

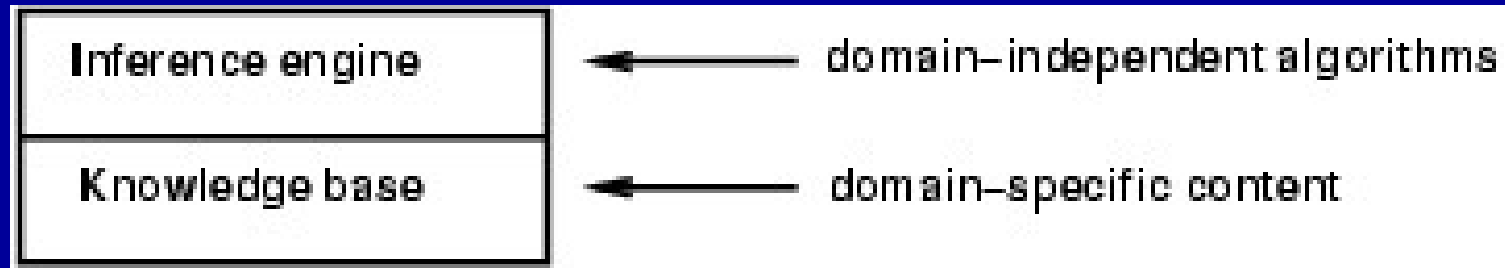
# Chapter 7

# Logical Agents

# Outline

- ❖ Knowledge-based agents
- ❖ Wumpus world
- ❖ Logic in general - models and entailment
- ❖ Propositional (Boolean) logic
- ❖ Equivalence, validity, satisfiability
- ❖ Inference rules and theorem proving

# Knowledge bases



- The central component of a knowledge-based agent is its knowledge base, or KB.
- Knowledge base(KB) = set of **sentences** in a **formal** language
- Each time the agent program is called, it does three things.
  - 1) it TELLS the knowledge base what it perceives.
  - 2) it ASKS the knowledge base what action it should perform.
  - 3) the agent program TELLS the knowledge base which action was chosen, and the agent executes the action.

