT(RoommateCarryingUmbrella1) AND
NOT(RoommateCarryingUmbrella2) AND ...)

Limitations of propositional logic

- No notion of **objects**
- No notion of **relations among objects**
- RoommateCarryingUmbrella0 is **instructive to us**, suggesting
 - there is an object we call Roommate,
 - there is an object we call Umbrella0,
 - there is a relationship Carrying between these two objects
- **Formally, none of this meaning is there**
 - Might as well have replaced RoommateCarryingUmbrella0 by P

Elements of first-order logic

- **Objects:** can give these names such as Umbrella0, Person0, John, Earth, ...
- **Relations:** Carrying(., .), IsAnUmbrella(.)
 - Carrying(Person0, Umbrella0),
IsUmbrella(Umbrella0)
 - Relations with one object = **unary relations**
= **properties**
- **Functions:** Roommate(.)
 - Roommate(Person0)
 - **Equality:** Roommate(Person0) = Person1

Things to note about functions

- Can be used to encode integers and also data structures.
 - 0, succ(0), succ(succ(0)), ...
 - tree(value2, tree(value2, ...), tree(value3, ...))
- It could be that we have a separate name for Roommate(Person0)
 - E.g., **Roommate(Person0) = Person1**
 - ... but we do not **need** to have such a name
- A function can be applied to any object
- E.g., **Roommate(Umbrella0)**

OK, but how do we talk and reason about many objects at once?

- **Variables:** x, y, z, ... can refer to multiple objects
- **We need two new operators** “for all” and “there exists” (**Universal quantifier** and **existential quantifier**)
- for all x: CompletelyWhite(x) => NOT(PartiallyBlack(x))
 - Completely white objects are never partially black
- there exists x: PartiallyWhite(x) AND PartiallyBlack(x)
 - There exists some object in the world that is partially white and partially black

Practice converting English to first-order logic

- “John has Jane’s umbrella”
- Has(John, Umbrella(Jane))
- “John has an umbrella”
- there exists y: (Has(John, y) AND IsUmbrella(y))
- “Anything that has an umbrella is not wet”
- for all x: ((there exists y: (Has(x, y) AND IsUmbrella(y))) => NOT(IsWet(x)))
- “Any person who has an umbrella is not wet”
- for all x: (IsPerson(x) => ((there exists y: (Has(x, y) AND IsUmbrella(y))) => NOT(IsWet(x))))

More practice converting English to first-order logic

- “John has at least two umbrellas”
- there exists x: (there exists y: (Has(John, x) AND IsUmbrella(x) AND Has(John, y) AND IsUmbrella(y) AND NOT(x=y)))
- “John has at most two umbrellas”
- for all x, y, z: ((Has(John, x) AND IsUmbrella(x) AND Has(John, y) AND IsUmbrella(y) AND Has(John, z) AND IsUmbrella(z)) => (x=y OR x=z OR y=z))

Even more practice converting English to first-order logic...

- “TUDa’s basketball team defeats any other basketball team”
- for all x: ((IsBasketballTeam(x) AND NOT(x=BasketballTeamOf(TUDa))) => Defeats(BasketballTeamOf(TUDa), x))
- “Every team defeats some other team”
- for all x: (IsTeam(x) => (there exists y: (IsTeam(y) AND NOT(x=y) AND Defeats(x,y))))

More realistically...

- “Any basketball team that defeats TUDa’s basketball team in one year will be defeated by TUDa’s basketball team in a future year”
- for all x,y: $(\text{IsBasketballTeam}(x) \text{ AND } \text{IsYear}(y) \text{ AND } \text{DefeatsIn}(x, \text{BasketballTeamOf}(\text{TUDa}), y)) \Rightarrow \text{there exists } z: (\text{IsYear}(z) \text{ AND } \text{IsLaterThan}(z, y) \text{ AND } \text{DefeatsIn}(\text{BasketballTeamOf}(\text{TUDa}), x, z))$

Relationship between universal and existential

There is a tight connection between universal and existential quantifiers that we will make use of for a systematic inference procedure

for all x : a

is equivalent to

NOT(there exists x : NOT(a))

Axioms and theorems

Before doing so, let us clarify two common terms (that you can also find in many scientific papers)

- **Axioms**: basic facts about the domain, our “initial” knowledge base
- **Theorems**: statements that are logically derived from axioms

since everything settles on axioms.

