

knowledge-based agent



# 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
- Or at the **implementation level**
  - i.e., data structures in KB and algorithms that manipulate them
  -

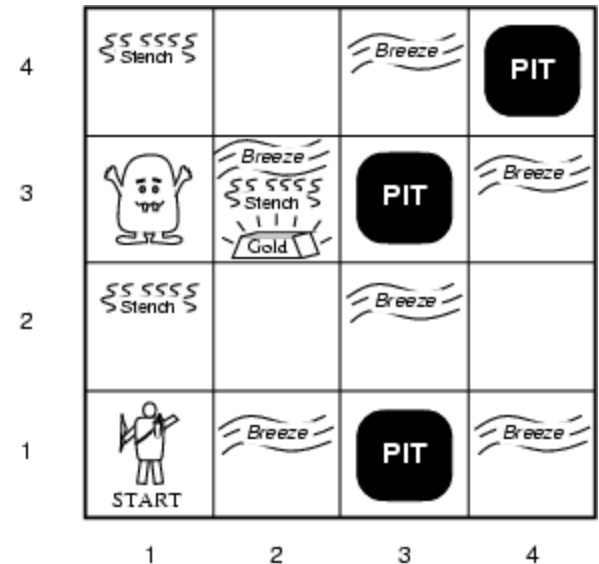
# A simple knowledge-based agent

```
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 ← ASK(KB, MAKE-ACTION-QUERY(t))  
  TELL(KB, MAKE-ACTION-SENTENCE(action, t))  
  t ← t + 1  
  return action
```

- The agent must be able to:
  - - Represent states, actions, etc.
    - 
    - Incorporate new percepts
    - 
    - Update internal representations of the world

# 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



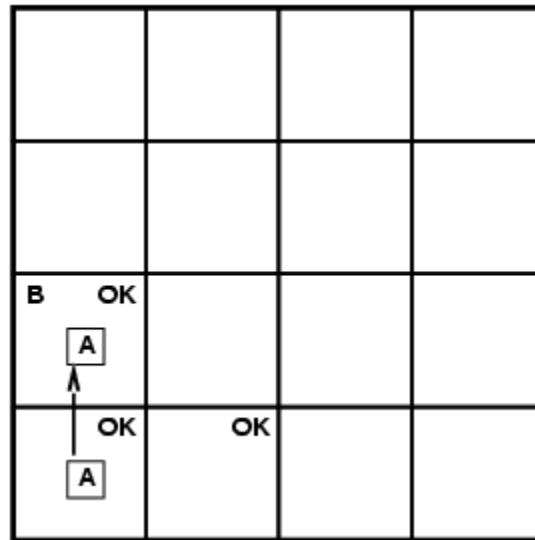
# Wumpus world characterization

- 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
- 
- Single-agent? Yes – Wumpus is essentially a natural feature