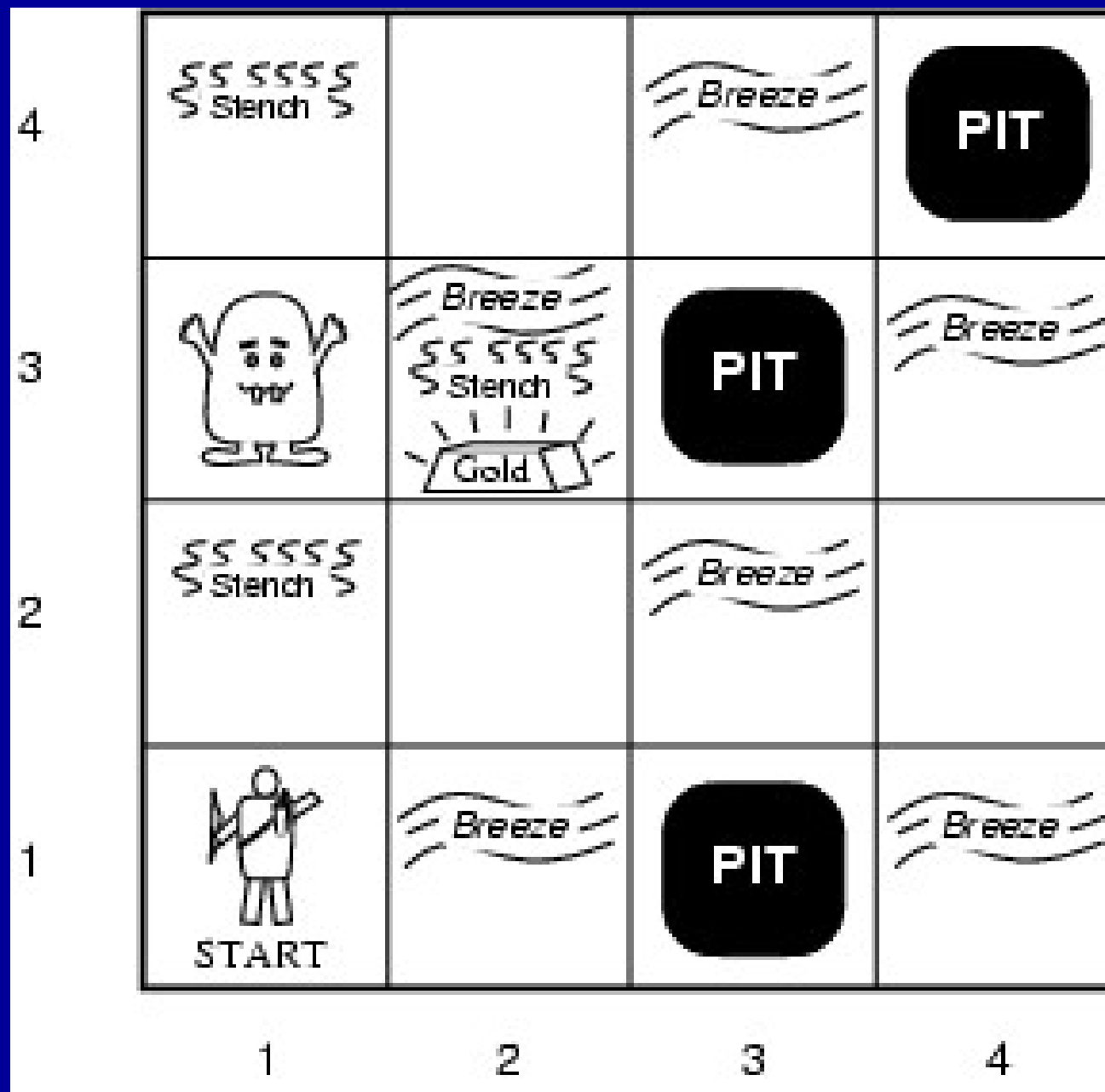
# A simple knowledge-based agent

```
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
  - Deduce hidden properties of the world
  - Deduce appropriate actions

# A Wumpus World



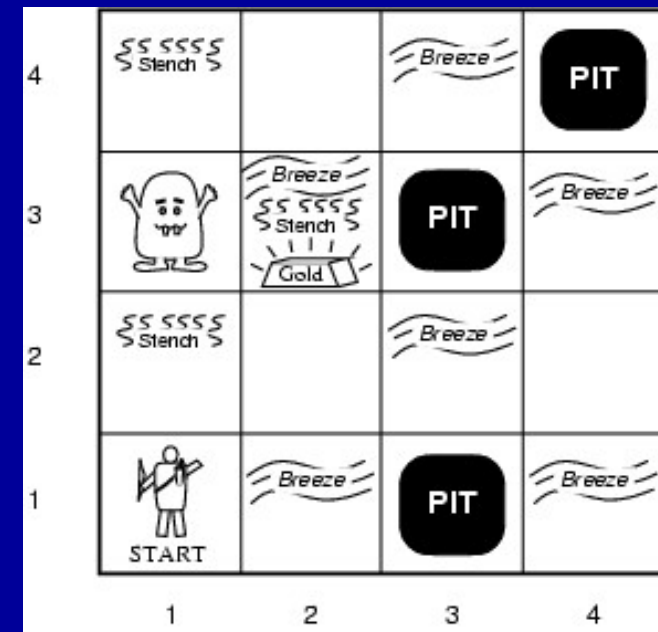
# Wumpus World PEAS description

## □ Performance measure

- gold +1000, death -1000
- -1 per step, -10 for using the arrow

## □ Environment: 4 x 4 grid of rooms

- 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 (shot Wumpus)

□ **Actuators:** Left turn, Right turn, Forward, Grab, Release, Shoot

# Wumpus world characterization

☐ Fully Observable

☐ Deterministic

☐ Episodic

☐ Static

☐ Discrete

☐ Single-agent?

# Wumpus world characterization

- ☐ Fully Observable No – only **local** perception
- ☐ Deterministic
- ☐ Episodic
- ☐ Static
- ☐ Discrete
- ☐ Single-agent?



# Wumpus world characterization

- ☐ Fully Observable No – only **local** perception
- ☐ Deterministic Yes – outcomes exactly specified
- ☐ Episodic
- ☐ Static
- ☐ Discrete
- ☐ Single-agent?

# Wumpus world characterization

- ☐ Fully Observable No – only **local** perception
- ☐ Deterministic Yes – outcomes exactly specified
- ☐ Episodic No – sequential at the level of actions
- ☐ Static
- ☐ Discrete
- ☐ Single-agent?

# 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
- ☐ Single-agent?

# 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?

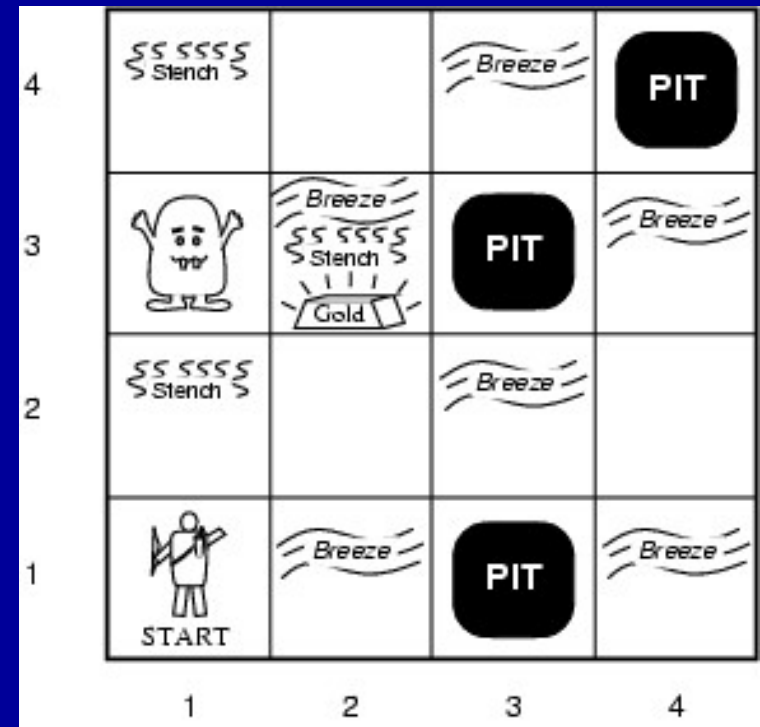
