

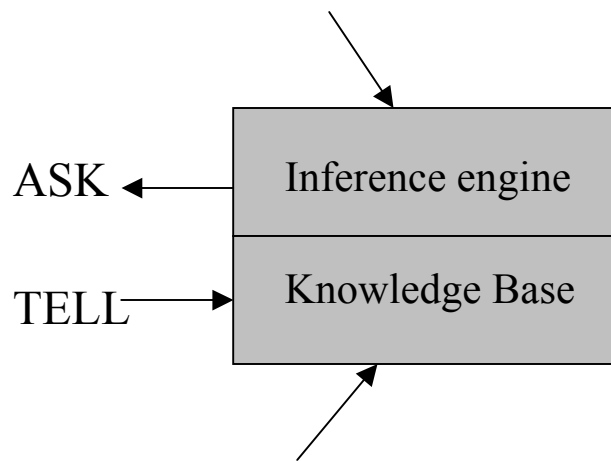
# Knowledge and reasoning – second part



- Knowledge representation
- Logic and representation
- Propositional (Boolean) logic
- Normal forms
- Inference in propositional logic
- Wumpus world example

# Knowledge-Based Agent

Domain independent algorithms



Domain specific content

- Agent that uses **prior** or **acquired** knowledge to achieve its goals
  - Can make more efficient decisions
  - Can make informed decisions
- Knowledge Base (KB): contains a set of representations of facts about the Agent's environment
- Each representation is called a **sentence**
- Use some **knowledge representation language**, to TELL it what to know e.g., (temperature 72F)
- ASK agent to query what to do
- Agent can use inference to deduce new facts from TELLED facts

# Generic 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
```

1. TELL KB what was perceived  
Uses a KRL to insert new sentences, representations of facts, into KB
2. ASK KB what to do.  
Uses logical reasoning to examine actions and select best.

# Wumpus world example

Percepts Breeze, Glitter, Smell

Actions Left turn, Right turn,  
Forward, Grab, Release, Shoot

Goals Get gold back to start  
without entering pit or wumpus square

## Environment

Squares adjacent to wumpus are smelly

Squares adjacent to pit are breezy

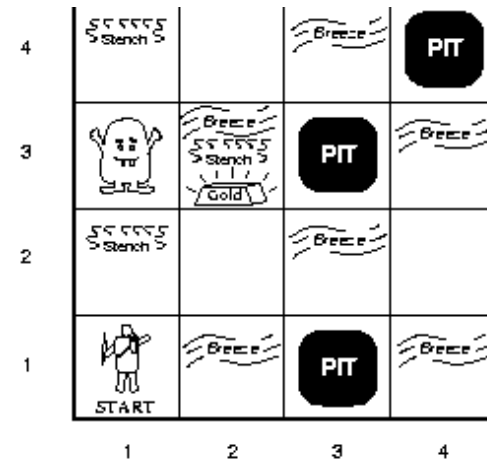
Glitter if and only if gold is in the same square

Shooting kills the wumpus if you are facing it

Shooting uses up the only arrow

Grabbing picks up the gold if in the same square

Releasing drops the gold in the same square



# Wumpus world characterization



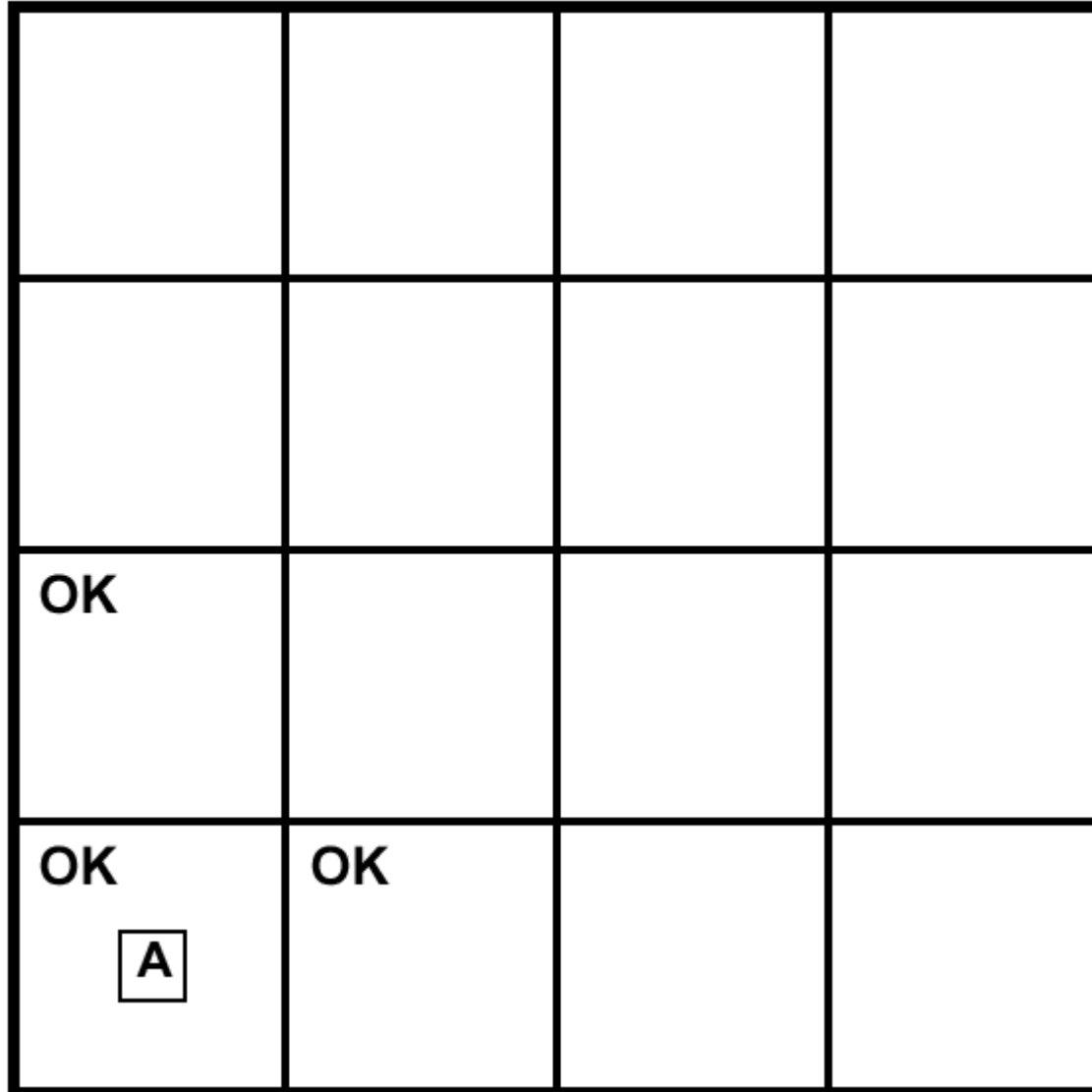
- Deterministic?