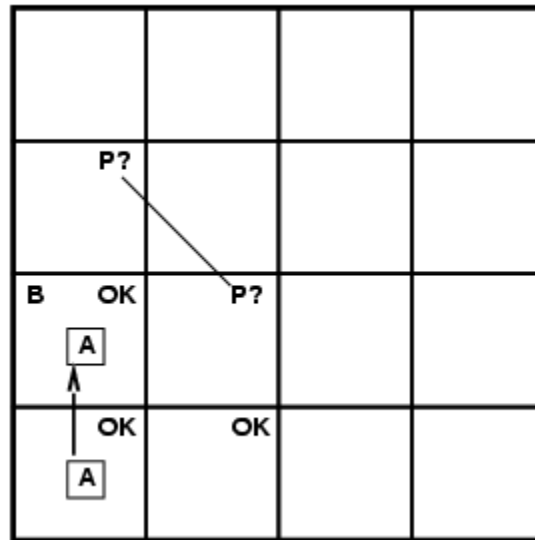
# Exploring a wumpus world

OK			
OK A	OK		

# Exploring a wumpus world

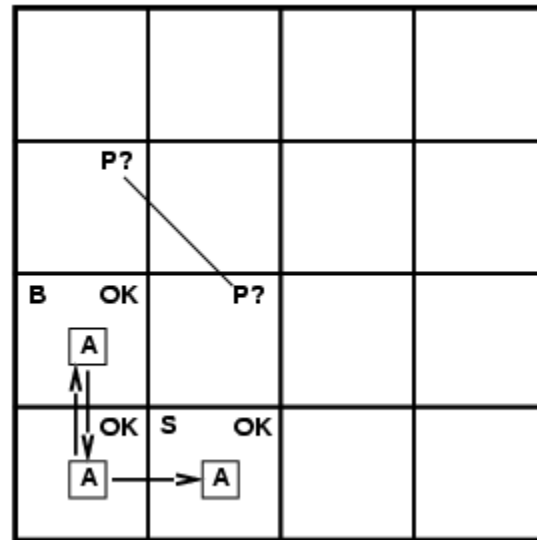


# Exploring a wumpus world

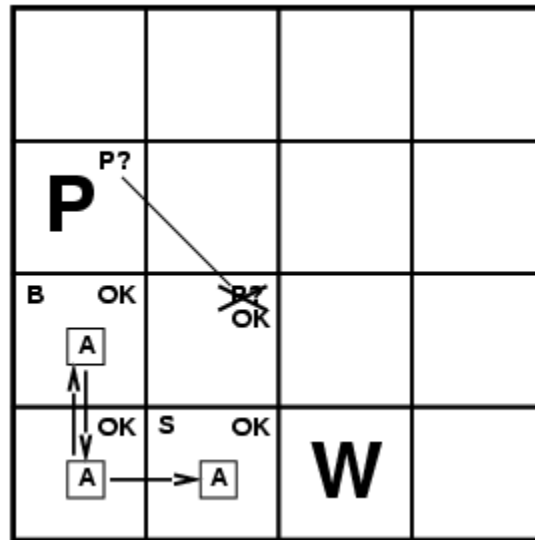




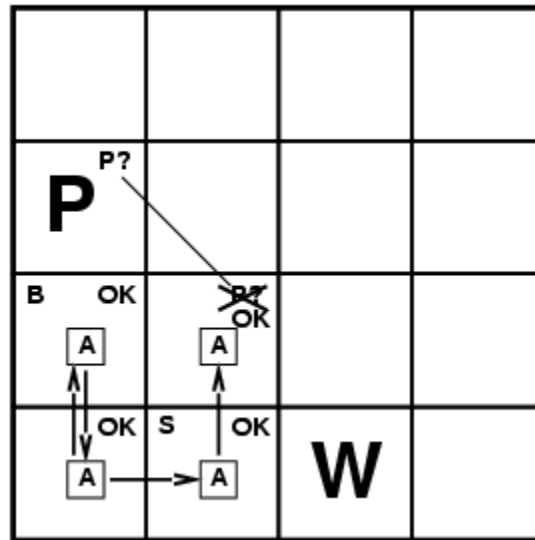
# Exploring a wumpus world



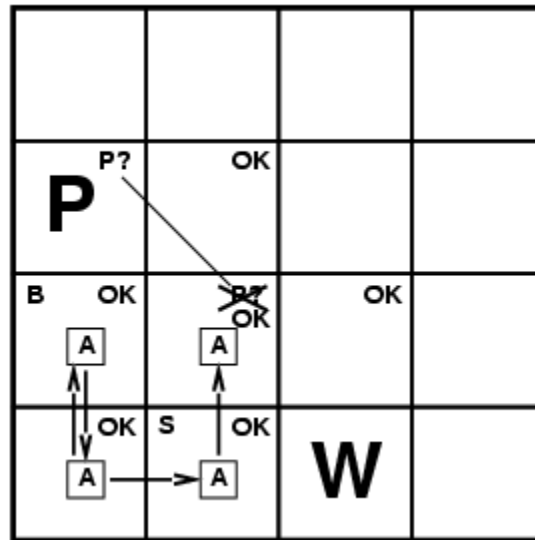
# Exploring a wumpus world



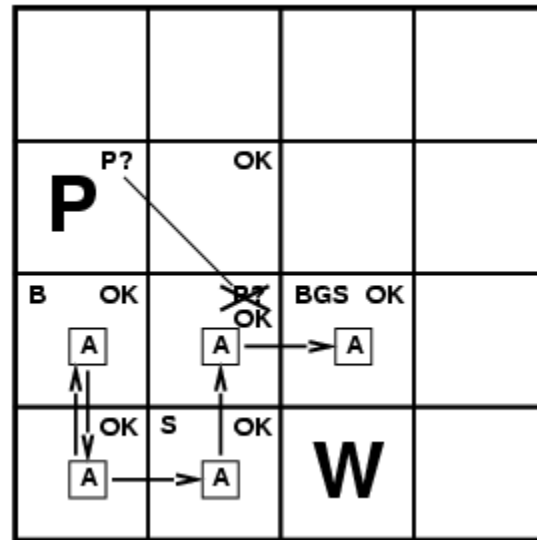
# Exploring a wumpus world



# Exploring a wumpus world



# Exploring a wumpus world



# Logic in general

- **Logics** are formal languages for representing information such that conclusions can be drawn
- 
- **Syntax** defines the sentences in the language
- 
- **Semantics** define the "meaning" of sentences;
- i.e., define **truth** of a sentence in a world
- E.g., the language of arithmetic
- $x+2 \geq y$  is a sentence;  $x^2+y > \{ \}$  is not a sentence