# Wumpus world characterization

- ☐ 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

# Wumpus World

## □ Percepts given to the agent

1. Stench
2. Breeze
3. Glitter/Gold
4. Bump (ran into a wall)
5. Scream (wumpus has been hit by arrow)



- ## □ Principle Difficulty: agent is initially ignorant of the configuration of the environment – going to have to reason to figure out where the gold is without getting killed!

# Exploring the Wumpus World

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 OK	2,2	3,2	4,2
1,1 A OK	2,1 OK	3,1	4,1

(a)

**A** = Agent  
**B** = Breeze  
**G** = Glitter, Gold  
**OK** = Safe square  
**P** = Pit  
**S** = Stench  
**V** = Visited  
**W** = Wumpus

**Initial situation:**

**Agent in 1,1 and percept is [None, None, None, None, None]**

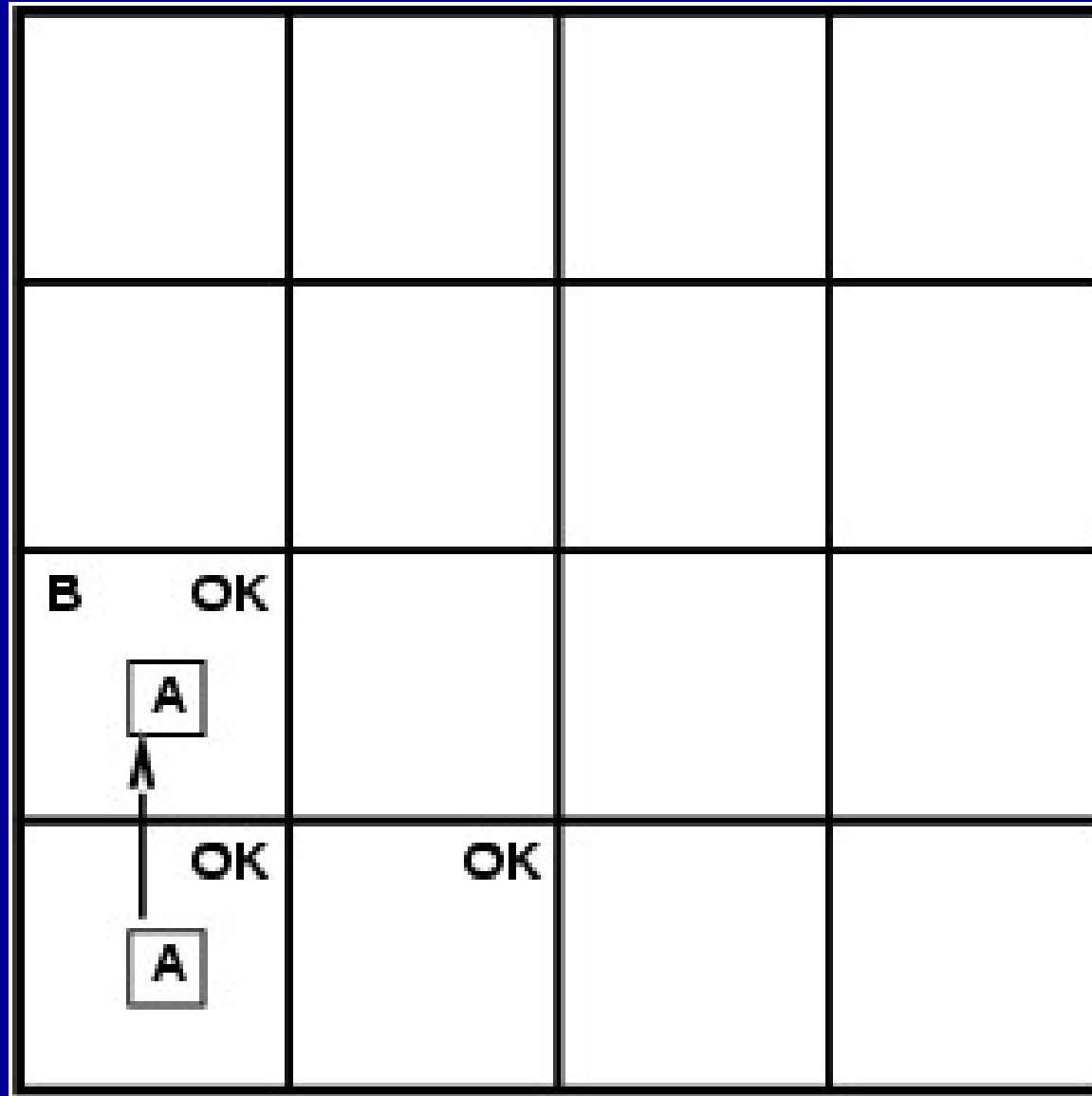
**From this the agent can infer the neighboring squares are safe (otherwise there would be a breeze or a stench)**

# Exploring a wumpus world

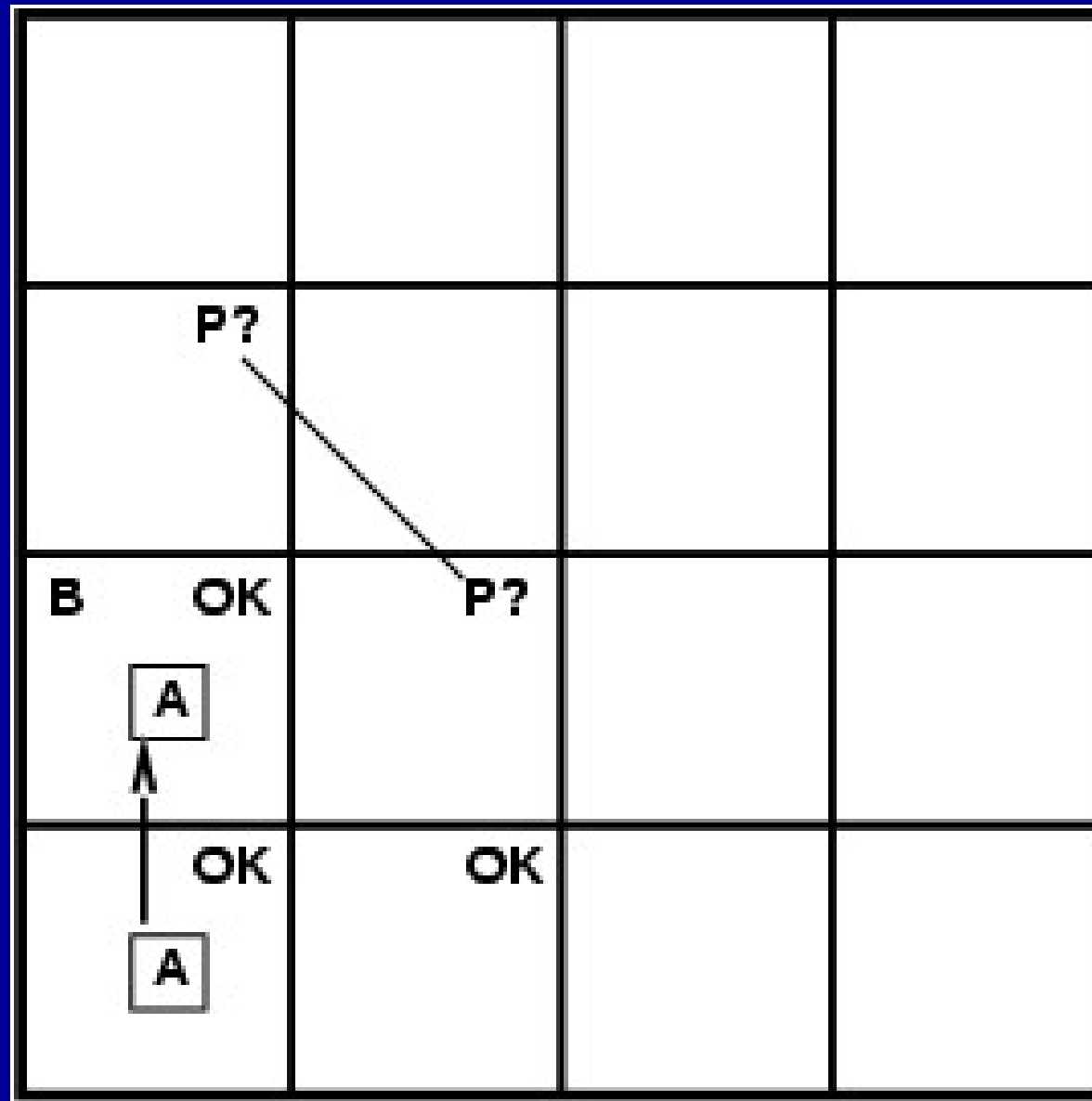
OK			
OK <div>A</div>	OK		



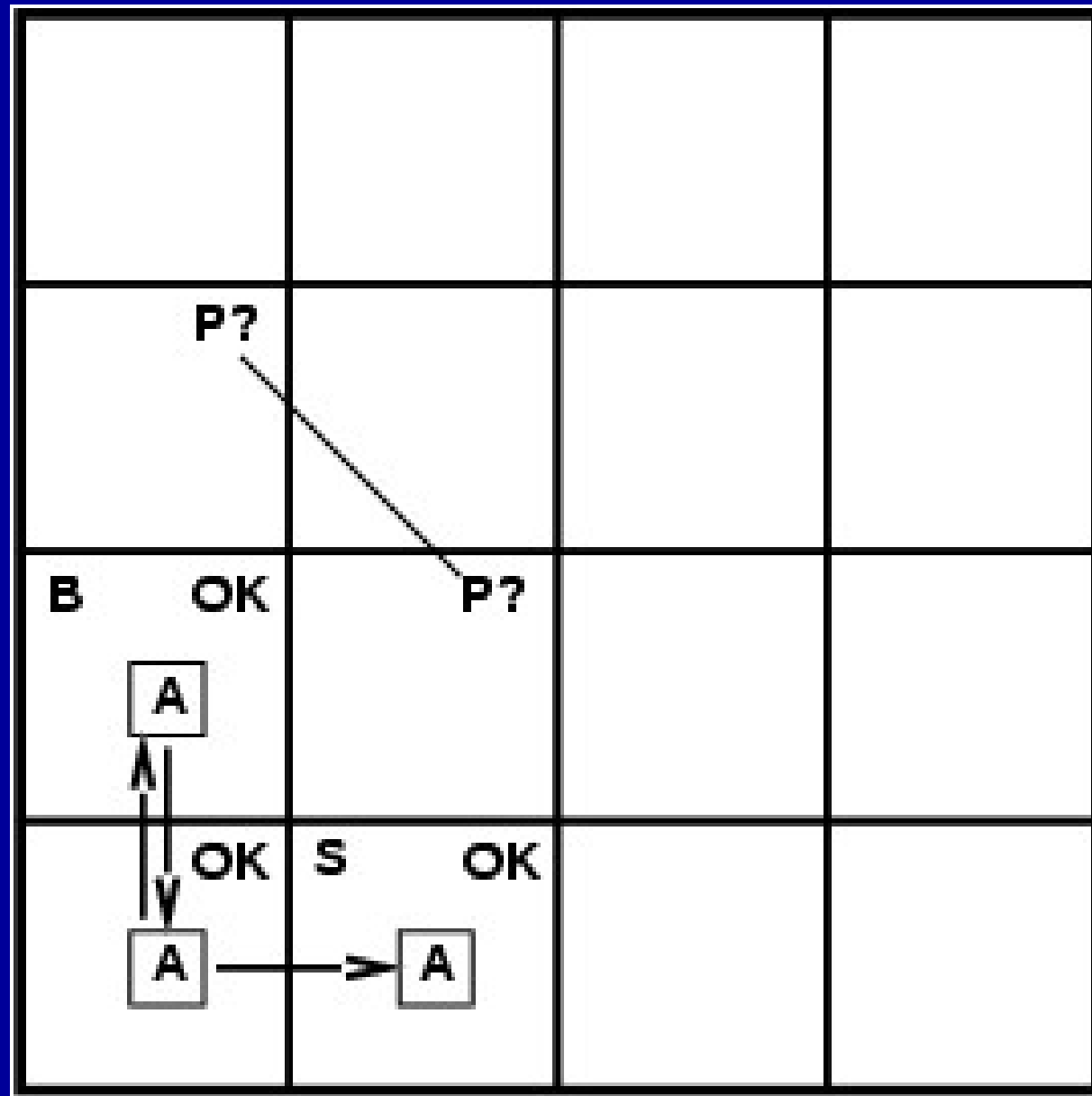
# Exploring a wumpus world



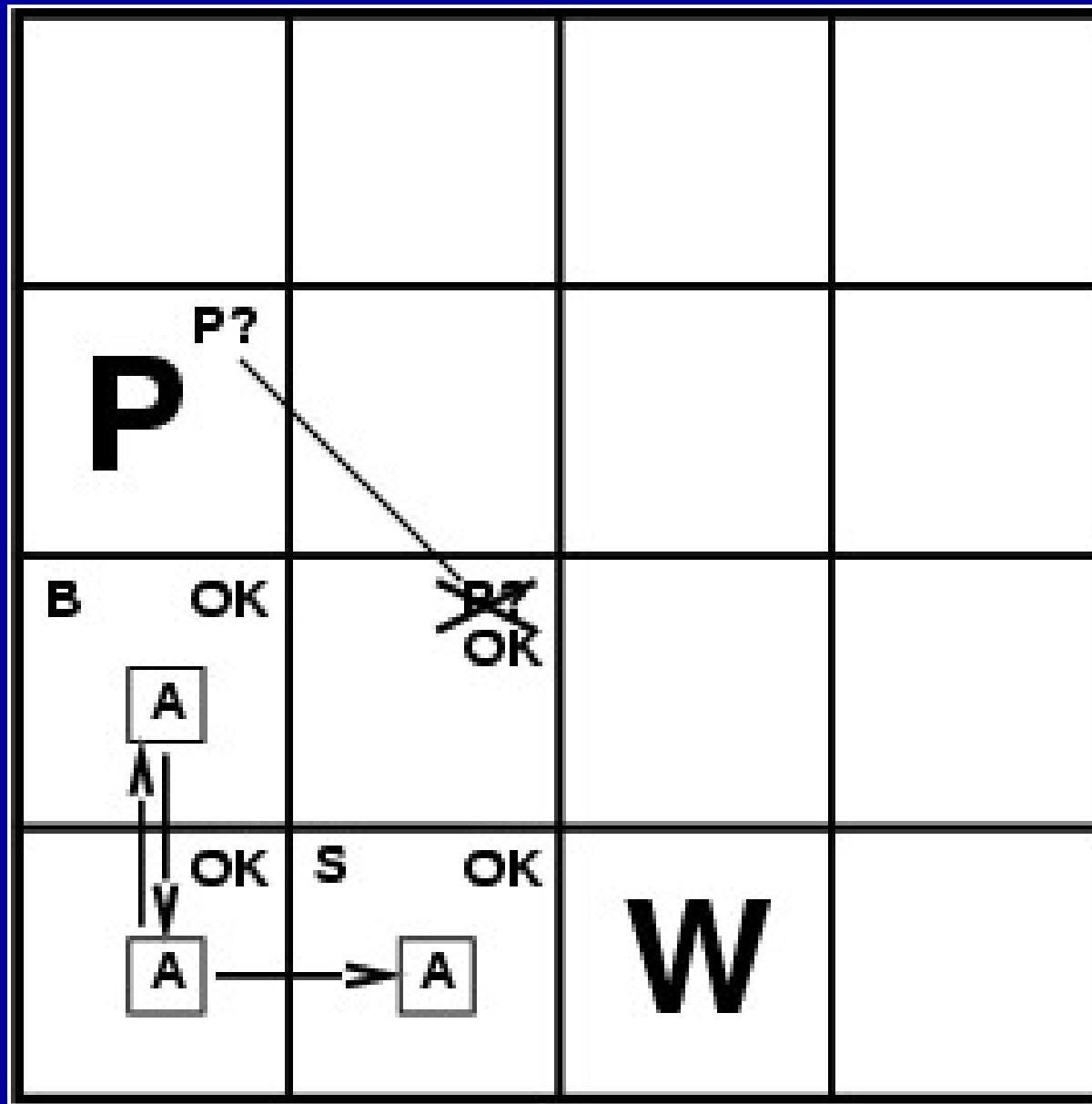
# Exploring a wumpus world



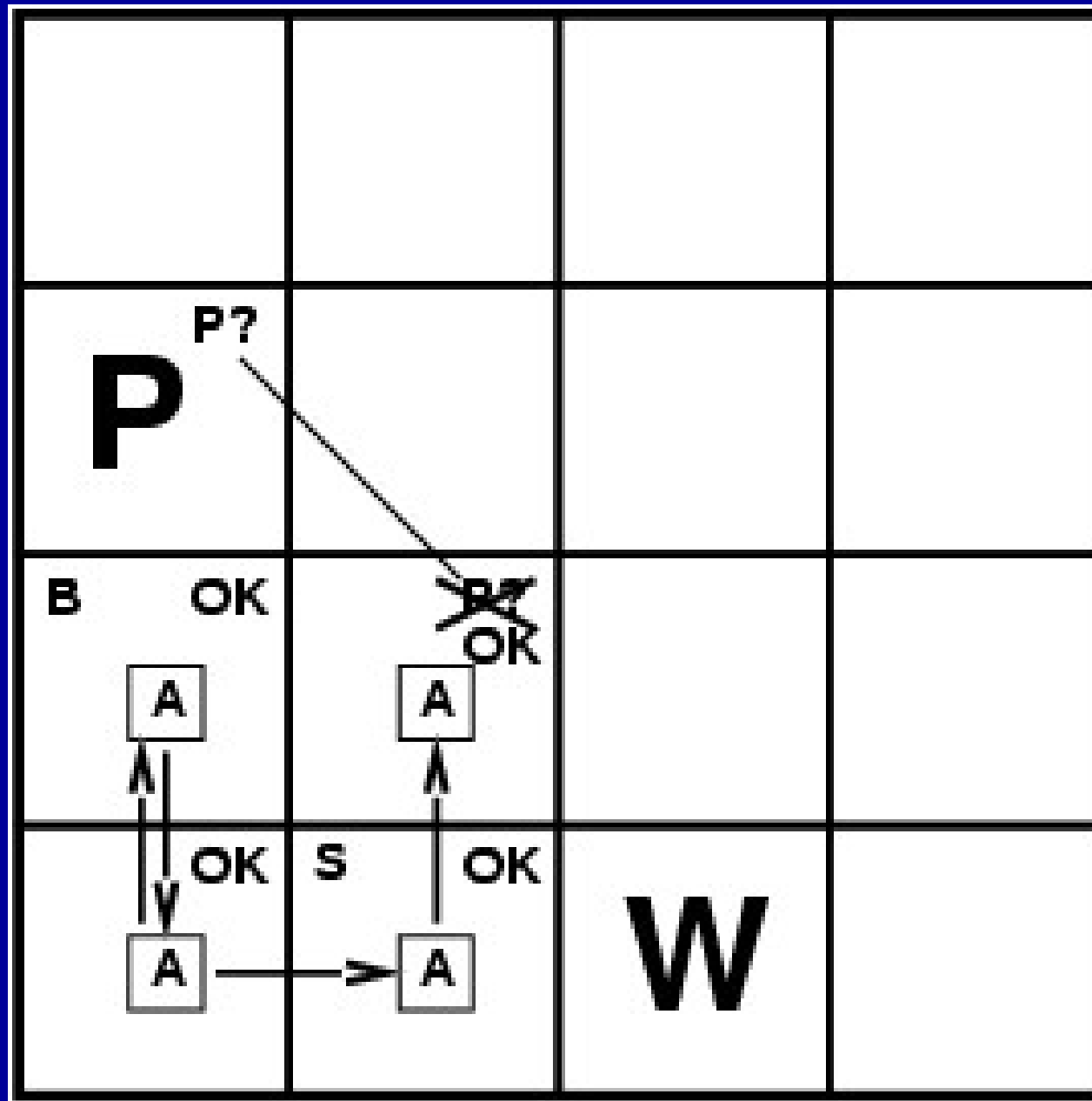
# Exploring a wumpus world



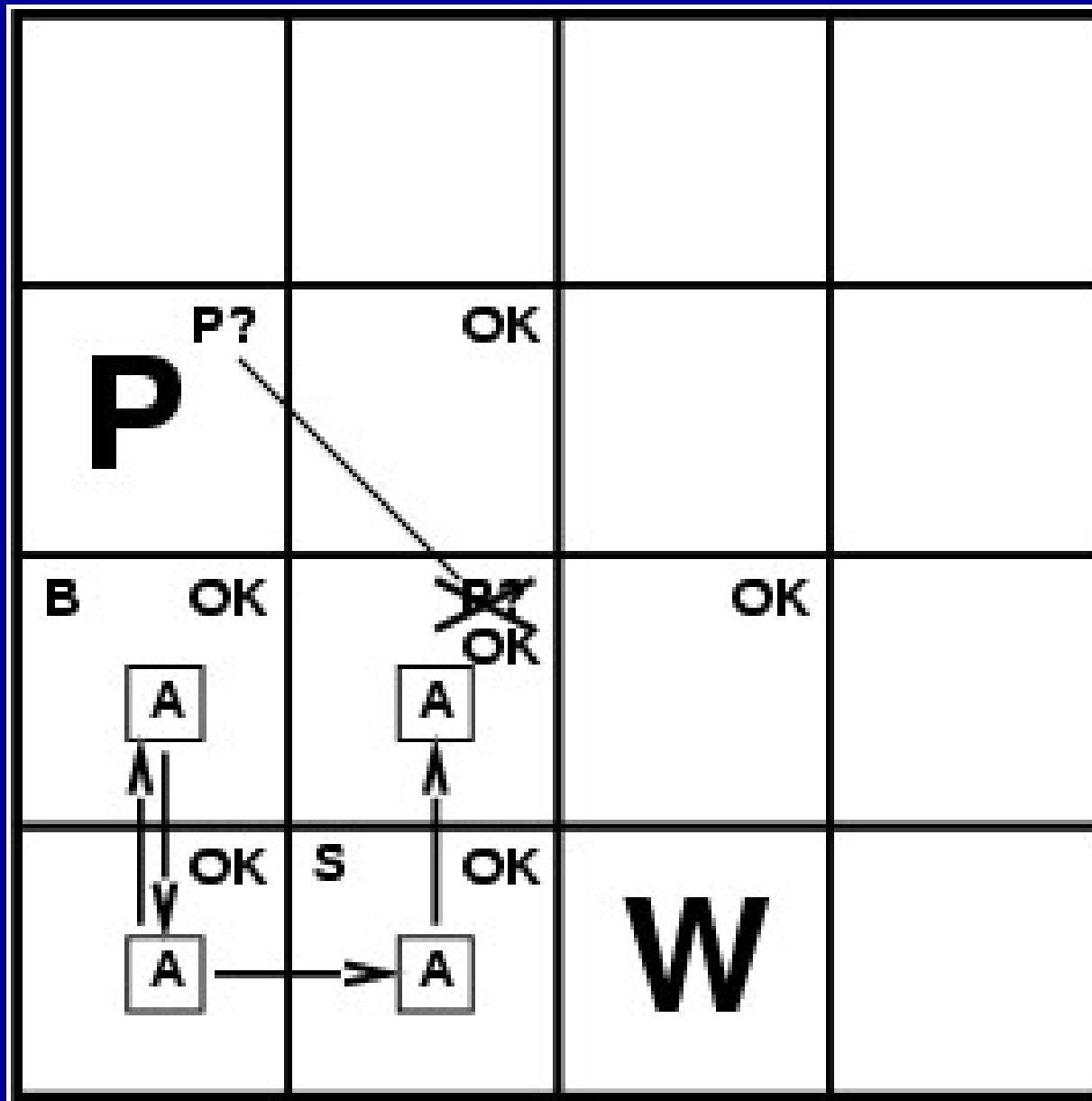
# Exploring a wumpus world



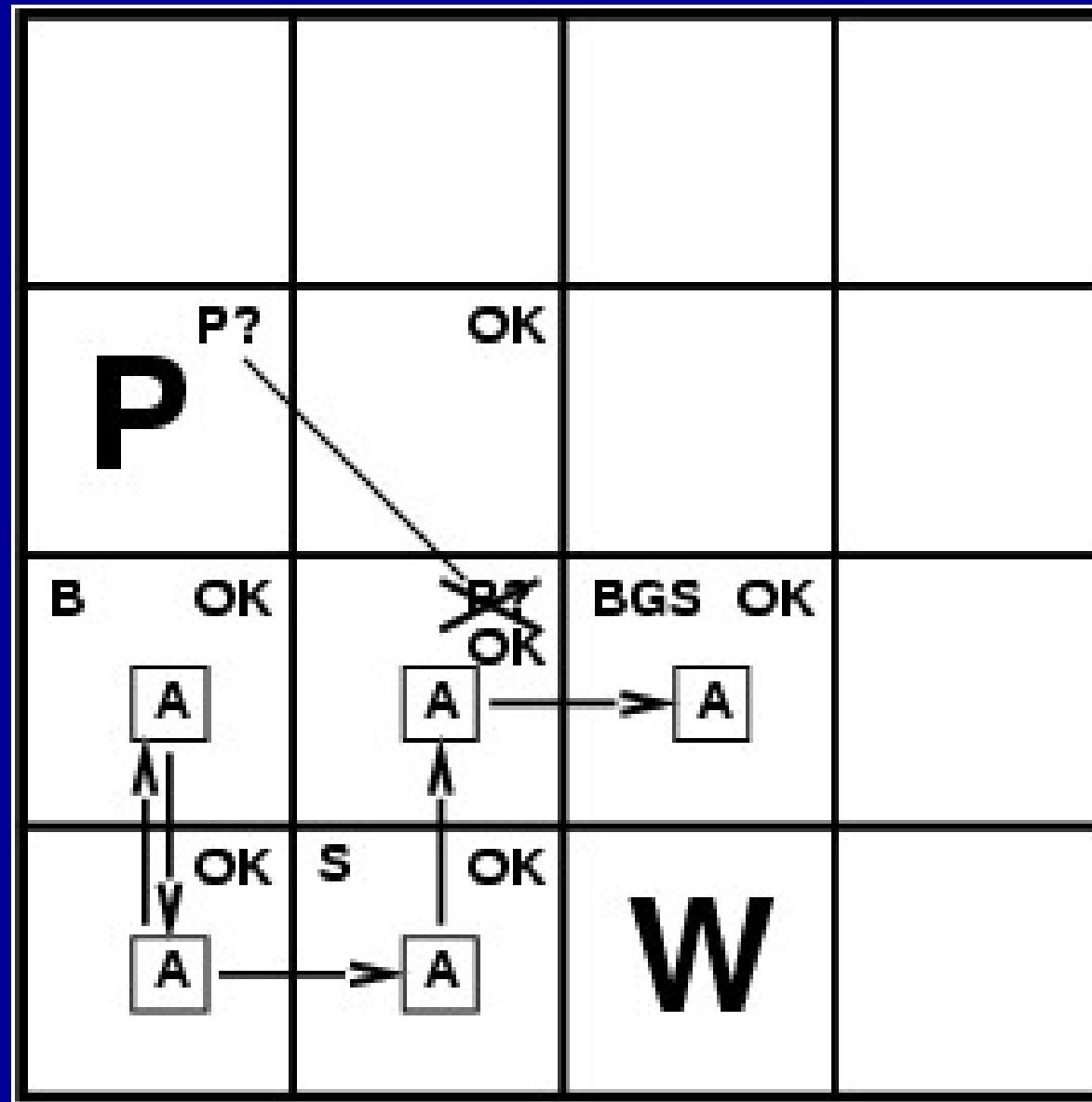