- Accessible?
- Static?
- Discrete?
- Episodic?

## Wumpus world characterization



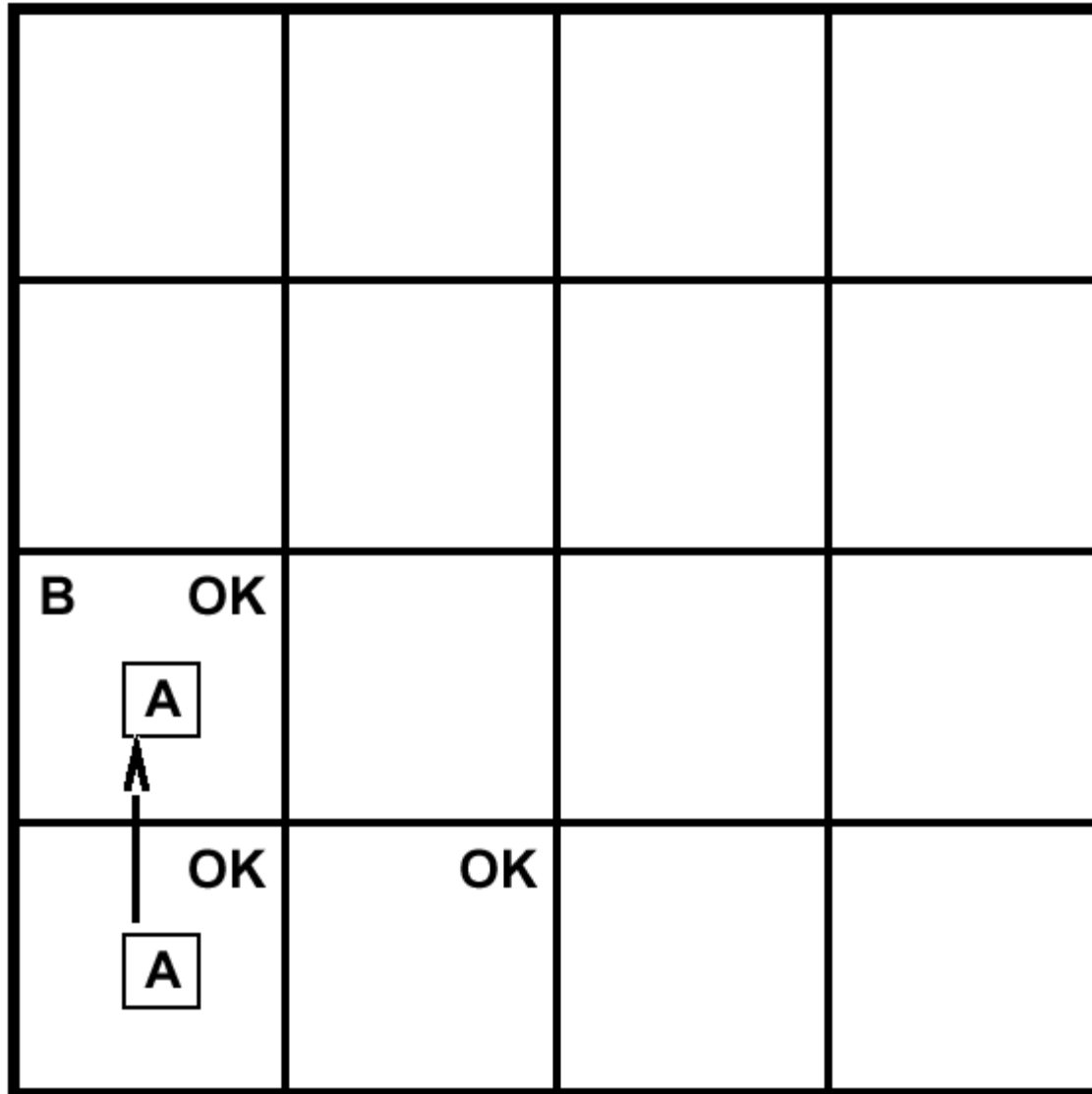
- Deterministic?            Yes – outcome exactly specified.
- Accessible?                No – only local perception.
- Static?                      Yes – Wumpus and pits do not move.
- Discrete?                  Yes
- Episodic?                  (Yes) – because static.

# Exploring a Wumpus world



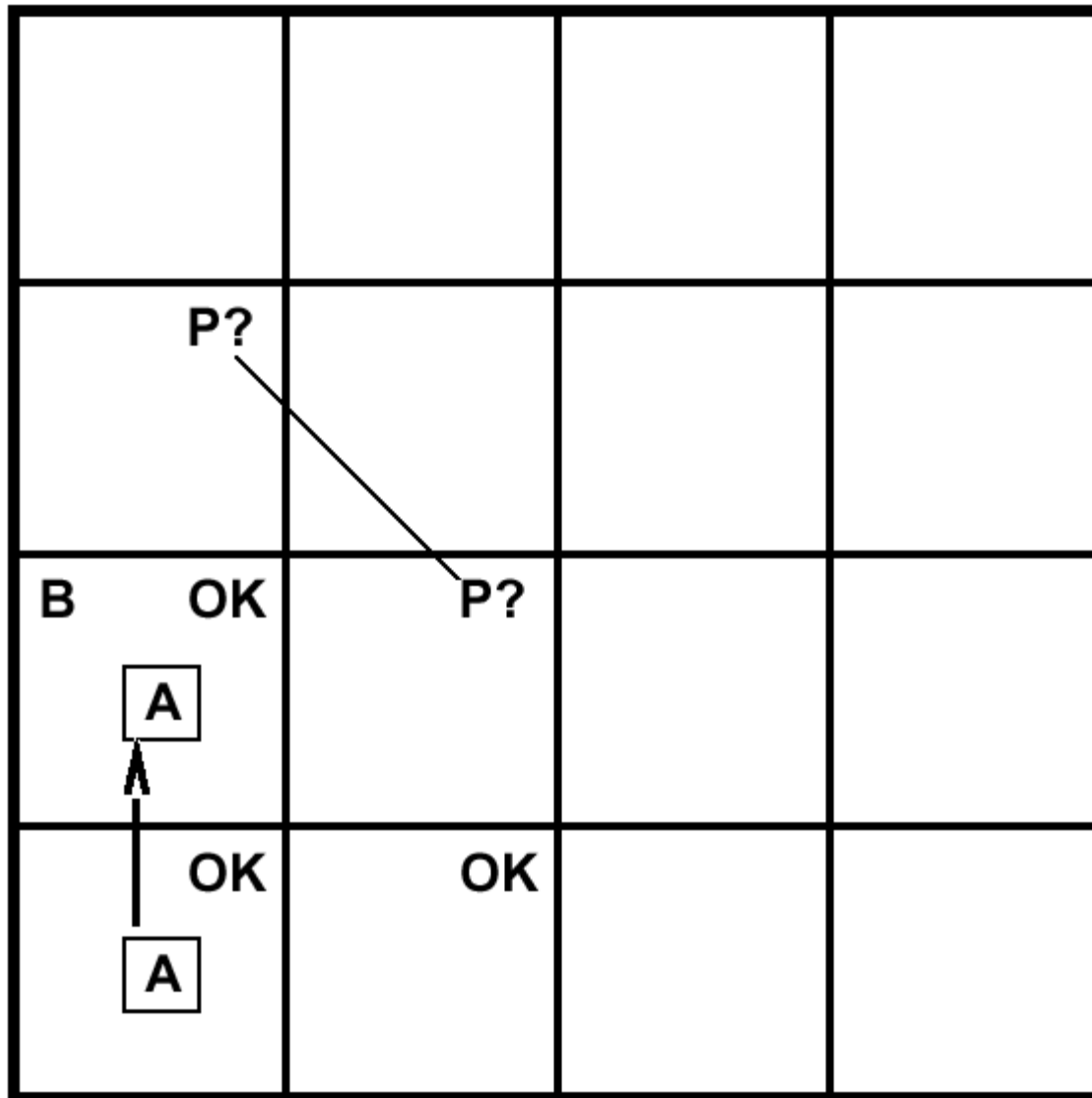
OK			
OK A	OK		

# Exploring a Wumpus world

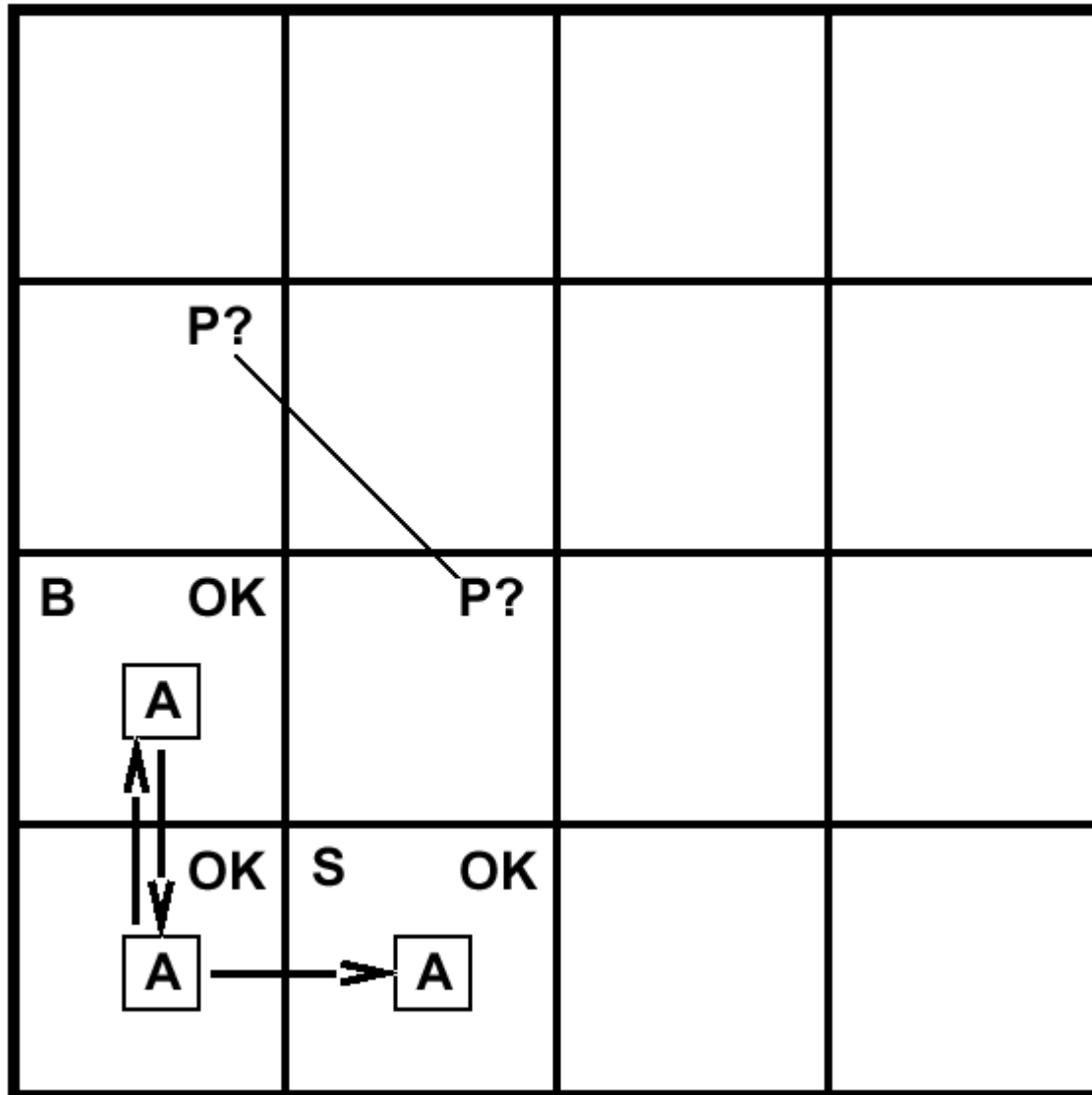




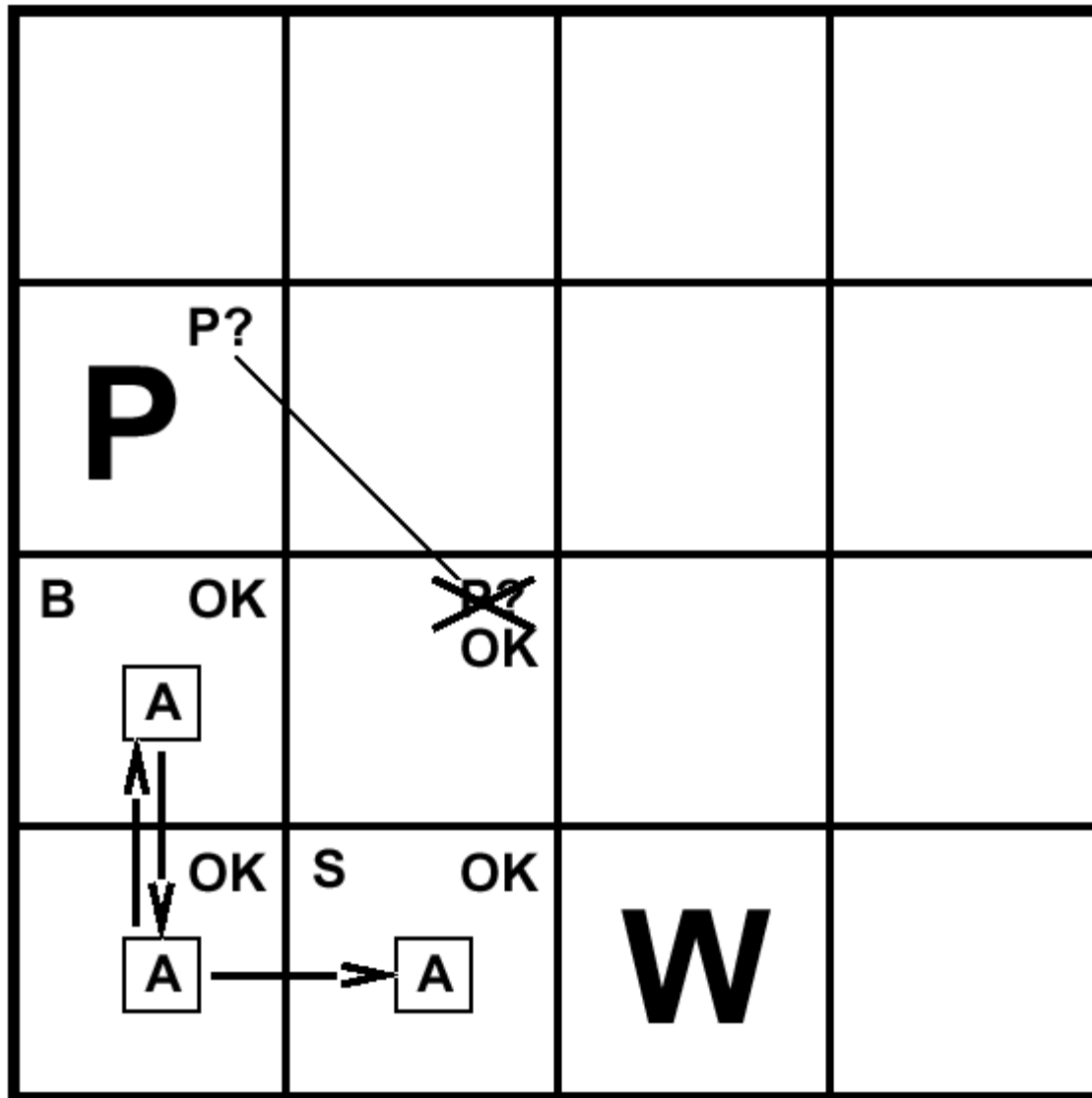
# Exploring a Wumpus world



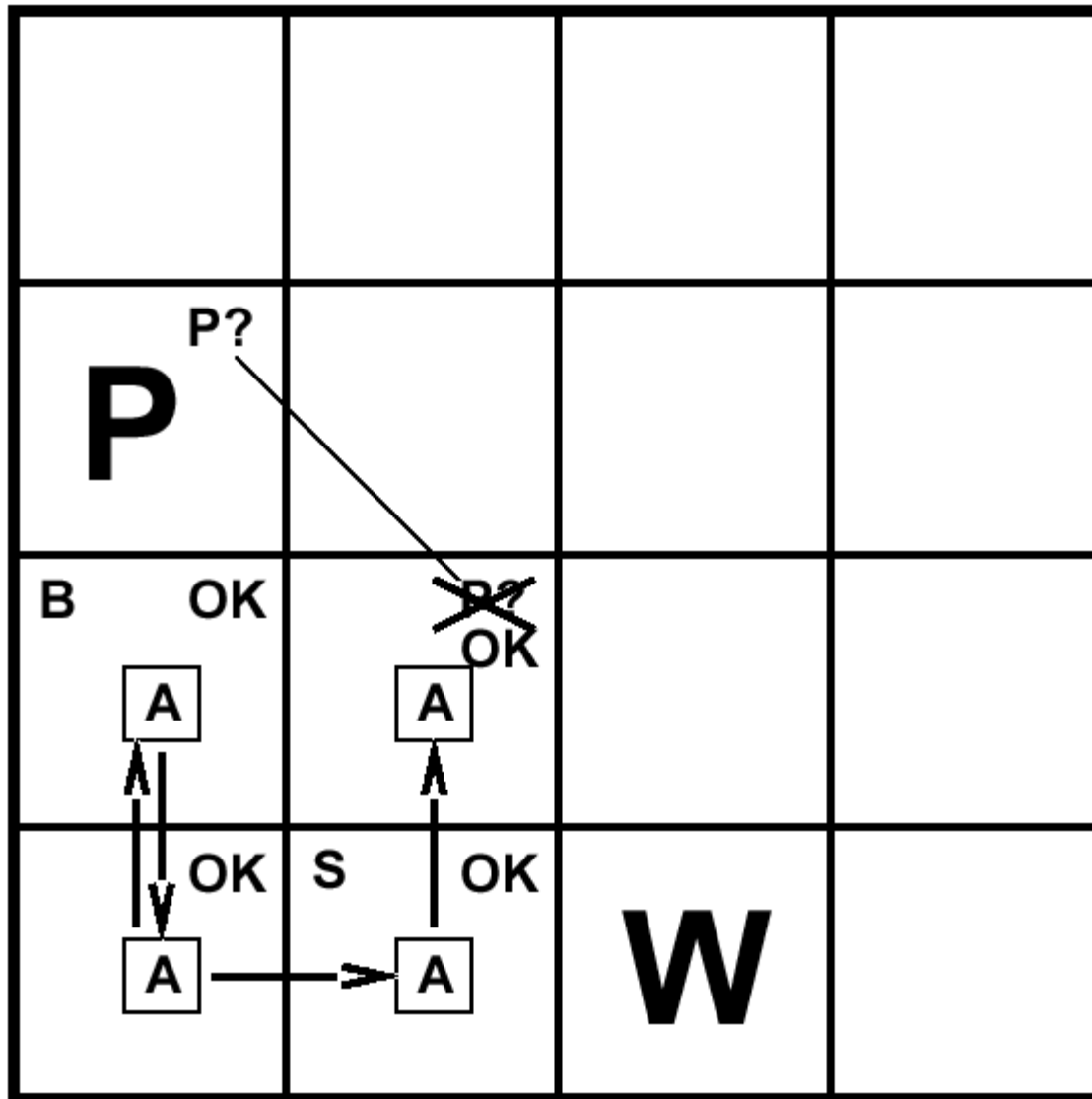