Anyway, how do we arrive at a systematic inference approach?

Since we can now talk about sets of objects, we must be able to talk about smaller, more specific sets:

SUBST

- SUBST replaces one or more variables with something else in a sentence
- For example:
 - $\text{SUBST}(\{x/\text{John}\}, \text{IsHealthy}(x) \Rightarrow \text{NOT}(\text{HasACold}(x)))$ gives us
 - $\text{IsHealthy}(\text{John}) \Rightarrow \text{NOT}(\text{HasACold}(\text{John}))$

Instantiating quantifiers

- From
- **for all x : a**
- we can obtain
- **SUBST($\{x/g\}$, a)**

- From
- **there exists x : a**
- we can obtain
- **SUBST($\{x/k\}$, a)**
- where k is a constant that does not appear elsewhere in the knowledge base (**Skolem constant**)
- **Don't need original sentence anymore**

If we want to instantiate existentials after universals, we have to be a bit more careful

- for all x : there exists y : $\text{IsParentOf}(y, x)$
- WRONG: for all x : $\text{IsParentOf}(k, x)$
- RIGHT: for all x : $\text{IsParentOf}(k(x), x)$
- Introduces a new function (Skolem function)
 - ... again, assuming k has not been used previously

Now, we can start to develop our systematic inference approach

Generalized modus ponens

- for all x: Loves(John, x)
 - John loves every thing
- for all y: (Loves(y, Jane) => FeelsAppreciatedBy(Jane, y))
 - Jane feels appreciated by every thing that loves her
- Can infer from this:
- FeelsAppreciatedBy(Jane, John)

- Here, we used the substitution {x/Jane, y/John}
 - Note we used different variables for the different sentences
- General UNIFY algorithms for finding a good substitution
(not touched here)

UNIFY: Keeping things as general as possible in unification

- Consider EdibleByWith
 - e.g., EdibleByWith(Soup, John, Spoon) – John can eat soup with a spoon
- for all x: for all y: EdibleByWith(Bread, x, y)
 - Anything can eat bread with anything
- for all u: for all v: (EdibleByWith(u, v, Spoon) => CanBeServedInBowlTo(u,v))
 - Anything that is edible with a spoon by something can be served in a bowl to that something
- Substitution: {x/z, y/Spoon, u/Bread, v/z}
- Gives: for all z: CanBeServedInBowlTo(Bread, z)
- Alternative substitution {x/John, y/Spoon, u/Bread, v/John} would only have given CanBeServedInBowlTo(Bread, John), which is not as general

Resolution for first-order logic

- for all x: (NOT(Knows(John, x)) OR IsMean(x) OR Loves(John, x))
 - John loves everything he knows, with the possible exception of mean things
- for all y: (Loves(Jane, y) OR Knows(y, Jane))
 - Jane loves everything that does not know her
- What can we unify? What can we conclude?
- Use the substitution: {x/Jane, y/John}
- Get: IsMean(Jane) OR Loves(John, Jane) OR Loves(Jane, John)
- Complete (i.e., if not satisfiable, will find a proof of this), **if we can remove literals that are duplicates after unification**
 - Also need to put everything in **canonical form** first

Canonical Form (First-order CNF)

- Convert to **Negation Normal Form**, i.e., negation symbols only occur immediately before predicate symbols
- If variable names are used twice within scopes of different quantifiers, **rename** one of the such that the name is not used anywhere else.
- **Skolemize the statements**: Move quantifiers out so that we get ForAll X₁, X₂,... Exists Y₁, Y₂, ... Formula. Quantifiers only appear in the initial prefix but never inside not, AND, or OR. Here, we have to make use of skolem functions! That is we preserve satisfiability rather than equivalence
- **Drop all universal quantifiers**
- **Distribute Ors inwards over ANDs**

Notes on inference in first-order logic

- Deciding whether a sentence is entailed is **semidecidable**: there are algorithms that will eventually produce a proof of any entailed sentence
- It is **not decidable**: we cannot always conclude that a sentence is not entailed

(Extremely informal statement of) Gödel's Incompleteness Theorem

- First-order logic is not rich enough to model basic arithmetic
- For any consistent system of axioms that is rich enough to capture basic arithmetic (in particular, mathematical induction), there exist true sentences that cannot be proved from those axioms
- This shows that we have not the chance to prove every truth in the “universe”!