# Exploring a wumpus world



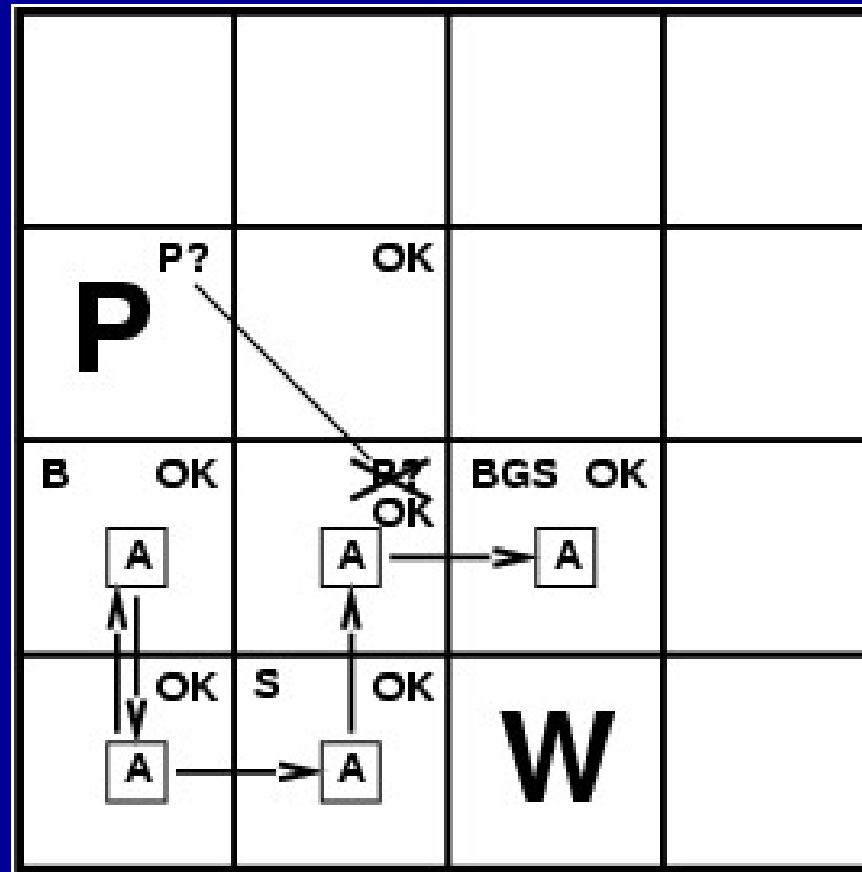
# Exploring a wumpus world



# Exploring a wumpus world



# Exploring a wumpus world



In each case where the agent draws a conclusion from the available Information, that conclusion is guaranteed to be correct if the available Information is correct...

This is a fundamental property of logical reasoning



# Logic in general

- **Logics** are formal languages for representing information such that conclusions can be drawn
- **Syntax** defines how symbols can be put together to form the sentences in the language
- **Semantics** define the "meaning" of sentences;
  - i.e., define **truth** of a sentence in a world (given an interpretation)
- 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  $x = 0, y = 6$

# Entailment

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

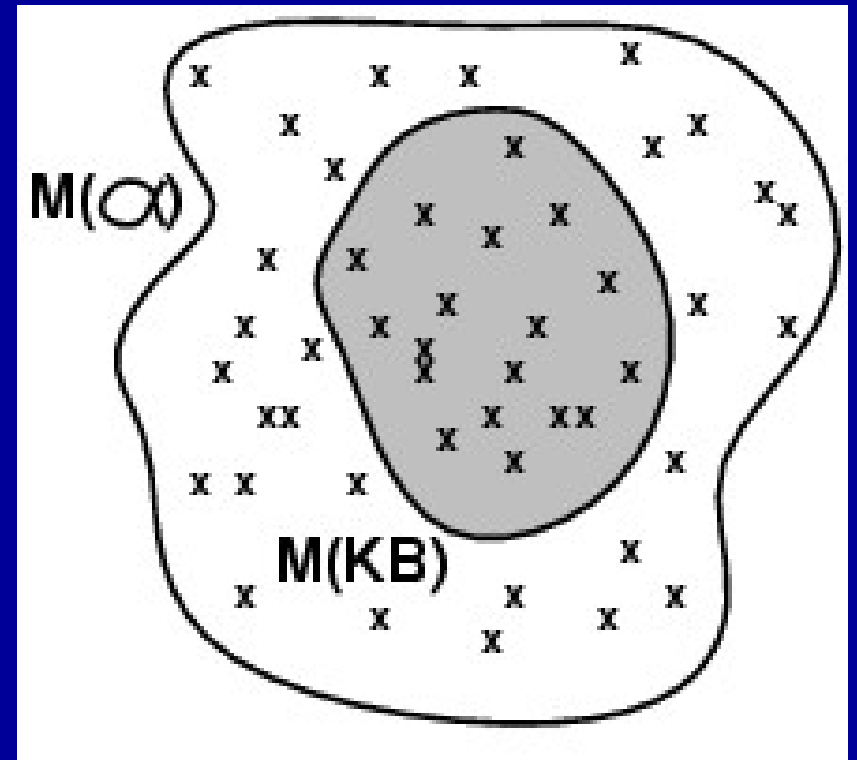
$$KB \models \alpha$$

- 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 Phillies won” and “the Reds won” entails “Either the Phillies 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.  $KB$  = Phillies won and Yankees won  
 $\alpha$  = Phillies won

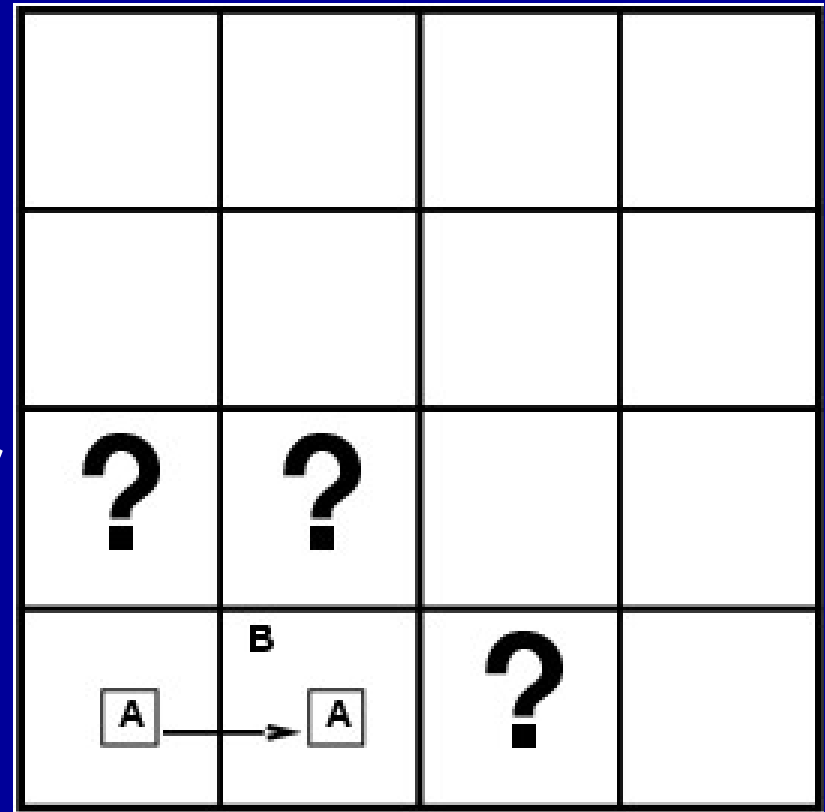


# Entailment in the wumpus world

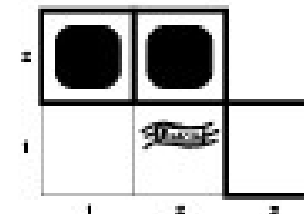
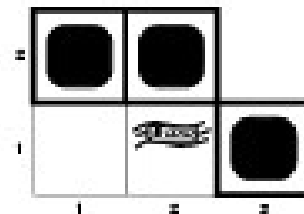
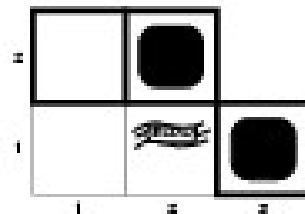
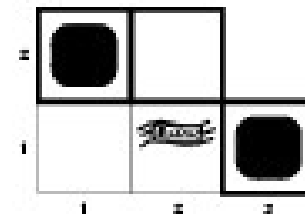
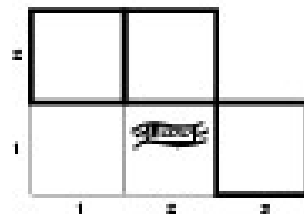