- $x+2 \geq y$  is true iff the number  $x+2$  is no less than the number  $y$
- $x+2 \geq y$  is true in a world where  $x = 7, y = 1$ 
  - $x+2 \geq y$  is false in a world where  $x = 0, y = 6$
  -

# Formal languages

Language	Ontology	Epistemology
<b>1- certain languages</b>		
Propositional logic	Facts	True/False/Unknown
First-order logic	Facts, objects, relations	True/False/Unknown
Temporal logic	Facts, objects, relations, time	True/False/Unknown
<b>2- uncertainty languages</b>		
Probability theory	Facts	Degree of belief [0,1]
Fuzzy logic	Facts	Degree of truth [0,1]

# Entailment

- **Entailment** means that one thing **follows from** another:

- 

$$KB \models \alpha$$

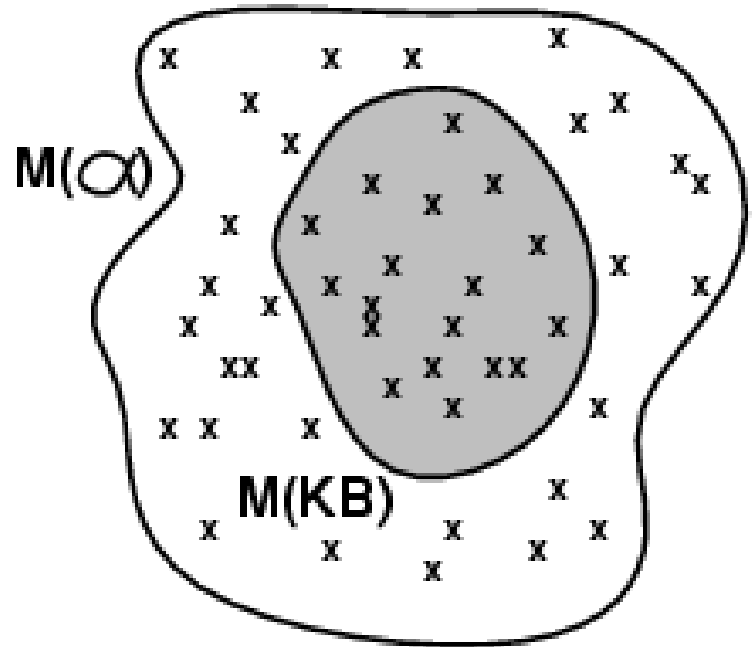
Knowledge base  $KB$  entails sentence  $\alpha$  if and only if  $\alpha$  is true in all worlds where  $KB$  is true

