# Artificial Intelligence
## Logical AI

Christopher Simpkins

# Logic and AI

In AI, **knowledge-based** agents use a process of **reasoning** over an internal **representation** of knowledge to decide what actions to take.

- ▶ **Knowledge base**: a set of sentences.
- ▶ **Sentence**: an assertion about the world expressed in a **knowledge represeantation language**, like propositional logic.
- ▶ **Axioms**: sentences taken as given – not derived from other sentences, assumptions.
- ▶ **Inference**: deriving new sentences from old sentences.
- ▶ **Background knowledge**: sentences present in an agent's knowledge base before it starts perceiving and acting.

# Knowledge-Based Agents

- MAKE-PERCEPT-SENTENCE constructs sentence asserting that agent perceived the given percept at the given time.
- MAKE-ACTION-QUERY constructs sentence that asks what action to take at current time.
- MAKE-ACTION-SENTENCE constructs sentence asserting chosen action was executed.

> **function** KB-AGENT(*percept*) **returns** an *action*
> **persistent**: *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*))
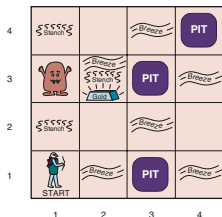> TELL(*KB*, MAKE-ACTION-SENTENCE(*action*, *t*))
> *t* ← *t* + 1
> **return** *action*

- Logical agents are described at the **knowledge level** using **declarative** statements of knowledge and goals.
- At the **implemetation level** we use a **procedureal** approach, encoding behaviors directly in program code.

# The Wumpus World



A cave that you can drop into or climb out of at square [1, 1].

▶ **Performance measure**: +1000 for climbing out of the cave with the gold, −1000 for falling into a pit or being eaten by the wumpus, −1 for each action taken, −10 for using the arrow. The game ends either when the agent dies or when the agent climbs out of the cave.

▶ **Environment**: A 4×4 grid of rooms, agent always starts at [1,1], facing east. Gold and the wumpus placed uniformly randomly from the squares other than the start square. Each non-start square can be a pit with probability 0.2.

▶ **Actuators**: Forward, TurnLeft by $90°$, or TurnRight by $90°$. The agent dies if it enters pit or a live wumpus square. Moves into walls have no effect. Grab picks up gold if agent in gold square. Shoot can fires arrow in direction agent is facing. The arrow continues until it either hits (and hence kills) the wumpus or hits a wall. The agent has only one arrow, so only the first Shoot action has any effect. Climb climbs out of the cave if at [1,1].

# First Steps Wumpus World

- ▶ **Sensors**: The agent has five sensors, each of which gives a single bit of information:
    - ▶ In squares directly (not diagonally) adjacent to wumpus, agent perceives a Stench.
    - ▶ In squares directly adjacent to a pit, the agent perceives a Breeze.
    - ▶ In the square with gold, agent perceives a Glitter.
    - ▶ When an agent walks into a wall, it perceives a Bump.
    - ▶ When the wumpus is killed, it emits a Scream perceivable anywhere in the cave.

Percepts encoded as list of five symbols indicating presense or absence (by None) of:
[Stench,Breeze,Glitter,Bump,Scream] (a bit vector).



(a)                                                    (b)

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

- ▶ (a) after percept [None,None,None,None,None]
- ▶ (b) after moving to [2,1] and perceiving [None,Breeze,None,None,None]

# Later Steps Wumpus World

| 1,4 | 2,4 | 3,4 | 4,4 |
|---|---|---|---|
| 1,3 **W!** | 2,3 | 3,3 | 4,3 |
| 1,2 **A** **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 **A** **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 |

(a)                                      (b)

KENNESAW STATE UNIVERSITY

# Logic

Basics:

- ▶ **Syntax** specifies the form of sentences.
    - ▶ $x + y = 4$ is *well-formed*, but $x4y+ =$ is not.
- ▶ **Semantics** specifies the meaning of sentences.
    - ▶ $x + y = 4$ is **true** in a world where $x = 1$ and $y = 3$.
- ▶ **Model**: a formal specification of a *possible world*, that is, a set of assignments of values to the variables in the sentences of a knowledge base.
    - ▶ Given a model $\{x = 3, \ y = 2\}$, the sentence $x + y = 4$ is false.