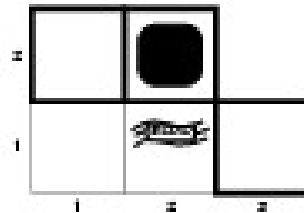
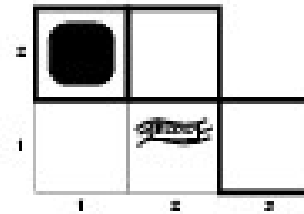
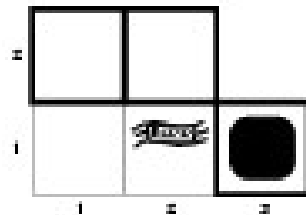
Situation after detecting  
nothing in [1,1], moving  
right, breeze in [2,1]

Consider possible models for  
*KB* assuming only pits

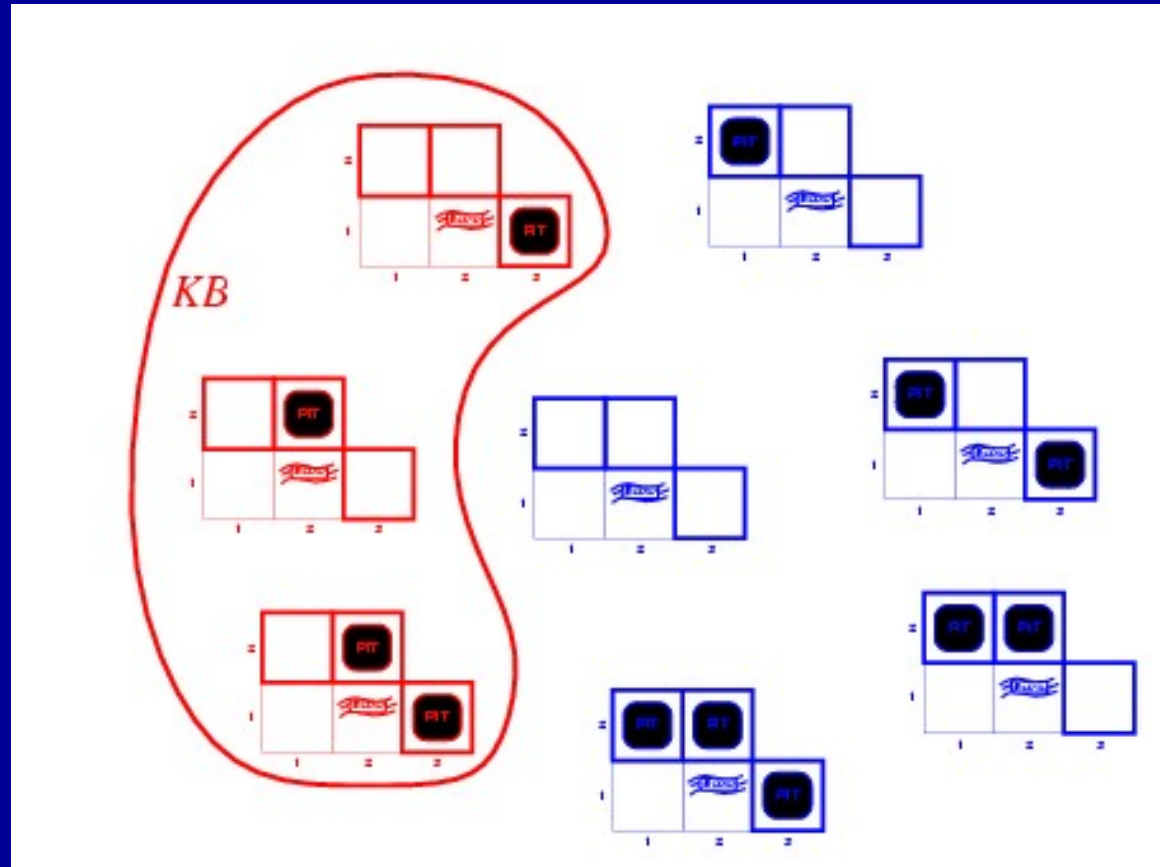
3 Boolean choices  $\Rightarrow$  8  
possible models



# Wumpus possible models

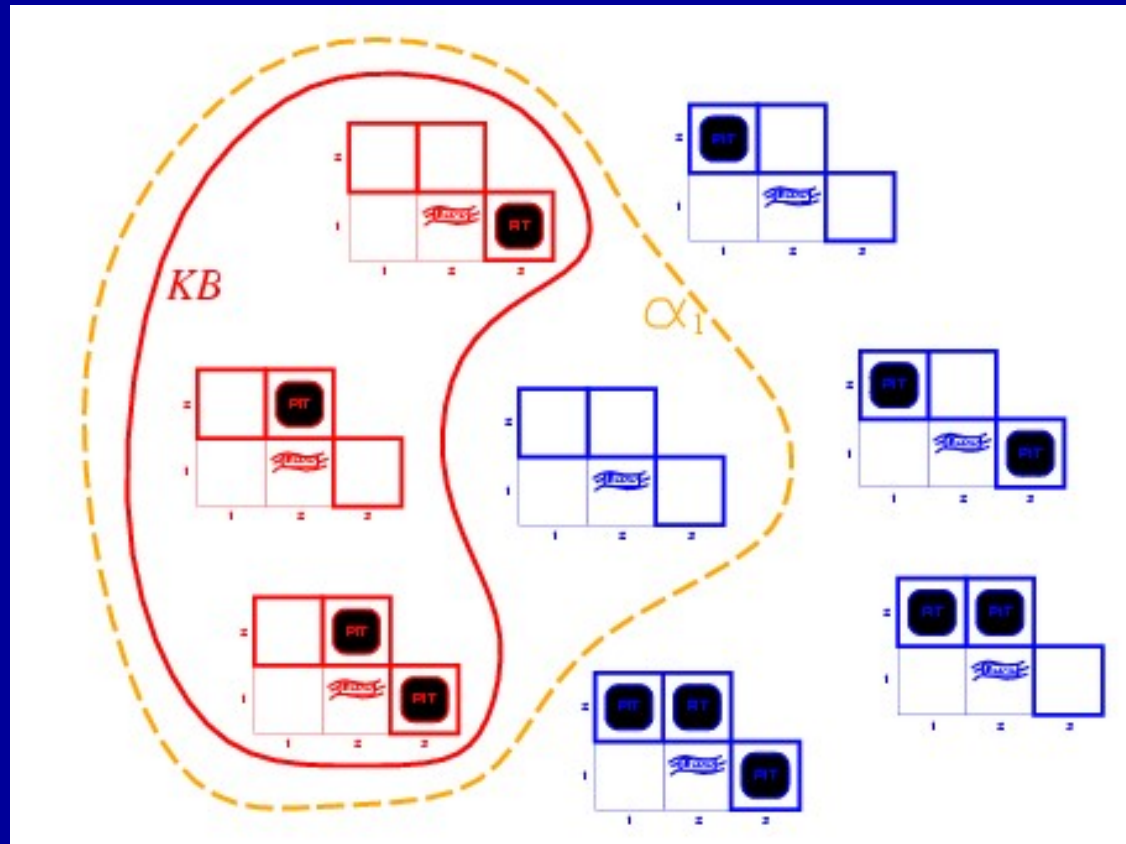


# Wumpus models



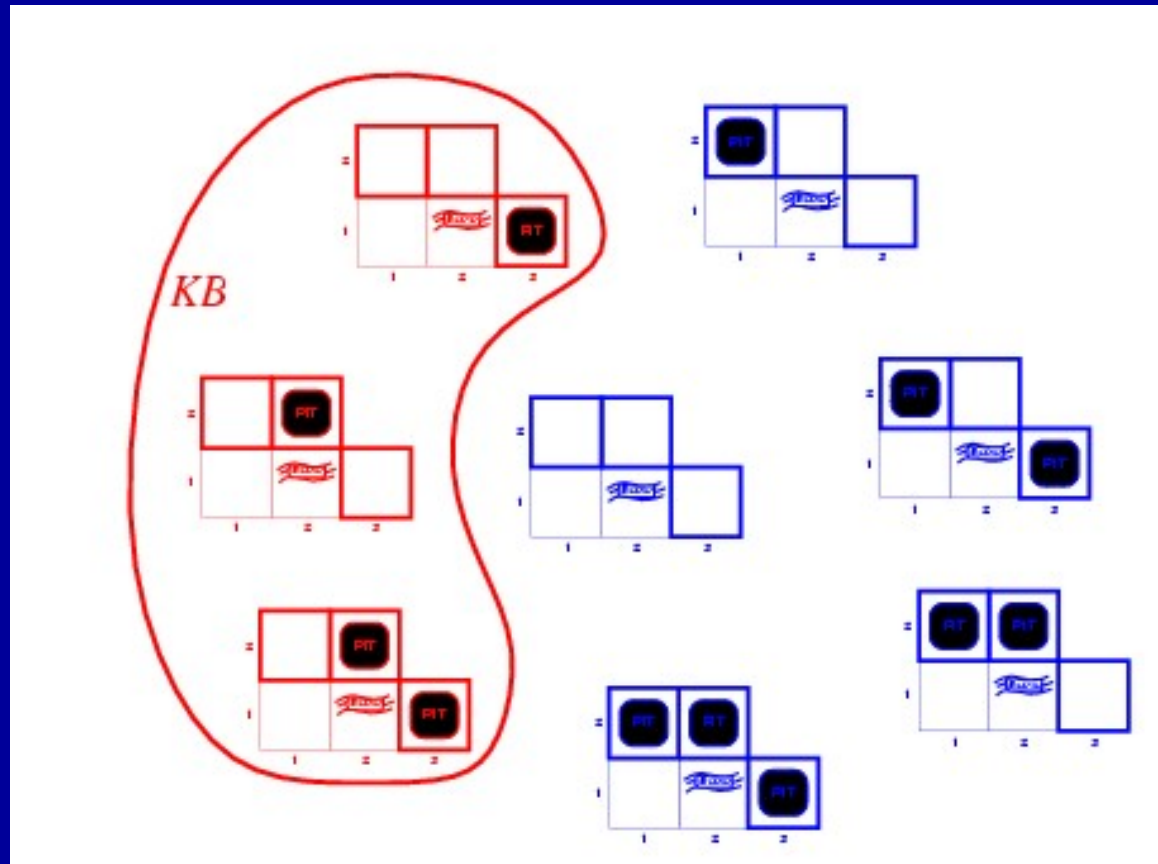
- *KB* = wumpus-world rules + observations

# Wumpus models



- $KB$  = wumpus-world rules + observations
- $\alpha_1$  = "there is no pit in [1,2]",  $KB \models \alpha_1$ , proved by **model checking**

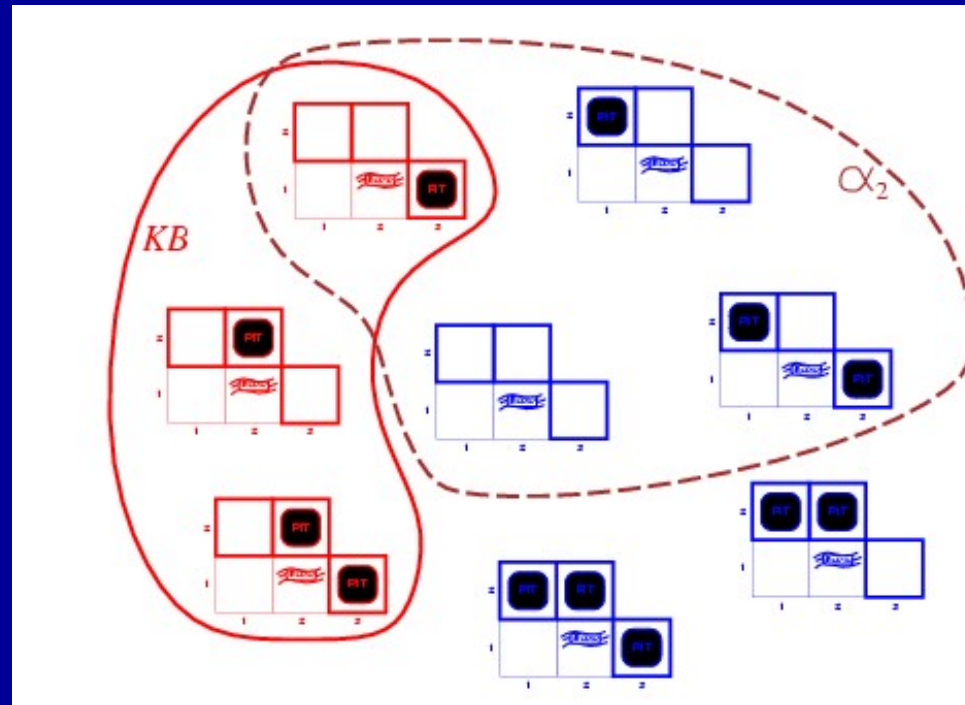
# Wumpus models



- *KB* = wumpus-world rules + observations



# Wumpus models



- $KB$  = wumpus-world rules + observations
- $\alpha_2$  = "there is no pit in  $[2,2]$ ",  $KB \models \alpha_2$

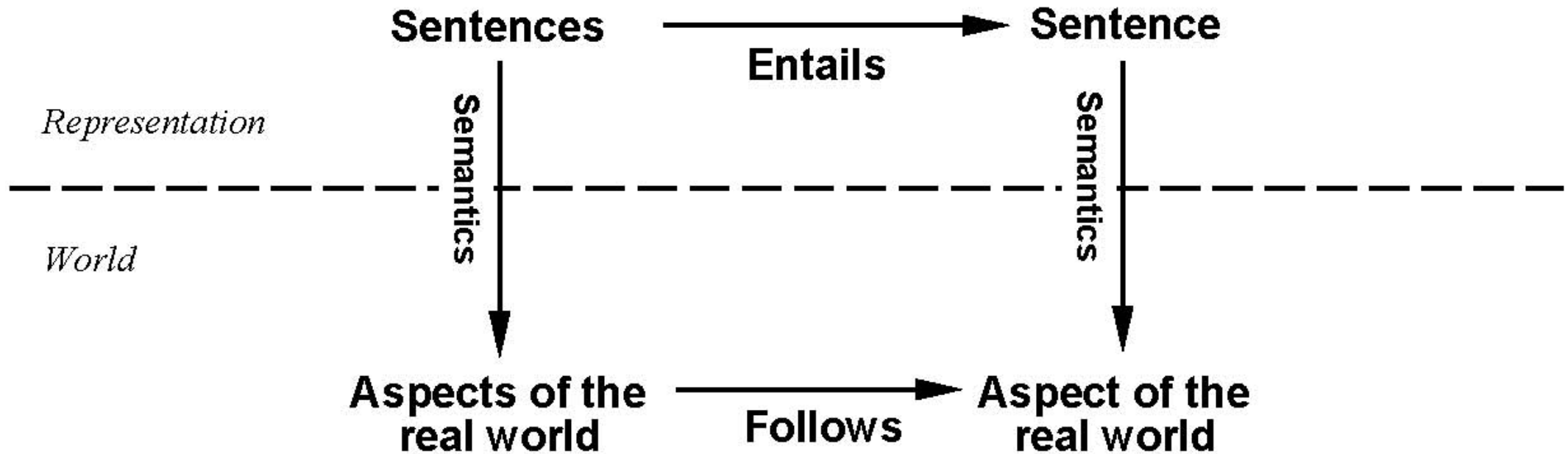
# Inference and Entailment

- Inference is a procedure that allows new sentences to be derived from a knowledge base.
- Understanding inference and entailment: think of
  - Set of all consequences of a KB as a haystack
  - $\alpha$  as the needle
- Entailment is like the needle being in the haystack
- Inference is like finding it

# Inference

- $KB \vdash_i \alpha$  = sentence  $\alpha$  can be derived from  $KB$  by inference 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  $KB \vdash_i \alpha$