In “Dangerous Knowledge”, David Malone looks at four brilliant mathematicians - Georg Cantor, Ludwig Boltzmann, **Kurt Gödel** and Alan Turing - whose genius has profoundly affected us ...

The screenshot shows the BBC Two website page for the documentary. At the top, the BBC Two logo is displayed. To the right, there is information about the broadcast: "Last on TWO Wed 11 Jun 2008 23:30 BBC TWO EXCEPT EAST & YORKSHIRE". Below this, a summary of the program is given: "Are there, in the matrix of rules which make reality work, some questions too dangerous to ask? David Malone looks at the tragic stories of 4 great thinkers whose obsessive pursuit of the deepest knowledge led them to madness and suicide." A duration of "1 hour, 30 minutes" is also mentioned. Further down, there is a "Credit" section listing "Presenter David Malone" and a "Broadcasts" section showing three broadcast times: "Wed 8 Aug 2007 22:05 BBC FOUR", "Thu 9 Aug 2007 02:05 BBC FOUR", and "Thu 16 Aug 2007 01:00 BBC FOUR".



<https://www.dailymotion.com/video/xdoj8c>

A more challenging exercise

- Suppose:
 - There are exactly 3 objects in the world,
 - If x is the spouse of y , then y is the spouse of x (spouse is a function, i.e., everything has a spouse)
- Prove:
 - Something is its own spouse

More challenging exercise

- there exist x, y, z : $(\text{NOT}(x=y) \text{ AND } \text{NOT}(x=z)) \text{ AND } \text{NOT}(y=z)$
- for all w, x, y, z : $(w=x \text{ OR } w=y \text{ OR } w=z \text{ OR } x=y \text{ OR } x=z \text{ OR } y=z)$
- for all x, y : $((\text{Spouse}(x)=y) \Rightarrow (\text{Spouse}(y)=x))$
- for all x, y : $((\text{Spouse}(x)=y) \Rightarrow \text{NOT}(x=y))$ (*for the sake of contradiction*)
- *Try to do this on the board...*

Umbrellas in first-order logic

- You know the following things:
 - You have exactly one other person living in your house, who is wet
 - If a person is wet, it is because of the rain, the sprinklers, or both
 - If a person is wet because of the sprinklers, the sprinklers must be on
 - If a person is wet because of rain, that person must not be carrying any umbrella
 - There is an umbrella that “lives in” your house, which is not in its house
 - An umbrella that is not in its house must be carried by some person who lives in that house
 - You are not carrying any umbrella
- **Can you conclude that the sprinklers are on?**