- E.g., the KB containing “the Giants won” and “the Reds won” entails “Either the Giants won or the Reds won”
- E.g.,  $x+y = 4$  entails  $4 = x+y$
- Entailment is a relationship between sentences (i.e., **syntax**) that is based on **semantics**



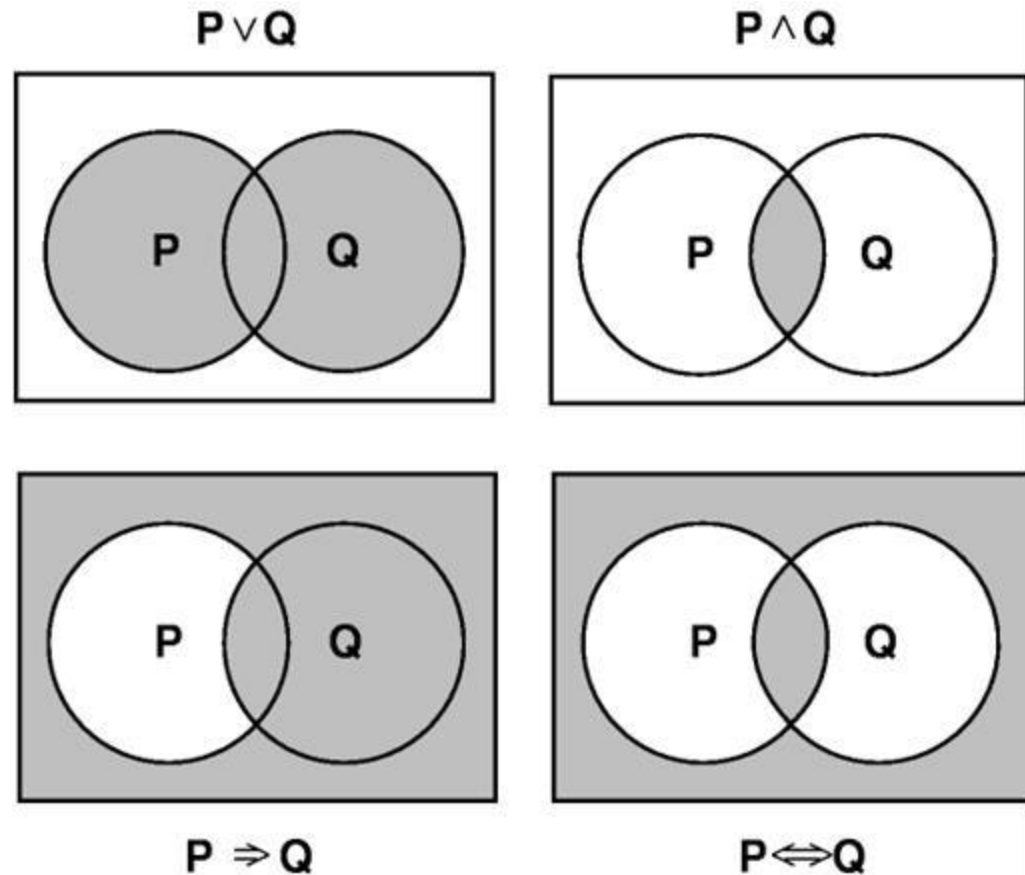
# Models

- Logicians typically think in terms of **models**, which are formally structured worlds with respect to which truth can be evaluated
- 
- We say  $m$  is a **model** of a sentence  $\alpha$  if  $\alpha$  is true in  $m$
- $M(\alpha)$  is the set of all models of  $\alpha$
- 
- Then  $KB \models \alpha$  iff  $M(KB) \subseteq M(\alpha)$
- - E.g.  $KB =$  Giants won  $\alpha =$  Giants won
  -



# Models

- We could define the meaning of a sentence by means of set operations on sets of models.
- For example, the set of models of  $P \wedge Q$  is the intersection of models of  $P$  and models of  $Q$ .