Satisfaction:

- ▶ "$m$ satisfies $\alpha$": sentence $\alpha$ is true in model $m$, also "$m$ is a model of $\alpha$."
- ▶ $M(\alpha)$ the set of all models of $\alpha$, i.e., the set of all models in which $\alpha$ is true.
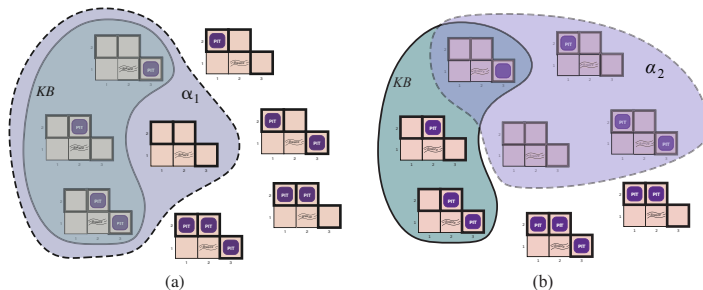
Entailment: $\alpha \models \beta$: $\beta$ *follows logically* from $\alpha$

Formal definition of entailment:

$$\alpha \models \beta \text{ if and only if } M(\alpha) \subseteq M(\beta)$$

KENNESAW STATE
UNIVERSITY

# Possible Models of Pits in Wumpus World

The presence of pits in squares [1, 2], [2, 2] and [3, 1] gives rise to $2^3 = 8$ possible models.



(a)          (b)

Solid line delineates KB based on percept [None, None, None, None, None] in [1,1].

▶ (a). $\alpha_1$ = "There is no pit in [1, 2]." Here, $KB \models \alpha_1$
▶ (b). $\alpha_2$ = "There is no pit in [2, 2]." Here, $KB \not\models \alpha_2$

Logical inference via model checking: because of (b), cannot conclude $\alpha_2$ (or $\neg\alpha_2$).

▶ $M(KB) \models \alpha_1$ but $M(KB) \not\models \alpha_2$

# Inference Algorithms

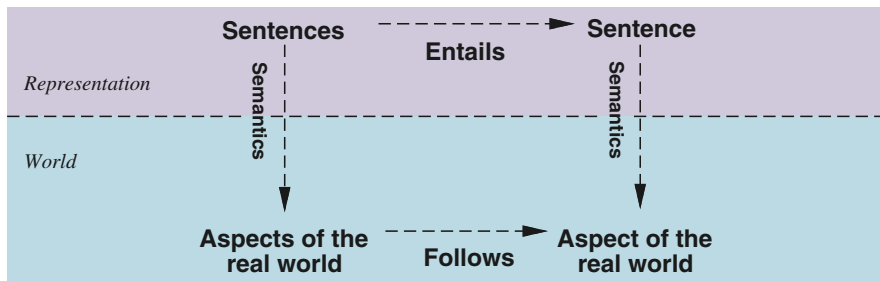If an inference algorithm $i$ can derive $\alpha$ from $KB$, we write

$$KB \vdash_i \alpha$$

which is pronounced "$\alpha$ is derived from $KB$ by $i$" or "$i$ derives $\alpha$ from $KB$."

**Important properties of inference algorithms**:

▶ An inference algorithm that derives only entailed sentences is called **sound** or **truth-preserving**.

▶ An inference algorithm is **complete** if it can derive any sentence that is entailed.

# Representation vs World

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?



**Grounding**: connection between logical reasoning and the real environment. How do we know that KB is true in the real world.

▶ Subject of volumes of philosophical investigation.
▶ For us: if agent perceives it, it is true.