# Exploring a Wumpus world



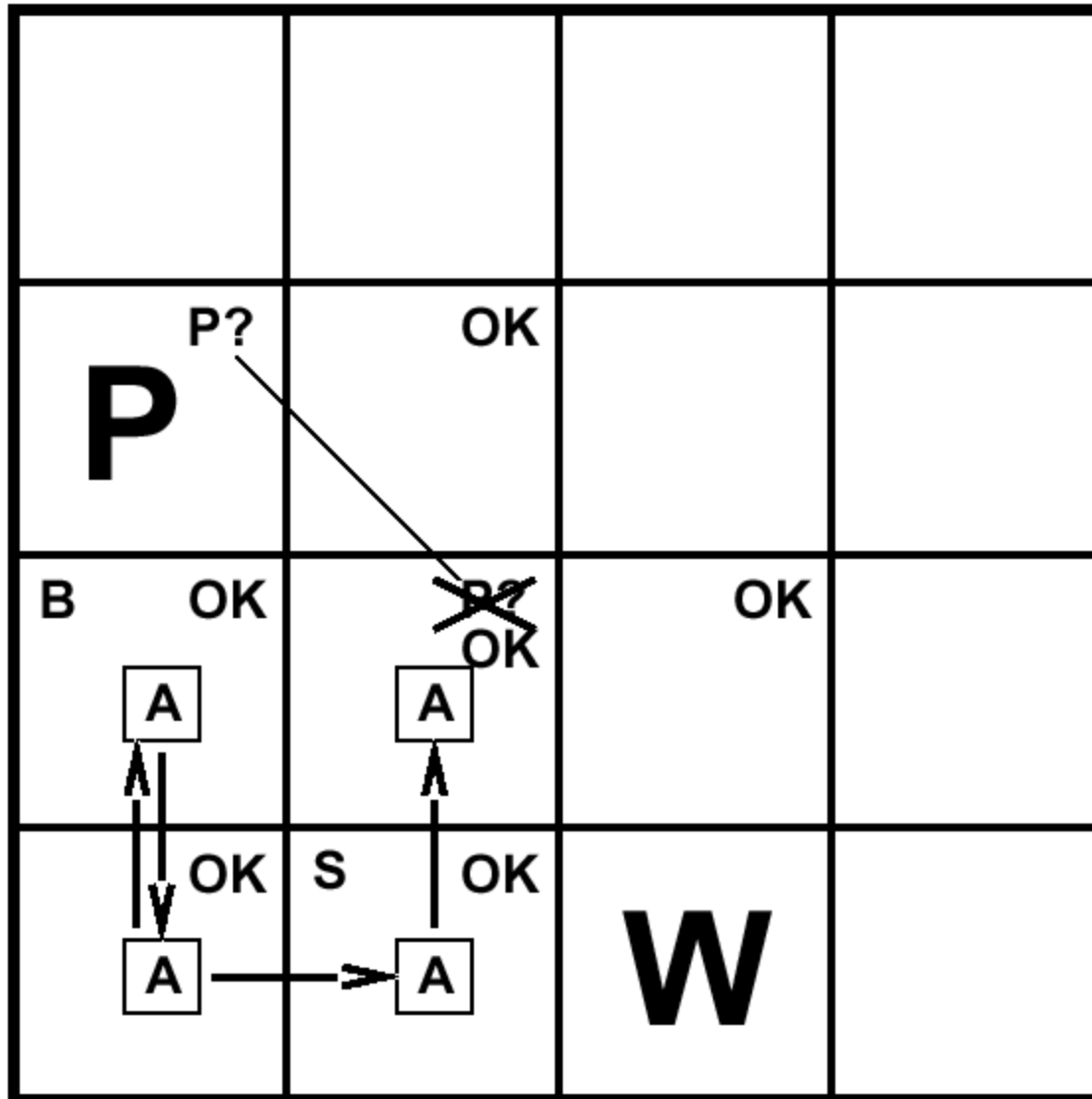
# Exploring a Wumpus world



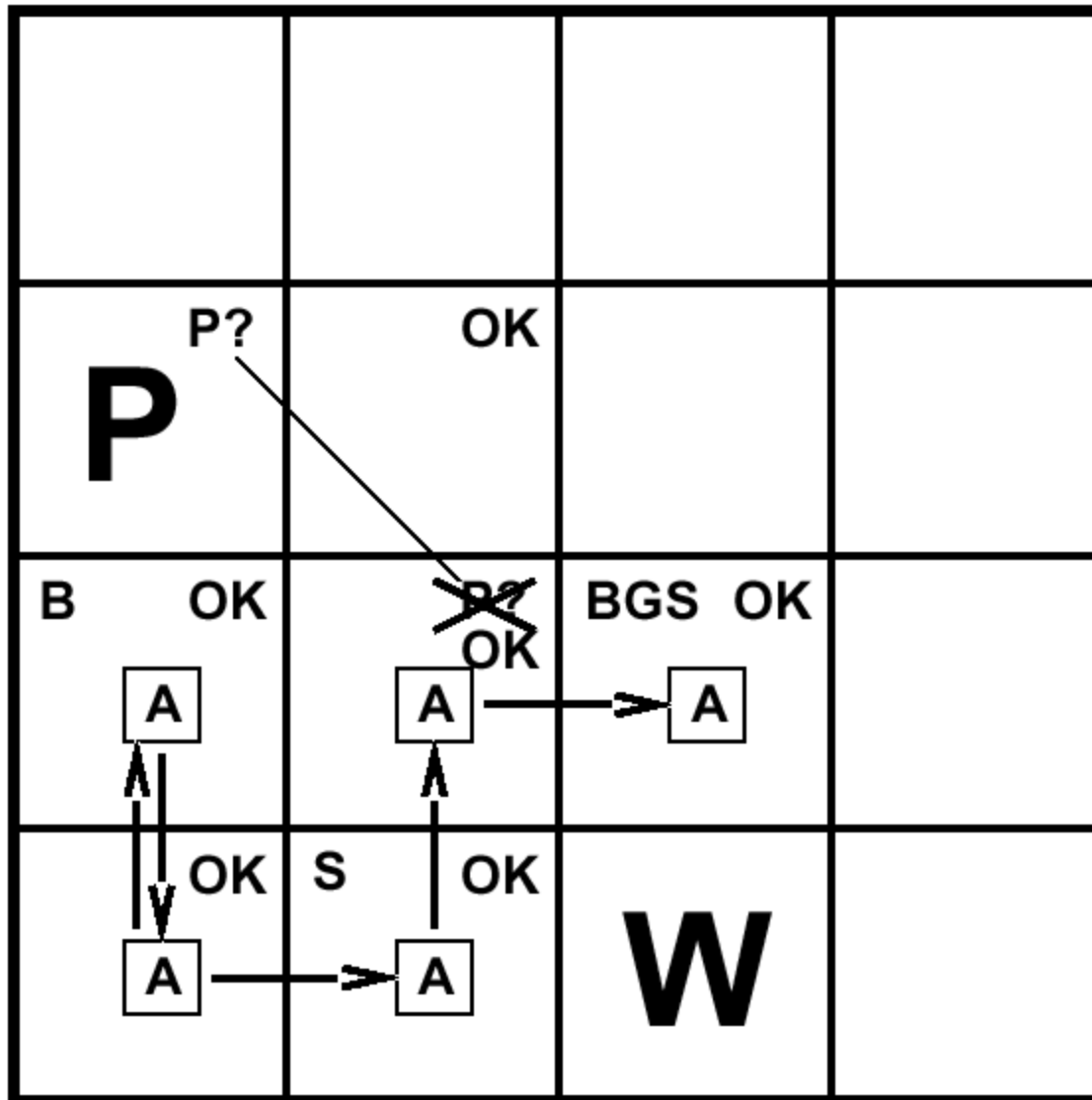
# Exploring a Wumpus world



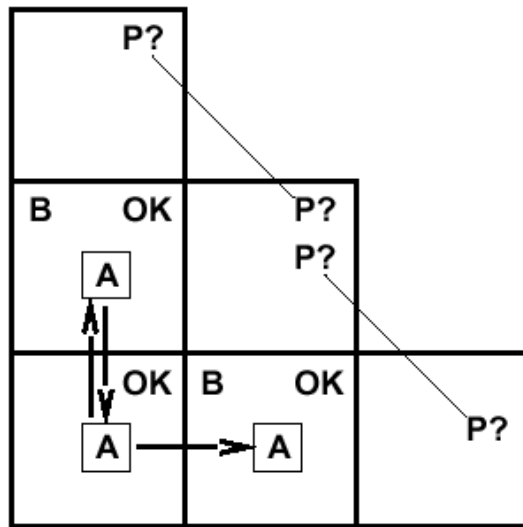
# Exploring a Wumpus world



# Exploring a Wumpus world

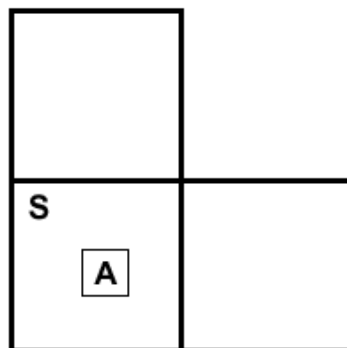


## Other tight spots



Breeze in (1,2) and (2,1)  
 $\Rightarrow$  no safe actions

Assuming pits uniformly distributed,  
 (2,2) is most likely to have a pit



Smell in (1,1)

$\Rightarrow$  cannot move

Can use a strategy of coercion:

shoot straight ahead

wumpus was there  $\Rightarrow$  dead  $\Rightarrow$  safe

wumpus wasn't there  $\Rightarrow$  safe