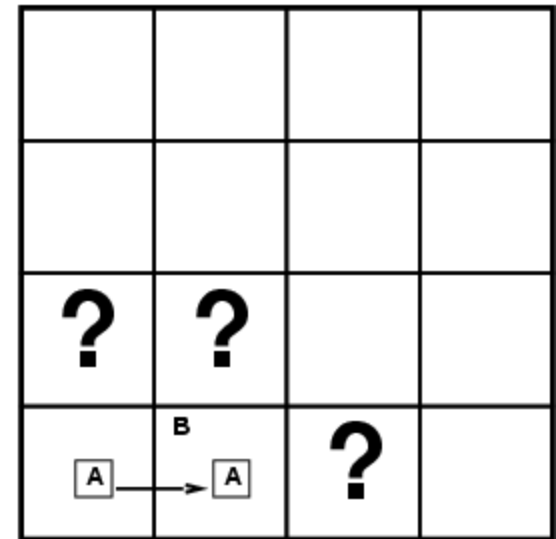


# Entailment in the wumpus world

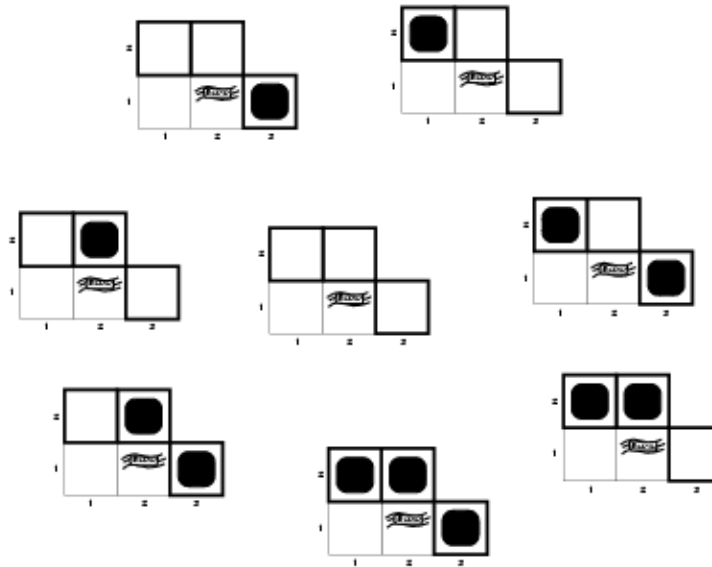
Situation after detecting nothing in [1,1], moving right, breeze in [2,1]

Consider possible models for *KB* assuming only pits

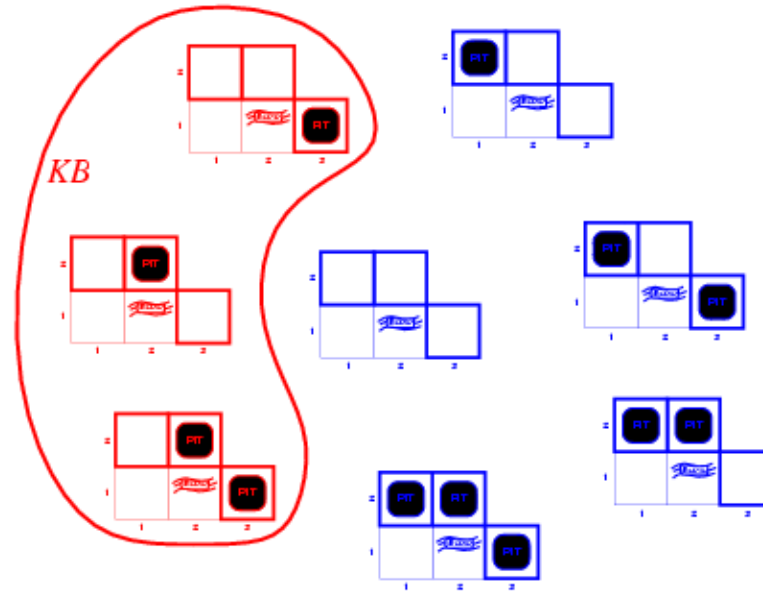
3 Boolean choices  $\Rightarrow$  8 possible models