- Preview: we will define a logic (first-order logic) which is expressive enough to say almost anything of interest, and for which there exists a sound and complete inference procedure.
- That is, the procedure will answer any question whose answer follows from what is known by the  $KB$ .

# Step Back...



This is an inference procedure whose conclusions are guaranteed to be true  
In any world where the premises are true.

If KB is true in the real world, then any sentence  $\alpha$  derived  
from KB by a sound inference procedure is also true in  
the real world.

# Propositional logic: Syntax

- Propositional logic is the simplest logic – illustrates basic ideas
- The proposition symbols  $P_1, P_2$  etc are (atomic) 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: 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.,

$$\neg P_{1,2} \wedge (P_{2,2} \vee P_{3,1}) = \text{true} \wedge (\text{true} \vee \text{false}) = \text{true} \wedge \text{true} = \text{true}$$

# 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>

# 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>

John likes football and John likes baseball.

John likes football or John likes baseball.

(English or is a bit different...)



# 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>

John likes football and John likes baseball.

John likes football or John likes baseball.

If John likes football then John likes baseball.

(Note different from English – if John likes football maps to false, then the sentence is true.)

(Implication seems to be if antecedent is true then I claim the consequence is, otherwise I make no claim.)

# 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"

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

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

# Simple Inference Procedure

- $KB \models \alpha$ ?
- Model checking – enumerate the models, and check if  $\alpha$  is true in every model in which  $KB$  is true. Size of truth table depends on # of atomic symbols.
- Remember – a model is a mapping of all atomic symbols to true or false – use semantics of connectives to come to an interpretation for them.

# Truth tables for inference

$B_{1,1}$	$B_{2,1}$	$P_{1,1}$	$P_{1,2}$	$P_{2,1}$	$P_{2,2}$	$P_{3,1}$	$KB$	$\alpha_1$
<i>false</i>	<i>false</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>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	<i>true</i>
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
<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>
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
<i>true</i>	<i>true</i>	<i>true</i>	<i>true</i>	<i>true</i>	<i>true</i>	<i>true</i>	<i>false</i>	<i>false</i>

# Inference by enumeration

- Depth-first enumeration of all models is sound and complete

