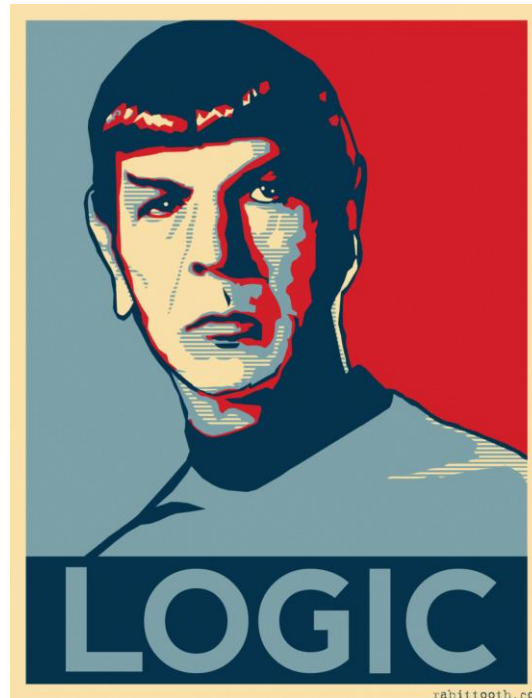


# AI: Representation and Problem Solving



## Logical Agent Algorithms



Instructor: Pat Virtue

Slide credits: CMU AI, <http://ai.berkeley.edu>

# Plan

## Last Time:

- Propositional logic
- Models and Knowledge Bases
- Satisfiability and Entailment

## Today: Logical Agent Algorithms

- Entailment
  - Model checking: Truth table entailment
  - Theorem proving: (Forward chaining), resolution
- Satisfiability: DPLL
- Planning with logic

# Plan

## Last Time:

- Propositional logic
- Models and Knowledge Bases
- Satisfiability and Entailment

# Propositional Logic Vocab

## Literal

- Atomic sentence: True, False, Symbol,  $\neg$ Symbol



## Clause

- Disjunction (OR) of literals:  $A \vee B \vee \neg C$



## Definite clause

- Disjunction (OR) of literals, *exactly one* is positive

- $\neg A \vee B \vee \neg C$



## Horn clause

- Disjunction of literals, *at most one* is positive
- All definite clauses are Horn clauses

# PL: Conjunctive Normal Form (CNF)

Every sentence can be expressed as a conjunction (AND) of clauses

Each clause is a disjunction (OR) of literals

Each literal is a symbol or a negated symbol

- Example:  $(\neg A \vee \neg C \vee B) \wedge (\neg A \vee \neg B \vee C) \wedge (x \vee x \vee x) \wedge \dots$
-

# PL: Conjunctive Normal Form (CNF)

Every sentence can be expressed as a **conjunction of clauses**

Each clause is a **disjunction of literals**

Each literal is a symbol or a negated symbol

Conversion to CNF by a sequence of standard transformations:

- $At_{1,1_0} \Rightarrow (Wall_{0,1} \Leftrightarrow Blocked\_W\_0)$
- $At_{1,1_0} \Rightarrow ((Wall_{0,1} \Rightarrow Blocked\_W\_0) \wedge (Blocked\_W\_0 \Rightarrow Wall_{0,1}))$
- $\neg At_{1,1_0} \vee ((\neg Wall_{0,1} \vee Blocked\_W\_0) \wedge (\neg Blocked\_W\_0 \vee Wall_{0,1}))$
- $(\neg At_{1,1_0} \vee \neg Wall_{0,1} \vee Blocked\_W\_0) \wedge (\neg At_{1,1_0} \vee \neg Blocked\_W\_0 \vee Wall_{0,1})$

# PL: Conjunctive Normal Form (CNF)

Every sentence can be expressed

Replace biconditional by two implications

Each clause is a **disjunction** of **literal**

Replace  $\alpha \Rightarrow \beta$  by  $\neg\alpha \vee \beta$

Each literal is a symbol or a negated symbol

Distribute  $\vee$  over  $\wedge$

Conversion to CNF by a sequence of standard transformations:

- $At_{1,1,0} \Rightarrow (Wall_{0,1} \Leftrightarrow Blocked_{W_0})$
- $At_{1,1,0} \Rightarrow ((Wall_{0,1} \Rightarrow Blocked_{W_0}) \wedge (Blocked_{W_0} \Rightarrow Wall_{0,1}))$
- $\neg At_{1,1,0} \vee ((\neg Wall_{0,1} \vee Blocked_{W_0}) \wedge (\neg Blocked_{W_0} \vee Wall_{0,1}))$
- $(\neg At_{1,1,0} \vee \neg Wall_{0,1} \vee Blocked_{W_0}) \wedge (\neg At_{1,1,0} \vee \neg Blocked_{W_0} \vee Wall_{0,1})$

# Logical Agent Vocab

## Model

- Complete assignment of symbols to True/False

## Sentence

- Logical statement
- Composition of logic symbols and operators

## KB

- Collection of sentences representing facts and rules we know about the world

## Query

- Sentence we want to know if it is *provably* True, *provably* False, or *unsure*.



# Provably True, Provably False, or Unsure



<http://thiagodnf.github.io/wumpus-world-simulator/>

# Logical Agent Vocab

- Entailment *KB* *z* *definitely*
- Input: sentence1, sentence2
  - Each model that satisfies sentence1 must also satisfy sentence2
  - "If I know 1 holds, then I know 2 holds"
  - (ASK), TT-ENTAILS, FC-ENTAILS, RESOLUTION-ENTAILS

- Satisfy *possible*
- Input: model, sentence → *T or F*
  - Is this sentence true in this model?
  - Does this model *satisfy* this sentence
  - "Does this particular state of the world work?"
  - PL-TRUE ]

# Logical Agent Vocab

## Satisfiable

- Input: sentence KB  $\rightarrow$  model
- Can find at least one model that satisfies this sentence
  - (We often want to know what that model is)
- "Is it possible to make this sentence true?"
- DPLL  $\leftarrow$  CSP

CSP

## Valid

- Input: sentence
- sentence is true in all possible models

# Outline

## Logical Agent Algorithms

- Vocab
- PL\_TRUE
- Entailment
- ■ Model checking: Truth table entailment
  - Theorem proving:
    - (Forward chaining), resolution
- Satisfiability: DPLL CSP
- Planning with logic