# Propositional Logic

- ▶ **Atomic** sentences consist of a single proposition symbol.
- ▶ Proposition symbol stands for a proposition that can be true or false.
  - ▶ We use symbols that start with uppercase letter and may contain other letters or subscripts, e.g., : P, Q, R, $W_{1,3}$ and FacingEast
  - ▶ True and False have fixed meanings
- ▶ **Complex sentence**: one or more atomic sentences constructed from **logical connectives**.
- ▶ ¬ (not) Unary connective. $\neg W_{1,3}$ is the **negation** of $W_{1,3}$
  - ▶ A **positive literal** is an atomic sentence.
  - ▶ a \*\* negative literal\*\* is a negated atomic sentence.
- ▶ ∧ (and). Binary connective. **Conjunction**, e.g., $W_{1,3} \wedge P_{3,1}$
- ▶ ∨ (or). Binary connective. **Disjunction**, e.g., $(W_{1,3} \wedge P_{3,1}) \vee W_{2,2}$
- ▶ $\implies$ (implies). Binary connective. **Implication**, e.g., $(W_{1,3} \wedge P_{3,1}) \implies \neg W_{2,2}$
  - ▶ $(W_{1,3} \wedge P_{3,1})$ is the **premise** or **antecedent**.
  - ▶ $\neg W_{2,2}$ is the **conclusion** or **consequent**.
  - ▶ Also known as rules or if-then statements.
  - ▶ Some authors use ⊃ or →
- ▶ $\iff$ (if and only if). Binary connective. $W_{1,3} \iff \neg W_{2,2}$ is a **biconditional**

**KENNESAW STATE** UNIVERSITY

# Grammar of Propositional Logic

$$
\begin{aligned}
\textit{Sentence} \quad &\rightarrow \quad \textit{AtomicSentence} \mid \textit{ComplexSentence} \\
\textit{AtomicSentence} \quad &\rightarrow \quad \textit{True} \mid \textit{False} \mid P \mid Q \mid R \mid \ldots \\
\textit{ComplexSentence} \quad &\rightarrow \quad (\ \textit{Sentence}\ ) \\
&\quad\mid\quad \neg\ \textit{Sentence} \\
&\quad\mid\quad \textit{Sentence} \wedge \textit{Sentence} \\
&\quad\mid\quad \textit{Sentence} \vee \textit{Sentence} \\
&\quad\mid\quad \textit{Sentence} \Rightarrow \textit{Sentence} \\
&\quad\mid\quad \textit{Sentence} \Leftrightarrow \textit{Sentence}
\end{aligned}
$$

OPERATOR PRECEDENCE : $\neg, \wedge, \vee, \Rightarrow, \Leftrightarrow$

KENNESAW STATE UNIVERSITY

# Semantics of Propositional Logic

| $P$ | $Q$ | $\neg P$ | $P \wedge Q$ | $P \vee Q$ | $P \Rightarrow Q$ | $P \Leftrightarrow Q$ |
|---|---|---|---|---|---|---|
| *false* | *false* | *true* | *false* | *false* | *true* | *true* |
| *false* | *true* | *true* | *false* | *true* | *true* | *false* |
| *true* | *false* | *false* | *false* | *true* | *false* | *false* |
| *true* | *true* | *false* | *true* | *true* | *true* | *true* |

KENNESAW STATE
UNIVERSITY

# Propositional Theorem Proving

So far we've done **model checking**: enumerating models and showing that the sentence must hold in all models.

Now we turn to **theorem proving**: applying rules of inference directly to the sentences in our knowledge base to construct a proof of the desired sentence without consulting models.

Some basic concepts:

- **Logical equivalence**: two sentences $\alpha$ and $\beta$ are logically equivalent if they are true in the same set of models. $\alpha \equiv \beta$

  - $\alpha \equiv \beta$ if and only if $\alpha \models \beta$ and $\beta \models \alpha$

- **Validity**: A sentence is valid if it is true in all models. For example, the sentence $P \wedge \neg P$ is valid. Valid sentences are also known as **tautologies**.

**Deduction theorem**:

For any sentences $\alpha$ and $beta$ , $\alpha \models \beta$ if and only iff the sentence $(\alpha \implies \beta)$ is valid.

# Inference Rules

General form:

$$\frac{Givens}{Conclusions}$$

Modus ponens

# Propositional Theorem Proving

| $B_{1,1}$ | $B_{2,1}$ | $P_{1,1}$ | $P_{1,2}$ | $P_{2,1}$ | $P_{2,2}$ | $P_{3,1}$ | $R_1$ | $R_2$ | $R_3$ | $R_4$ | $R_5$ | $KB$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| false | false | false | false | false | false | false | true | true | true | true | false | false |
| false | false | false | false | false | false | true | true | true | false | true | false | false |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| false | true | false | false | false | false | false | true | true | false | true | true | false |
| false | true | false | false | false | false | true | true | true | true | true | true | <u>*true*</u> |
| false | true | false | false | false | true | false | true | true | true | true | true | <u>*true*</u> |
| false | true | false | false | false | true | true | true | true | true | true | true | <u>*true*</u> |
| false | true | false | false | true | false | false | true | false | false | true | true | false |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| true | true | true | true | true | true | true | false | true | true | false | true | false |