# Another example solution

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2	2,2	3,2	4,2
OK			
1,1	2,1	3,1	4,1
<b>A</b>			
OK	OK		

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

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2	2,2	3,2	4,2
OK	P?		
1,1	2,1	3,1	4,1
V	<b>A</b>	P?	
OK	B		
	OK		

No perception → 1,2 and 2,1 OK

Move to 2,1

B in 2,1 → 2,2 or 3,1 P?

1,1 V → no P in 1,1

Move to 1,2 (only option)



# Example solution

1,4	2,4	3,4	4,4
1,3 W!	2,3	3,3	4,3
1,2 <b>A</b> S OK	2,2  OK	3,2	4,2
1,1  V OK	2,1 B V OK	3,1 P!	4,1

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

1,4	2,4 P?	3,4	4,4
1,3 W!	2,3 <b>A</b> S G B	3,3 P?	4,3
1,2 S V OK	2,2  V OK	3,2	4,2
1,1  V OK	2,1 B V OK	3,1 P!	4,1

S and No S when in 2,1 → 1,3 or 1,2 has W

1,2 OK → 1,3 W

No B in 1,2 → 2,2 OK & 3,1 P

## 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  $x = 0$ ,  $y = 6$

# Types of logic

Logics are characterized by what they commit to as “primitives”

Ontological commitment: what exists—facts? objects? time? beliefs?

Epistemological commitment: what states of knowledge?

Language	Ontological Commitment	Epistemological Commitment
Propositional logic	facts	true/false/unknown
First-order logic	facts, objects, relations	true/false/unknown
Temporal logic	facts, objects, relations, times	true/false/unknown
Probability theory	facts	degree of belief 0...1