# Propositional Logic

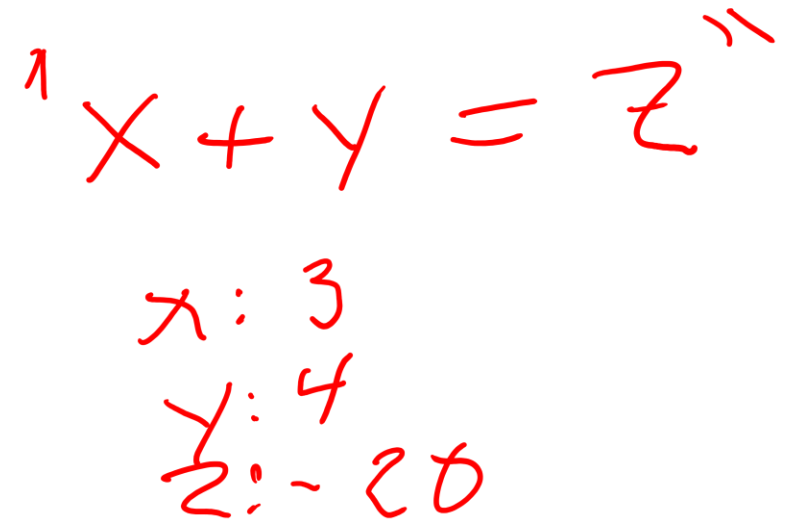
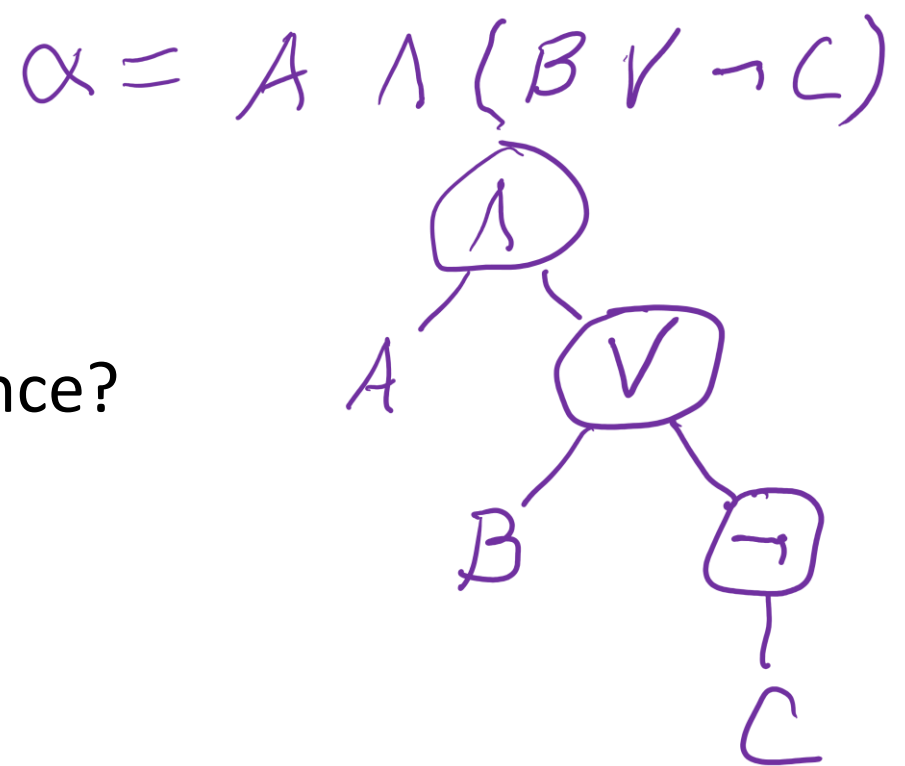
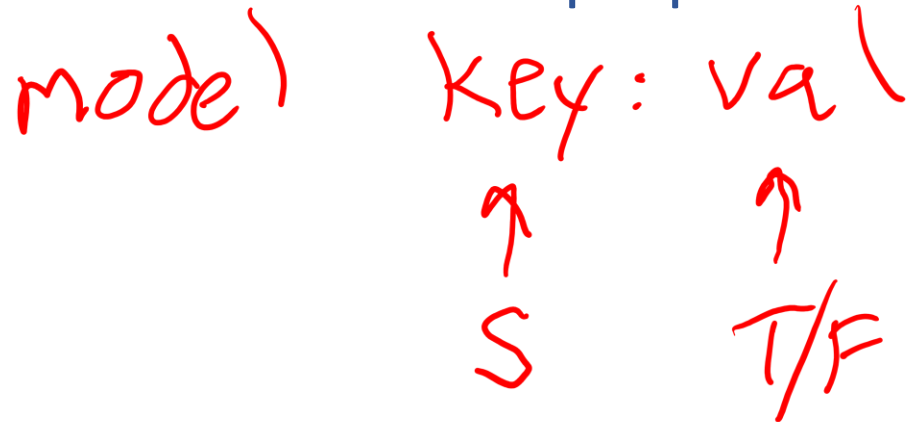
Check if sentence is true in given model

In other words, does the model *satisfy* the sentence?

function **PL-TRUE?**( $\alpha$ , model) returns true or false



But are models and propositional logic sentences  $\alpha$  represented?



# Propositional Logic

Check if sentence is true in given model

In other words, does the model *satisfy* the sentence?

function **PL-TRUE?**( $\alpha$ , model) returns true or false

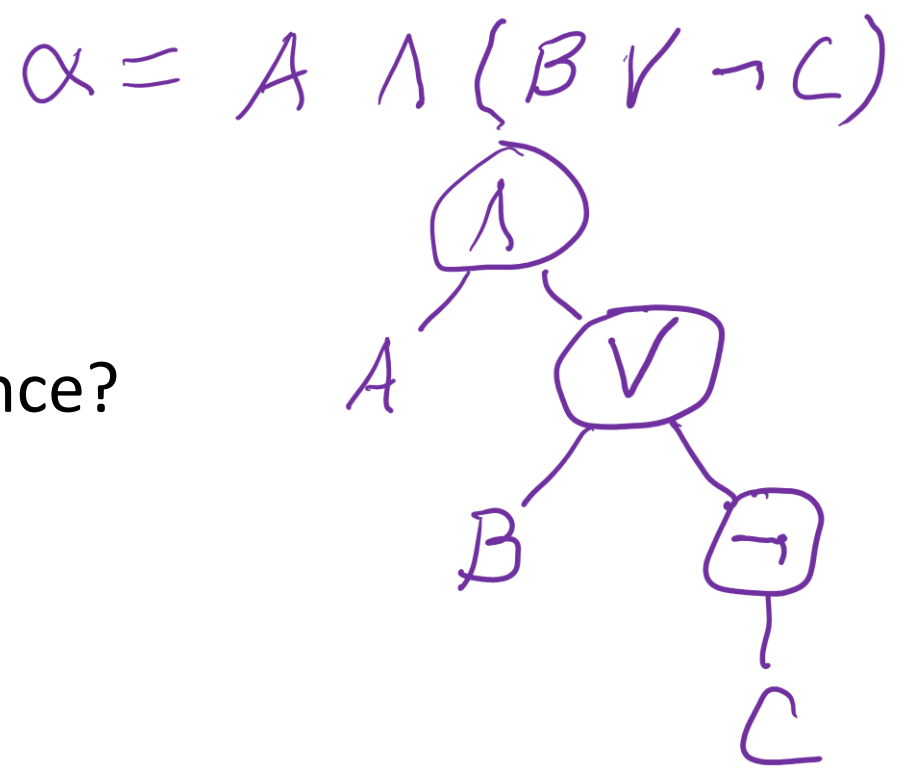
if  $\alpha$  is a symbol then return **Lookup**( $\alpha$ , model)