## Propositional Theorem Proving

**function** TT-ENTAILS?(*KB*, $\alpha$) **returns** *true* or *false*
  **inputs**: *KB*, the knowledge base, a sentence in propositional logic
           $\alpha$, the query, a sentence in propositional logic

  *symbols* ← 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*    // when KB is false, always return true
  **else**
    *P* ← FIRST(*symbols*)
    *rest* ← REST(*symbols*)
    **return** (TT-CHECK-ALL(*KB*, $\alpha$, *rest*, *model* ∪ {*P* = *true*})
        **and**
        TT-CHECK-ALL(*KB*, $\alpha$, *rest*, *model* ∪ {*P* = *false*})

KENNESAW STATE
UNIVERSITY

## Propositional Theorem Proving

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

## Propositional Theorem Proving

$$
\begin{aligned}
CNFSentence &\rightarrow Clause_1 \wedge \cdots \wedge Clause_n \\
Clause &\rightarrow Literal_1 \vee \cdots \vee Literal_m \\
Fact &\rightarrow Symbol \\
Literal &\rightarrow Symbol \mid \neg Symbol \\
Symbol &\rightarrow P \mid Q \mid R \mid \ldots \\
HornClauseForm &\rightarrow DefiniteClauseForm \mid GoalClauseForm \\
DefiniteClauseForm &\rightarrow Fact \mid (Symbol_1 \wedge \cdots \wedge Symbol_l) \Rightarrow Symbol \\
GoalClauseForm &\rightarrow (Symbol_1 \wedge \cdots \wedge Symbol_l) \Rightarrow False
\end{aligned}
$$

## Propositional Theorem Proving

**function** PL-RESOLUTION($KB, \alpha$) **returns** *true* or *false*
  **inputs**: *KB*, the knowledge base, a sentence in propositional logic
          $\alpha$, the query, a sentence in propositional logic

  *clauses* ← the set of clauses in the CNF representation of $KB \wedge \neg\alpha$
  *new* ← { }
  **while** *true* **do**
    **for each** pair of clauses $C_i, C_j$ in *clauses* **do**
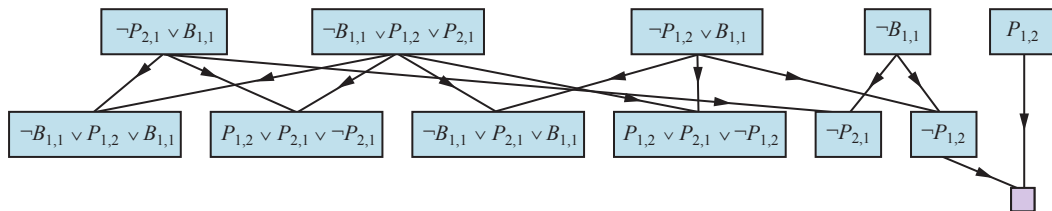      *resolvents* ← PL-RESOLVE($C_i, C_j$)
      **if** *resolvents* contains the empty clause **then return** *true*
      *new* ← *new* ∪ *resolvents*
    **if** *new* ⊆ *clauses* **then return** *false*
    *clauses* ← *clauses* ∪ *new*

# Propositional Theorem Proving

## Propositional Theorem Proving

**function** PL-FC-ENTAILS?($KB, q$) **returns** *true* or *false*
  **inputs**: $KB$, the knowledge base, a set of propositional definite clauses
        $q$, the query, a proposition symbol
  *count* ← a table, where *count*[$c$] is initially the number of symbols in clause $c$'s premise
  *inferred* ← a table, where *inferred*[$s$] is initially *false* for all symbols
  *queue* ← a queue of symbols, initially symbols known to be true in $KB$

  **while** *queue* is not empty **do**
    $p$ ← POP(*queue*)
    **if** $p = q$ **then return** *true*
    **if** *inferred*[$p$] = *false* **then**
      *inferred*[$p$] ← *true*
      **for each** clause $c$ in $KB$ where $p$ is in $c$.PREMISE **do**
        decrement *count*[$c$]
        **if** *count*[$c$] = 0 **then** add $c$.CONCLUSION to *queue*
  **return** *false*