# Wumpus models

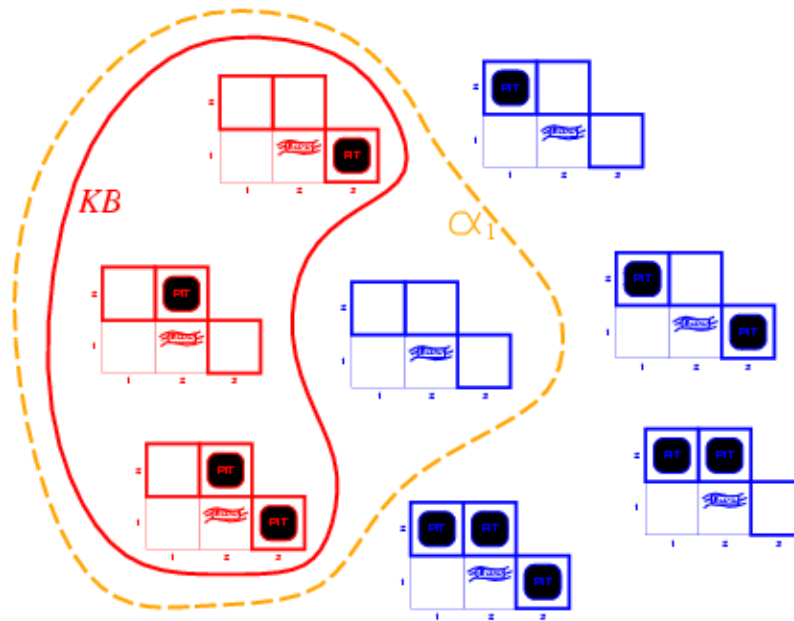


# Wumpus models



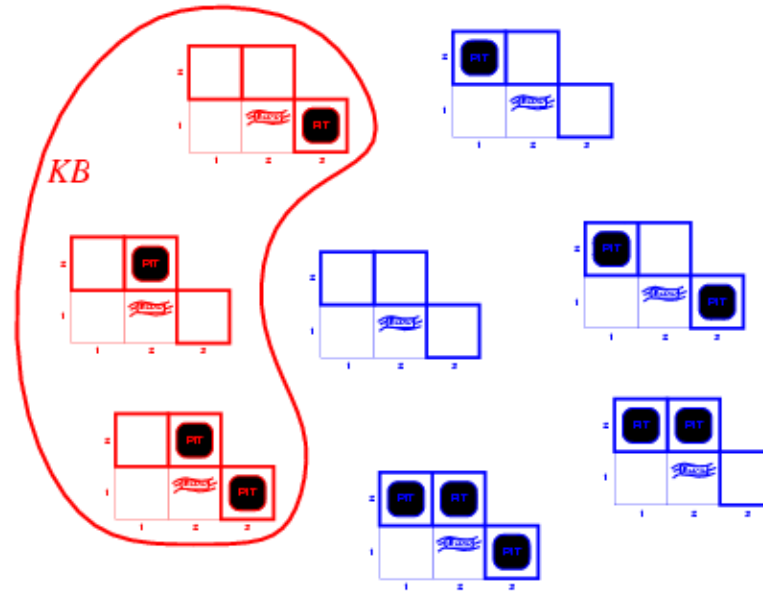
- $KB =$  wumpus-world rules + observations
-