if  $\text{Op}(\alpha) = \neg$  then return **not**(**PL-TRUE?**(**Arg1**( $\alpha$ ), model))

if  $\text{Op}(\alpha) = \wedge$  then return **and**(**PL-TRUE?**(**Arg1**( $\alpha$ ), model),  
**PL-TRUE?**(**Arg2**( $\alpha$ ), model))


etc.

(Sometimes called “recursion over syntax”)



# Outline

## Logical Agent Algorithms

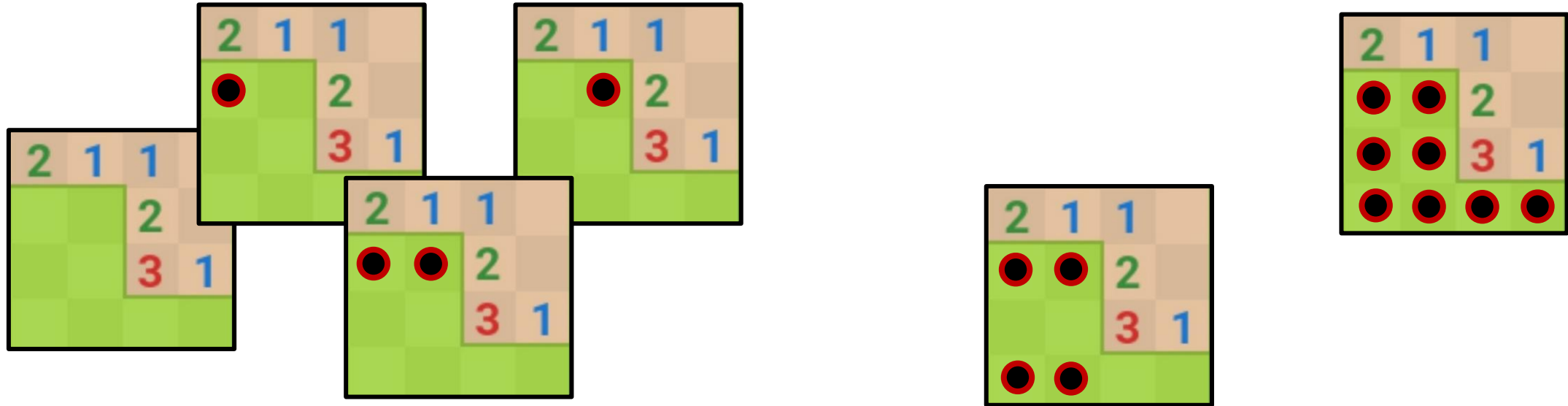
- Vocab
-  PL\_TRUE
- Entailment
  - Model checking: Truth table entailment
  - Theorem proving:
    - (Forward chaining), resolution
- Satisfiability: DPLL
- Planning with logic

# Inference: Proofs

A proof is a *demonstration* of entailment between  $\alpha$  and  $\beta$

Method 1: *model-checking*

- For every possible world, if  $\alpha$  is true make sure that is  $\beta$  true too
- OK for propositional logic (finitely many worlds); not easy for first-order logic