KENNESAW STATE
UNIVERSITY

# Propositional Theorem Proving

$P \Rightarrow Q$

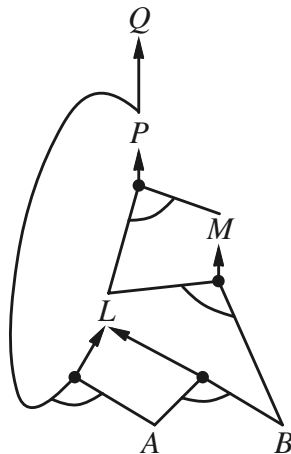$L \wedge M \Rightarrow P$

$B \wedge L \Rightarrow M$

$A \wedge P \Rightarrow L$

$A \wedge B \Rightarrow L$

$A$

$B$

(a)



(b)

## Propositional Model Checking

**function** DPLL-SATISFIABLE?(*s*) **returns** *true* or *false*
  **inputs**: *s*, a sentence in propositional logic

  *clauses* ← the set of clauses in the CNF representation of *s*
  *symbols* ← a list of the proposition symbols in *s*
  **return** DPLL(*clauses*, *symbols*, { })

**function** DPLL(*clauses*, *symbols*, *model*) **returns** *true* or *false*

  **if** every clause in *clauses* is true in *model* **then return** *true*
  **if** some clause in *clauses* is false in *model* **then return** *false*
  *P*, *value* ← FIND-PURE-SYMBOL(*symbols*, *clauses*, *model*)
  **if** *P* is non-null **then return** DPLL(*clauses*, *symbols* – *P*, *model* ∪ {*P=value*})
  *P*, *value* ← FIND-UNIT-CLAUSE(*clauses*, *model*)
  **if** *P* is non-null **then return** DPLL(*clauses*, *symbols* – *P*, *model* ∪ {*P=value*})
  *P* ← FIRST(*symbols*); *rest* ← REST(*symbols*)
  **return** DPLL(*clauses*, *rest*, *model* ∪ {*P=true*}) **or**
       DPLL(*clauses*, *rest*, *model* ∪ {*P=false*})

## Propositional Model Checking

**function** WALKSAT(*clauses*, *p*, *max_flips*) **returns** a satisfying model or *failure*
  **inputs**: *clauses*, a set of clauses in propositional logic
        *p*, the probability of choosing to do a "random walk" move, typically around 0.5
        *max_flips*, number of value flips allowed before giving up

  *model* ← a random assignment of *true/false* to the symbols in *clauses*
  **for each** $i = 1$ **to** *max_flips* **do**
    **if** *model* satisfies *clauses* **then return** *model*
    *clause* ← a randomly selected clause from *clauses* that is false in *model*
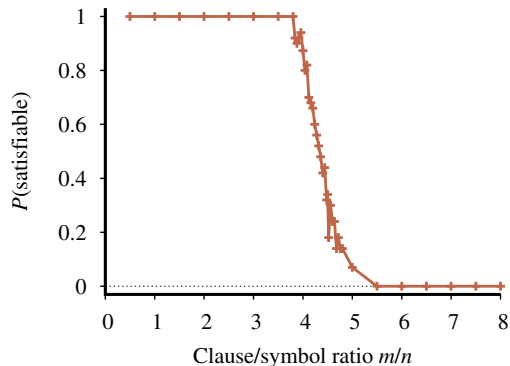    **if** RANDOM(0, 1) $\leq$ *p* **then**
      flip the value in *model* of a randomly selected symbol from *clause*
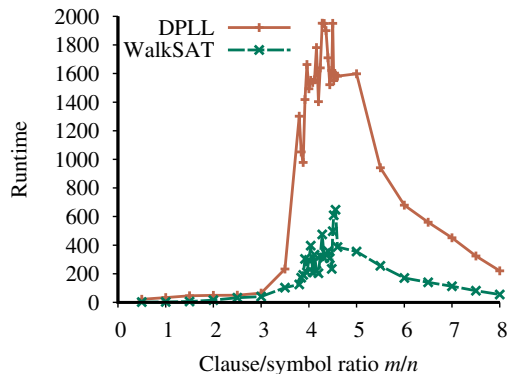    **else** flip whichever symbol in *clause* maximizes the number of satisfied clauses
  **return** *failure*