# Wumpus models



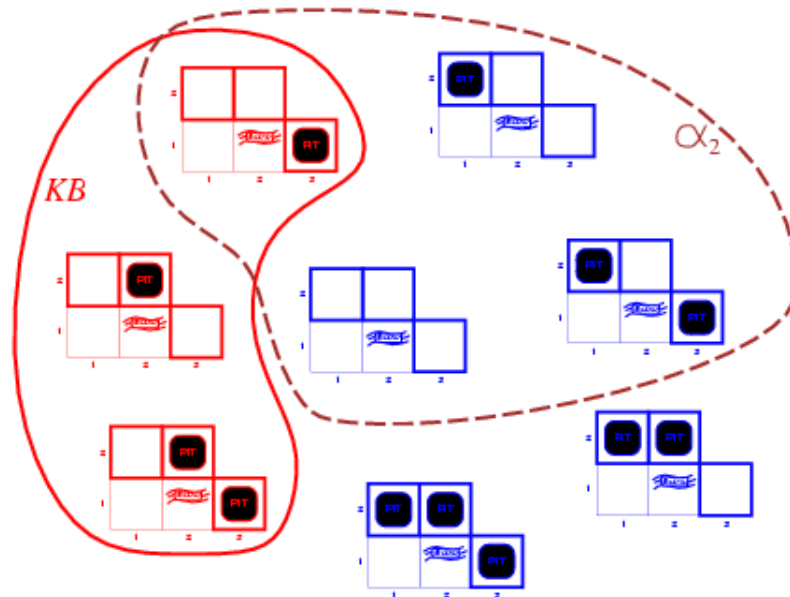
- $KB$  = wumpus-world rules + observations
- $\alpha_1 = "[1,2]$  is safe",  $KB \models \alpha_1$ , proved by model checking
-

# Wumpus models



- $KB = \text{wumpus-world rules} + \text{observations}$

# Wumpus models



- $KB$  = wumpus-world rules + observations
- $\alpha_2$  = "[2,2] is safe",  $KB \not\models \alpha_2$
-



# Inference

- $KB \vdash_i \alpha$  = sentence  $\alpha$  can be derived from  $KB$  by procedure  $i$
- 
- **Soundness:**  $i$  is sound if whenever  $KB \vdash_i \alpha$ , it is also true that  $KB \models \alpha$
- 
- **Completeness:**  $i$  is complete if whenever  $KB \models \alpha$ , it is also true that  $KB \vdash_i \alpha$
- 
- The procedure will answer any question whose answer follows from what is known by the  $KB$ .
-

# Propositional logic: Syntax

- Propositional logic is the simplest logic – illustrates basic ideas
- 
- The proposition symbols  $P_1, P_2$  etc are sentences
  - If  $S$  is a sentence,  $\neg S$  is a sentence (**negation**)
  - 
  - If  $S_1$  and  $S_2$  are sentences,  $S_1 \wedge S_2$  is a sentence (**conjunction**)
  - 
  - If  $S_1$  and  $S_2$  are sentences,  $S_1 \vee S_2$  is a sentence (**disjunction**)
  - 
  - If  $S_1$  and  $S_2$  are sentences,  $S_1 \Rightarrow S_2$  is a sentence (**implication**)
  - 
  - If  $S_1$  and  $S_2$  are sentences,  $S_1 \Leftrightarrow S_2$  is a sentence (**biconditional**)
  -

# Propositional logic: Syntax

*Sentence*  $\rightarrow$  *AtomicSentence* | *ComplexSentence*

*AtomicSentence*  $\rightarrow$  **True** | **False** | *Symbol*

*Symbol*  $\rightarrow$  **P** | **Q** | **R** | ...