```
function TT-ENTAILS?(KB,  $\alpha$ ) returns true or false
```

```
  symbols  $\leftarrow$  a list of the proposition symbols in KB and  $\alpha$ 
```

```
  return TT-CHECK-ALL(KB,  $\alpha$ , symbols, [])
```

---

```
function TT-CHECK-ALL(KB,  $\alpha$ , symbols, model) returns true or false
```

```
  if EMPTY?(symbols) then
```

```
    if PL-TRUE?(KB, model) then return PL-TRUE?( $\alpha$ , model)
```

```
    else return true
```

```
  else do
```

```
     $P \leftarrow$  FIRST(symbols);  $rest \leftarrow$  REST(symbols)
```

```
    return TT-CHECK-ALL(KB,  $\alpha$ , rest, EXTEND(P, true, model) and
```

```
      TT-CHECK-ALL(KB,  $\alpha$ , rest, EXTEND(P, false, model))
```

- For  $n$  symbols, time complexity is  $O(2^n)$ , space complexity is  $O(n)$

# 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

# Proof methods

- Proof methods divide into (roughly) two kinds:
  - Application of inference rules
    - Legitimate (sound) generation of new sentences from old
    - **Proof** = a sequence of inference rule applications  
Can use inference rules as operators in a standard search algorithm
    - Typically require transformation of sentences into a **normal form**
  - Model checking
    - truth table enumeration (always exponential in  $n$ )
    - improved backtracking, e.g., Davis--Putnam-Logemann-Loveland (DPLL)
    - heuristic search in model space (sound but incomplete)  
e.g., min-conflicts-like hill-climbing algorithms



# Conversion to CNF

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

1. Eliminate  $\Leftrightarrow$ , replacing  $\alpha \Leftrightarrow \beta$  with  $(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$ .

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

2. Eliminate  $\Rightarrow$ , replacing  $\alpha \Rightarrow \beta$  with  $\neg\alpha \vee \beta$ .

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

3. Move  $\neg$  inwards using de Morgan's rules and double-negation:

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

4. Apply distributivity law ( $\wedge$  over  $\vee$ ) and flatten:

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

# Resolution example

- $KB = (B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})) \wedge \neg B_{1,1} \quad \alpha = \neg P_{1,2}$