Fuzzy logic	degree of truth	degree of belief 0...1

# Entailment



$$KB \models \alpha$$

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”

# Models

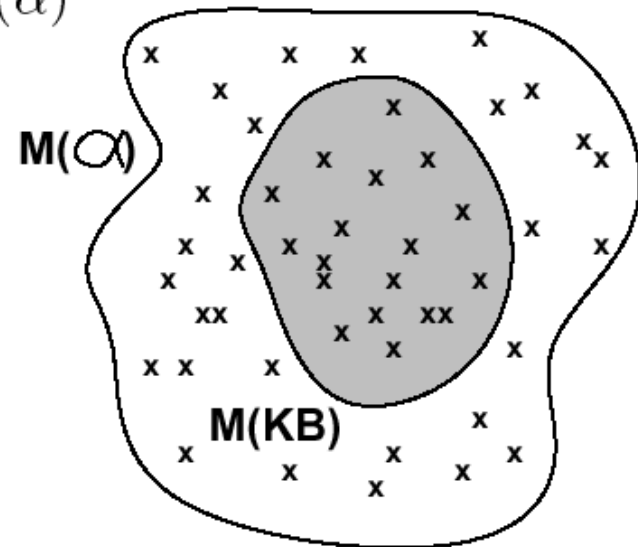
Logicians typically think in terms of models, which are formally structured worlds with respect to which truth can be evaluated

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  $KB \models \alpha$  if and only if  $M(KB) \subseteq M(\alpha)$

E.g.  $KB =$  Giants won and Reds won  
 $\alpha =$  Giants won



# 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  $KB \vdash_i \alpha$

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  $KB$ .