*ComplexSentence*  $\rightarrow$   $\neg$  *Sentence*

| ( *Sentence*  $\wedge$  *Sentence* )

| ( *Sentence*  $\vee$  *Sentence* )

| ( *Sentence*  $\Rightarrow$  *Sentence* )

| ( *Sentence*  $\Leftrightarrow$  *Sentence* )

# Propositional logic: Semantics

Each model specifies true/false for each proposition symbol

E.g.  $P_{1,2}$        $P_{2,2}$        $P_{3,1}$   
false            true            false

With these symbols, 8 possible models, can be enumerated automatically.

Rules for evaluating truth with respect to a model  $m$ :

$\neg S$	is true iff	$S$ is false	
$S_1 \wedge S_2$	is true iff	$S_1$ is true <b>and</b>	$S_2$ is true
$S_1 \vee S_2$	is true iff	$S_1$ is true <b>or</b>	$S_2$ is true
$S_1 \Rightarrow S_2$	is true iff	$S_1$ is false <b>or</b>	$S_2$ is true
i.e.,	is false iff	$S_1$ is true <b>and</b>	$S_2$ is false
$S_1 \Leftrightarrow S_2$	is true iff	$S_1 \Rightarrow S_2$ is true <b>and</b>	$S_2 \Rightarrow S_1$ is true

Simple recursive process evaluates an arbitrary sentence, e.g.,

# Truth tables for connectives

$P$	$Q$	$\neg P$	$P \wedge Q$	$P \vee Q$	$P \Rightarrow Q$	$P \Leftrightarrow Q$
<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>true</i>
<i>false</i>	<i>true</i>	<i>true</i>	<i>false</i>	<i>true</i>	<i>true</i>	<i>false</i>
<i>true</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	<i>false</i>
<i>true</i>	<i>true</i>	<i>false</i>	<i>true</i>	<i>true</i>	<i>true</i>	<i>true</i>

# 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"

- 

$$B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})$$

$$B_{2,1} \Leftrightarrow (P_{1,1} \vee P_{2,2} \vee P_{3,1})$$

# Logical equivalence

- Two sentences are **logically equivalent** iff true in same models:  $\alpha \equiv \beta$  iff  $\alpha \models \beta$  and  $\beta \models \alpha$

- $(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$  commutativity of  $\wedge$
- $(\alpha \vee \beta) \equiv (\beta \vee \alpha)$  commutativity of  $\vee$
- $((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma))$  associativity of  $\wedge$
- $((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma))$  associativity of  $\vee$
- $\neg(\neg\alpha) \equiv \alpha$  double-negation elimination
- $(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha)$  contraposition
- $(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta)$  implication elimination
- $(\alpha \Leftrightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha))$  biconditional elimination
- $\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta)$  de Morgan
- $\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta)$  de Morgan
- $(\alpha \wedge (\beta \vee \gamma)) \equiv ((\alpha \wedge \beta) \vee (\alpha \wedge \gamma))$  distributivity of  $\wedge$  over  $\vee$
- $(\alpha \vee (\beta \wedge \gamma)) \equiv ((\alpha \vee \beta) \wedge (\alpha \vee \gamma))$  distributivity of  $\vee$  over  $\wedge$

# Validity and satisfiability

A sentence is **valid** if it is true in **all** models,  
e.g., *True*,  $A \vee \neg A$ ,  $A \Rightarrow A$ ,  $(A \wedge (A \Rightarrow B)) \Rightarrow B$

Validity is connected to inference via the **Deduction Theorem**:  
 $KB \models \alpha$  if and only if  $(KB \Rightarrow \alpha)$  is valid

A sentence is **satisfiable** if it is true in **some** model  
e.g.,  $A \vee B$ ,  $C$

A sentence is **unsatisfiable** if it is true in **no** models  
e.g.,  $A \wedge \neg A$

Satisfiability is connected to inference via the following:  
 $KB \models \alpha$  if and only if  $(KB \wedge \neg \alpha)$  is unsatisfiable