Applications

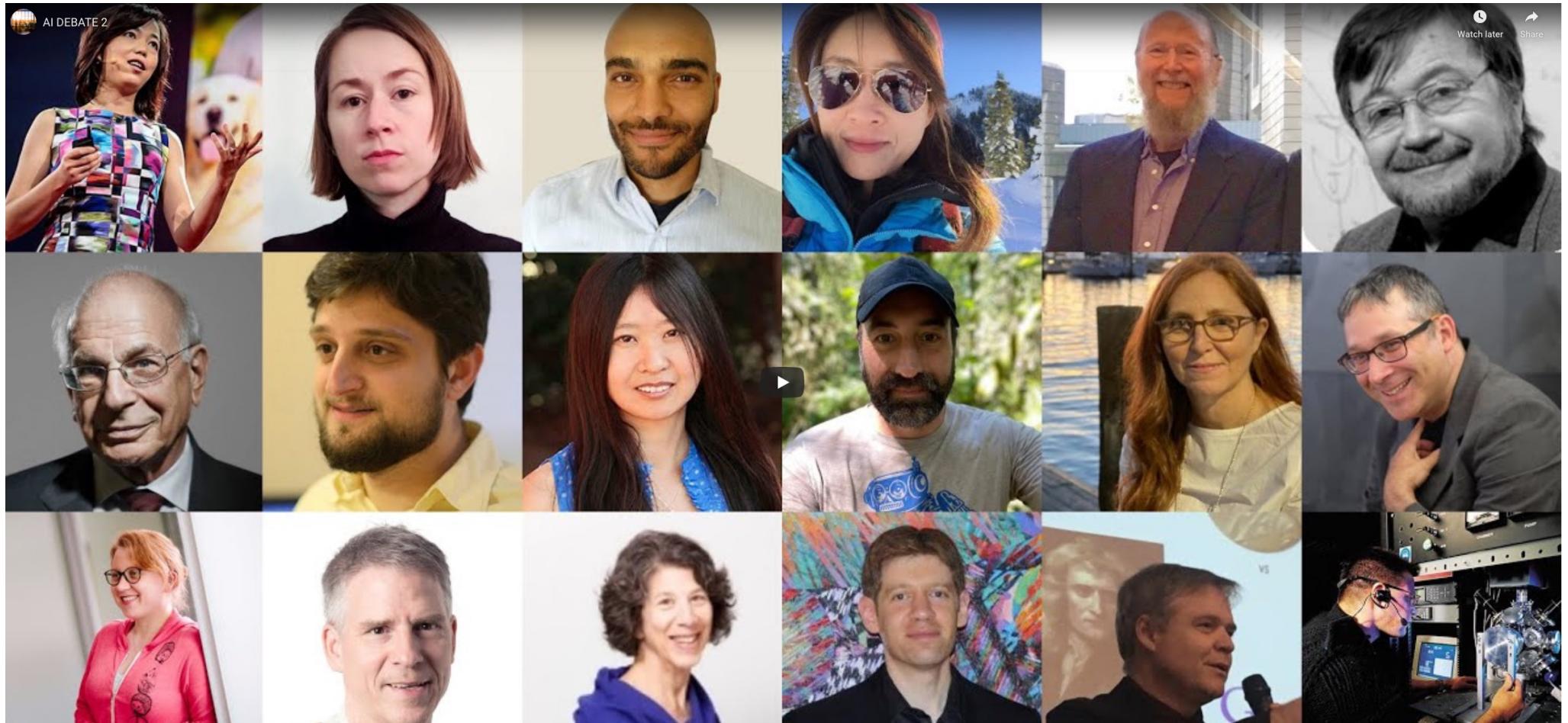
- Some serious novel mathematical results proved
- Verification of hardware and software
 - Prove outputs satisfy required properties for all inputs
- Synthesis of hardware and software
 - Try to prove that there exists a program satisfying such and such properties, **in a constructive way**

Something we cannot do in first-order logic

- We are **not** allowed to reason in general about relations and functions
- The following would correspond to **higher-order logic** (which is more powerful):
- “If John is Jack’s roommate, then any property of John is also a property of Jack’s roommate”
- $(\text{John} = \text{Roommate}(\text{Jack})) \Rightarrow \text{for all } p: (p(\text{John}) \Rightarrow p(\text{Roommate}(\text{Jack})))$
- “If a property is inherited by children, then for any thing, if that property is true of it, it must also be true for any child of it”
- $\text{for all } p: (\text{IsInheritedByChildren}(p) \Rightarrow (\text{for all } x, y: ((\text{IsChildOf}(x, y) \text{ AND } p(y)) \Rightarrow p(x))))$