## Basic symbols

- Expressions only evaluate to either “true” or “false.”
  
- $P$  “P is true”
- $\neg P$  “P is false” negation
- $P \vee Q$  “either P is true or Q is true or both” disjunction
- $P \wedge Q$  “both P and Q are true” conjunction
- $P \Rightarrow Q$  “if P is true, the Q is true” implication
- $P \Leftrightarrow Q$  “P and Q are either both true or both false” equivalence

## Propositional logic: syntax

Propositional logic is the simplest logic

The proposition symbols  $P_1, P_2$  etc are sentences

If  $S$  is a sentence,  $\neg S$  is a sentence

If  $S_1$  and  $S_2$  is a sentence,  $S_1 \wedge S_2$  is a sentence

If  $S_1$  and  $S_2$  is a sentence,  $S_1 \vee S_2$  is a sentence

If  $S_1$  and  $S_2$  is a sentence,  $S_1 \Rightarrow S_2$  is a sentence

If  $S_1$  and  $S_2$  is a sentence,  $S_1 \Leftrightarrow S_2$  is a sentence



# Propositional logic: semantics

Each model specifies true/false for each proposition symbol

E.g.  $A$     $B$     $C$   
*True True False*

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

# Truth tables

- Truth value: whether a statement is true or false.
- Truth table: complete list of truth values for a statement given all possible values of the individual atomic expressions.

Example:

<u>P</u>	<u>Q</u>	<u>P V Q</u>
T	T	T
T	F	T
F	T	T
F	F	F

## Truth tables for basic connectives

P	Q	$\neg P$	$\neg Q$	$P \vee Q$	$P \wedge Q$	$P \Rightarrow Q$	$P \Leftrightarrow Q$
T	T	F	F	T	T	T	T
T	F	F	T	T	F	F	F
F	T	T	F	T	F	T	F
F	F	T	T	F	F	T	T

# Propositional logic: basic manipulation rules

- $\neg(\neg A) = A$  Double negation
- $\neg(A \wedge B) = (\neg A) \vee (\neg B)$  Negated "and"
- $\neg(A \vee B) = (\neg A) \wedge (\neg B)$  Negated "or"
- $A \wedge (B \vee C) = (A \wedge B) \vee (A \wedge C)$  Distributivity of  $\wedge$  on  $\vee$
- $A \Rightarrow B = (\neg A) \vee B$  by definition
- $\neg(A \Rightarrow B) = A \wedge (\neg B)$  using negated or
- $A \Leftrightarrow B = (A \Rightarrow B) \wedge (B \Rightarrow A)$  by definition
- $\neg(A \Leftrightarrow B) = (A \wedge (\neg B)) \vee (B \wedge (\neg A))$  using negated and & or
- ...

# Propositional inference: enumeration method

Let  $\alpha = A \vee B$  and  $KB = (A \vee C) \wedge (B \vee \neg C)$

Is it the case that  $KB \models \alpha$ ?

Check all possible models— $\alpha$  must be true wherever  $KB$  is true

<i>A</i>	<i>B</i>	<i>C</i>	$A \vee C$	$B \vee \neg C$	$KB$	$\alpha$
<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>True</i>	<i>True</i>				
<i>True</i>	<i>False</i>	<i>False</i>				
<i>True</i>	<i>False</i>	<i>True</i>				
<i>True</i>	<i>True</i>	<i>False</i>				
<i>True</i>	<i>True</i>	<i>True</i>				

# Enumeration: Solution

<i>A</i>	<i>B</i>	<i>C</i>	$A \vee C$	$B \vee \neg C$	<i>KB</i>	$\alpha$
<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>True</i>	<i>True</i>	<i>False</i>	<i>False</i>	<i>False</i>
<i>False</i>	<i>True</i>	<i>False</i>	<i>False</i>	<i>True</i>	<i>False</i>	<i>True</i>
<i>False</i>	<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>	<i>True</i>	<i>True</i>	<i>True</i>	<i>True</i>
<i>True</i>	<i>False</i>	<i>True</i>	<i>True</i>	<i>False</i>	<i>False</i>	<i>True</i>
<i>True</i>	<i>True</i>	<i>False</i>	<i>True</i>	<i>True</i>	<i>True</i>	<i>True</i>
<i>True</i>	<i>True</i>	<i>True</i>	<i>True</i>	<i>True</i>	<i>True</i>	<i>True</i>

# Propositional inference: normal forms

Other approaches to inference use syntactic operations on sentences, often expressed in standardized forms

Conjunctive Normal Form (CNF—universal)

*conjunction of disjunctions of literals*  
*clauses*

E.g.,  $(A \vee \neg B) \wedge (B \vee \neg C \vee \neg D)$

“product of sums of simple variables or negated simple variables”

Disjunctive Normal Form (DNF—universal)

*disjunction of conjunctions of literals*  
*terms*

E.g.,  $(A \wedge B) \vee (A \wedge \neg C) \vee (A \wedge \neg D) \vee (\neg B \wedge \neg C) \vee (\neg B \wedge \neg D)$

“sum of products of simple variables or negated simple variables”

Horn Form (restricted)

*conjunction of Horn clauses (clauses with  $\leq 1$  positive literal)*

E.g.,  $(A \vee \neg B) \wedge (B \vee \neg C \vee \neg D)$

Often written as set of implications:

$B \Rightarrow A$  and  $(C \wedge D) \Rightarrow B$

## Deriving expressions from functions

- Given a boolean function in truth table form, find a propositional logic expression for it that uses only  $\vee$ ,  $\wedge$  and  $\neg$ .
- **Idea:** We can easily do it by disjoining the "T" rows of the truth table.

Example: XOR function

P	Q	RESULT	
T	T	F	
T	F	T	$P \wedge (\neg Q)$
F	T	T	$(\neg P) \wedge Q$
F	F	F	

$$\text{RESULT} = (P \wedge (\neg Q)) \vee ((\neg P) \wedge Q)$$



## A more formal approach



- To construct a logical expression in disjunctive normal form from a truth table:
  - Build a **“minterm”** for each row of the table, where:
    - For each variable whose value is T in that row, include the variable in the minterm
    - For each variable whose value is F in that row, include the negation of the variable in the minterm
    - Link variables in minterm by conjunctions
  - The expression consists of the **disjunction of all minterms**.