# Propositional Model Checking



(a)                                                                (b)
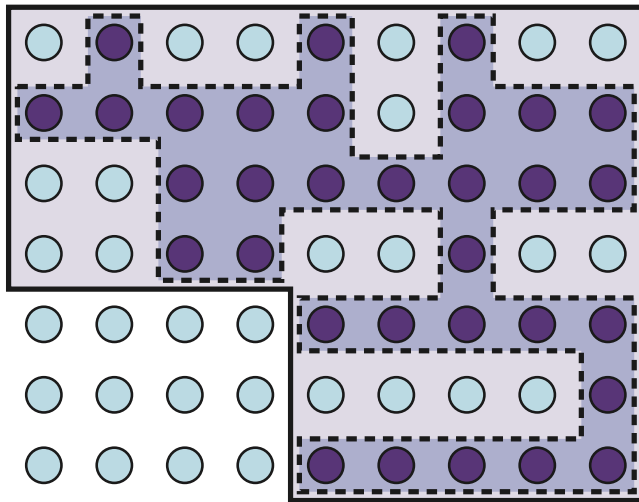
# Agents Based on Propositional Logic

**function** HYBRID-WUMPUS-AGENT(*percept*) **returns** an *action*
  **inputs**: *percept*, a list, [*stench,breeze,glitter,bump,scream*]
  **persistent**: *KB*, a knowledge base, initially the atemporal "wumpus physics"
        *t*, a counter, initially 0, indicating time
        *plan*, an action sequence, initially empty

  TELL(*KB*, MAKE-PERCEPT-SENTENCE(*percept*, *t*))
  TELL the *KB* the temporal "physics" sentences for time *t*
  *safe* ← {[*x*, *y*] : ASK(*KB*, $OK^t_{x,y}$) = *true*}
  **if** ASK(*KB*, *Glitter$^t$*) = *true* **then**
    *plan* ← [*Grab*] + PLAN-ROUTE(*current*, {[1,1]}, *safe*) + [*Climb*]
  **if** *plan* is empty **then**
    *unvisited* ← {[*x*, *y*] : ASK(*KB*, $L^{t'}_{x,y}$) = *false* for all *t'* ≤ *t*}
    *plan* ← PLAN-ROUTE(*current*, *unvisited* ∩ *safe*, *safe*)
  **if** *plan* is empty and ASK(*KB*, *HaveArrow$^t$*) = *true* **then**
    *possible_wumpus* ← {[*x*, *y*] : ASK(*KB*, ¬ $W_{x,y}$) = *false*}
    *plan* ← PLAN-SHOT(*current*, *possible_wumpus*, *safe*)
  **if** *plan* is empty **then**      // no choice but to take a risk
    *not_unsafe* ← {[*x*, *y*] : ASK(*KB*, ¬ $OK^t_{x,y}$) = *false*}
    *plan* ← PLAN-ROUTE(*current*, *unvisited* ∩ *not_unsafe*, *safe*)
  **if** *plan* is empty **then**
    *plan* ← PLAN-ROUTE(*current*, {[1,1]}, *safe*) + [*Climb*]
  *action* ← POP(*plan*)
  TELL(*KB*, MAKE-ACTION-SENTENCE(*action*, *t*))
  *t* ← *t* + 1
  **return** *action*

**function** PLAN-ROUTE(*current*, *goals*, *allowed*) **returns** an action sequence
  **inputs**: *current*, the agent's current position
      *goals*, a set of squares; try to plan a route to one of them
      *allowed*, a set of squares that can form part of the route

  *problem* ← ROUTE-PROBLEM(*current*, *goals*, *allowed*)
  **return** SEARCH(*problem*)    // Any search algorithm from Chapter 3

**KENNESAW STATE**
UNIVERSITY

# Agents Based on Propositional Logic

# Agents Based on Propositional Logic

**function** SATPLAN(*init*, *transition*, *goal*, $T_{max}$) **returns** solution or *failure*
  **inputs**: *init*, *transition*, *goal*, constitute a description of the problem
        $T_{max}$, an upper limit for plan length

  **for** $t = 0$ **to** $T_{max}$ **do**
    *cnf* ← TRANSLATE-TO-SAT(*init*, *transition*, *goal*, *t*)
    *model* ← SAT-SOLVER(*cnf*)
    **if** *model* is not null **then**
      **return** EXTRACT-SOLUTION(*model*)
  **return** *failure*