AI Debate 2

<https://montrealartificialintelligence.com/aidebate2.html>

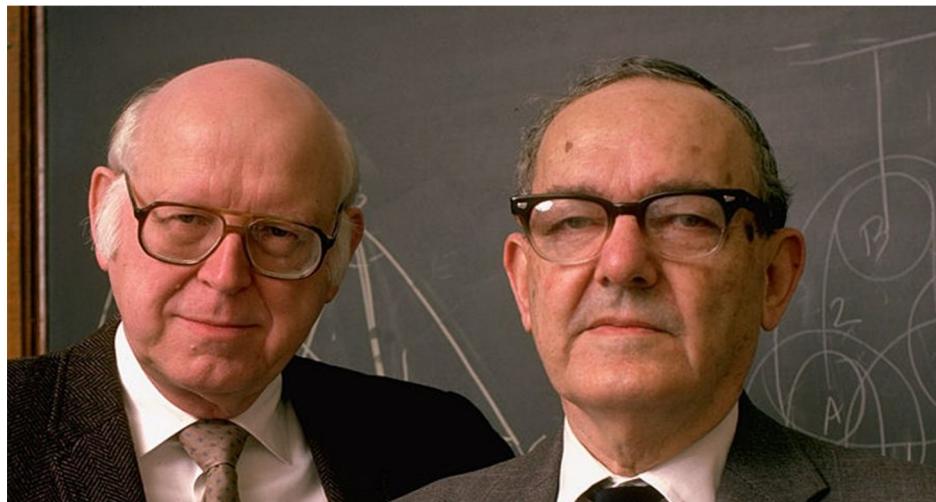


How to get logical reasoning and deep learning "under one umbrella"?

Symbols

“A physical symbol system has the necessary and sufficient means for general intelligent action.”

- Allen Newell & Herbert Simon
ACM Turing Awardees



Neurons

“Symbols are Luminiferous Aether of AI”

- Geoff Hinton
ACM Turing Awardee



Let us focus on the computational and algorithmic level not just the implementational one

Three levels of description (*David Marr, 1982*)

Computational

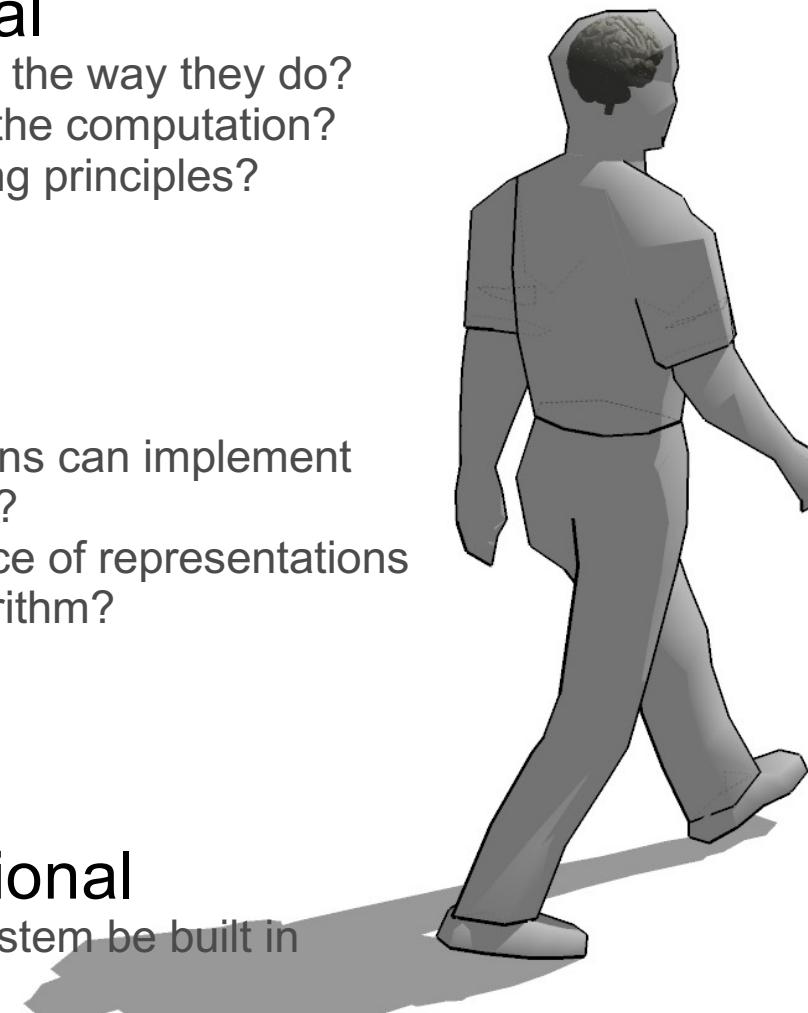
Why do things work the way they do?
What is the goal of the computation?
What are the unifying principles?

Algorithmic

What representations can implement such computations?
How does the choice of representations determine the algorithm?

Implementational

How can such a system be built in hardware?
How can neurons carry out the computations?



maximize:

$$R_t = r_{t+1} + r_{t+2} + \cdots + r_T$$