## Example: adder with carry

Takes 3 variables in:  $x$ ,  $y$  and  $ci$  (carry-in); yields 2 results: sum ( $s$ ) and carry-out ( $co$ ). To get you used to other notations, here we assume  $T = 1$ ,  $F = 0$ ,  $V = \text{OR}$ ,  $\wedge = \text{AND}$ ,  $\neg = \text{NOT}$ .

$x$	$y$	$ci$	$co$	$s$	
0	0	0	0	0	
0	0	1	0	1	$s : \text{NOT } x \text{ AND NOT } y \text{ AND } ci$
0	1	0	0	1	$s : \text{NOT } x \text{ AND } y \text{ AND NOT } ci$
0	1	1	1	0	$co: \text{NOT } x \text{ AND } y \text{ AND } ci$
1	0	0	0	1	$s : x \text{ AND NOT } y \text{ AND NOT } ci$
1	0	1	1	0	$co: x \text{ AND NOT } y \text{ AND } ci$
1	1	0	1	0	$co: x \text{ AND } y \text{ AND NOT } ci$
1	1	1	1	1	$co,s: x \text{ AND } y \text{ AND } ci$

The logical expression for  $co$  is:

$(\text{NOT } x \text{ AND } y \text{ AND } ci) \text{ OR } (x \text{ AND NOT } y \text{ AND } ci) \text{ OR}$   
 $(x \text{ AND } y \text{ AND NOT } ci) \text{ OR } (x \text{ AND } y \text{ AND } ci)$

The logical expression for  $s$  is:

$(\text{NOT } x \text{ AND NOT } y \text{ AND } ci) \text{ OR } (\text{NOT } x \text{ AND } y \text{ AND NOT } ci)$   
 $\text{OR } (x \text{ AND NOT } y \text{ AND NOT } ci) \text{ OR } (x \text{ AND } y \text{ AND } ci)$

# Tautologies

- Logical expressions that are always true. Can be simplified out.

Examples:

$T$

$T \vee A$

$A \vee (\neg A)$

$\neg(A \wedge (\neg A))$

$A \Leftrightarrow A$

$((P \vee Q) \Leftrightarrow P) \vee (\neg P \wedge Q)$

$(P \Leftrightarrow Q) \Rightarrow (P \Rightarrow Q)$

## Validity and satisfiability

A sentence is valid if it is true in all models

e.g.,  $A \vee \neg A$ ,  $A \Rightarrow A$ ,  $(A \wedge (A \Rightarrow B)) \Rightarrow$

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

i.e., prove  $\alpha$  by *reductio ad absurdum*

# Proof methods



Proof methods divide into (roughly) two kinds:

## Model checking

truth table enumeration (sound and complete for propositional)  
heuristic search in model space (sound but incomplete)  
e.g., the GSAT algorithm (Ex. 6.15)

## 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 alg.

# Inference rules

- ◇ **Modus Ponens** or **Implication-Elimination**: (From an implication and the premise of the implication, you can infer the conclusion.)

$$\frac{\alpha \Rightarrow \beta, \quad \alpha}{\beta}$$

- ◇ **And-Elimination**: (From a conjunction, you can infer any of the conjuncts.)

$$\frac{\alpha_1 \wedge \alpha_2 \wedge \dots \wedge \alpha_n}{\alpha_i}$$

- ◇ **And-Introduction**: (From a list of sentences, you can infer their conjunction.)

$$\frac{\alpha_1, \alpha_2, \dots, \alpha_n}{\alpha_1 \wedge \alpha_2 \wedge \dots \wedge \alpha_n}$$

- ◇ **Or-Introduction**: (From a sentence, you can infer its disjunction with anything else at all.)

$$\frac{\alpha_i}{\alpha_1 \vee \alpha_2 \vee \dots \vee \alpha_n}$$

- ◇ **Double-Negation Elimination**: (From a doubly negated sentence, you can infer a positive sentence.)

$$\frac{\neg\neg\alpha}{\alpha}$$

- ◇ **Unit Resolution**: (From a disjunction, if one of the disjuncts is false, then you can infer the other one is true.)

$$\frac{\alpha \vee \beta, \quad \neg\beta}{\alpha}$$

- ◇ **Resolution**: (This is the most difficult. Because  $\beta$  cannot be both true and false, one of the other disjuncts must be true in one of the premises. Or equivalently, implication is transitive.)

$$\frac{\alpha \vee \beta, \quad \neg\beta \vee \gamma}{\alpha \vee \gamma}$$

or equivalently

$$\frac{\neg\alpha \Rightarrow \beta, \quad \beta \Rightarrow \gamma}{\neg\alpha \Rightarrow \gamma}$$

## Wumpus world: example

- **Facts:** Percepts inject (TELL) facts into the KB
  - [stench at 1,1 and 2,1]  $\rightarrow S_{1,1} ; S_{2,1}$
- **Rules:** if square has no stench then neither the square or adjacent square contain the wumpus
  - R1:  $\neg S_{1,1} \Rightarrow \neg W_{1,1} \wedge \neg W_{1,2} \wedge \neg W_{2,1}$
  - R2:  $\neg S_{2,1} \Rightarrow \neg W_{1,1} \wedge \neg W_{2,1} \wedge \neg W_{2,2} \wedge \neg W_{3,1}$
  - ...
- **Inference:**
  - KB contains  $\neg S_{1,1}$  then using Modus Ponens we infer  $\neg W_{1,1} \wedge \neg W_{1,2} \wedge \neg W_{2,1}$
  - Using And-Elimination we get:  $\neg W_{1,1} \quad \neg W_{1,2} \quad \neg W_{2,1}$
  - ...

# Limitations of Propositional Logic



1. It is too weak, i.e., has very limited expressiveness:
  - Each rule has to be represented for each situation:  
e.g., “don’t go forward if the wumpus is in front of you” takes 64 rules
2. It cannot keep track of changes:
  - If one needs to track changes, e.g., where the agent has been before then we need a timed-version of each rule. To track 100 steps we’ll then need 6400 rules for the previous example.

Its **hard to write and maintain** such a huge rule-base  
**Inference becomes intractable**



# Summary



Logical agents apply inference to a knowledge base to derive new information and make decisions

Basic concepts of logic:

- syntax: formal structure of sentences
- semantics: truth of sentences wrt models
- entailment: necessary truth of one sentence given another
- inference: deriving sentences from other sentences
- soundness: derivations produce only entailed sentences
- completeness: derivations can produce all entailed sentences

Wumpus world requires the ability to represent partial and negated information, reason by cases, etc.

Propositional logic suffices for some of these tasks

Truth table method is sound and complete for propositional logic

## Next time



- First-order logic: [AIMA] Chapter 7