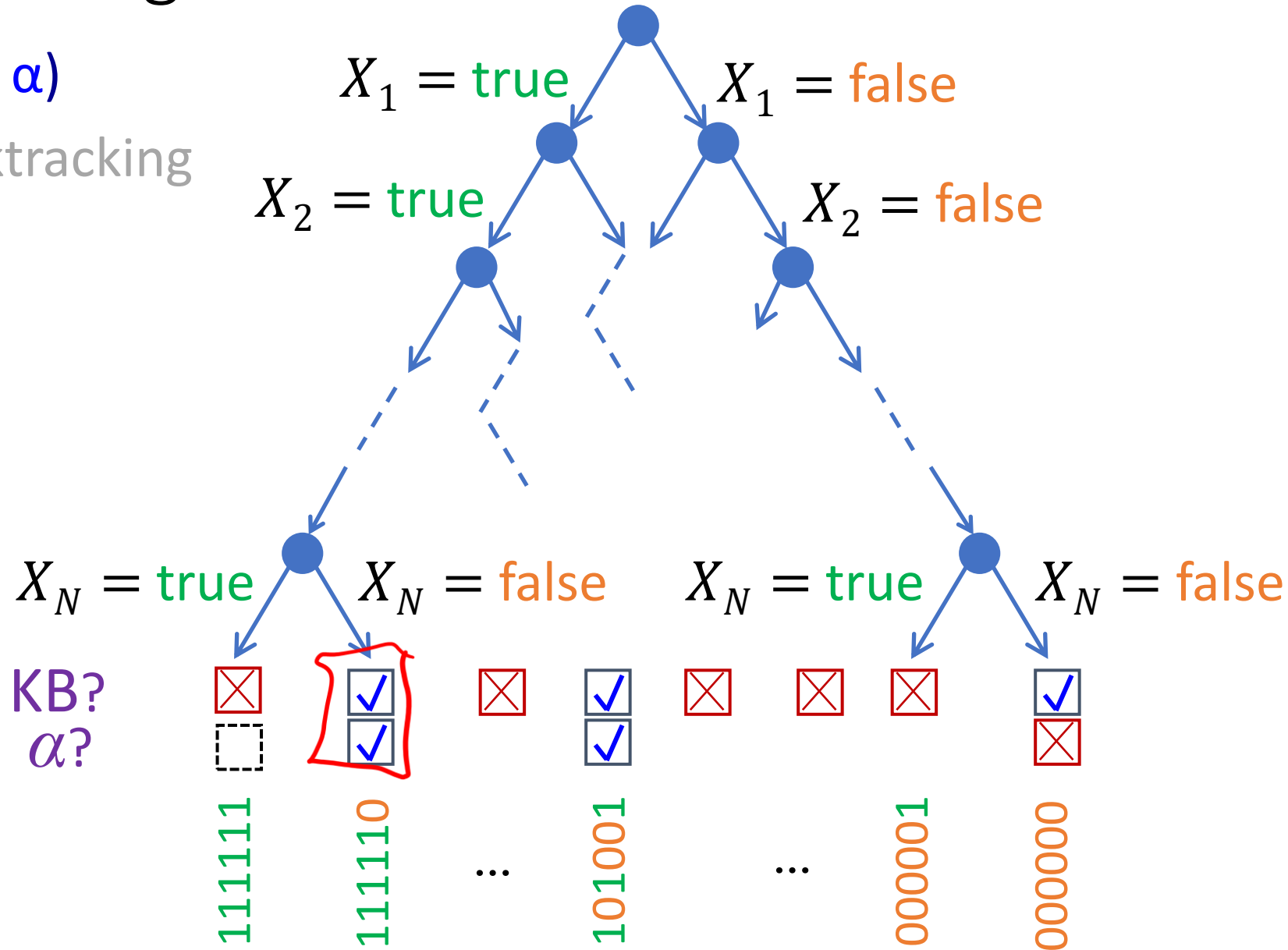
# Simple Model Checking

function **TT-ENTAILS?**(KB,  $\alpha$ ) Returns true or false

# Simple Model Checking

function **TT-ENTAILS?**(KB,  $\alpha$ )

Same recursion as backtracking



# Simple Model Checking

function **TT-ENTAILS?**(KB,  $\alpha$ ) Returns true or false

return **TT-CHECK-ALL**(KB,  $\alpha$ , symbols(KB)  $\cup$  symbols( $\alpha$ ), {})

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

Recursively check to make sure all models

 that satisfy the KB also satisfy  $\alpha$



# Simple Model Checking

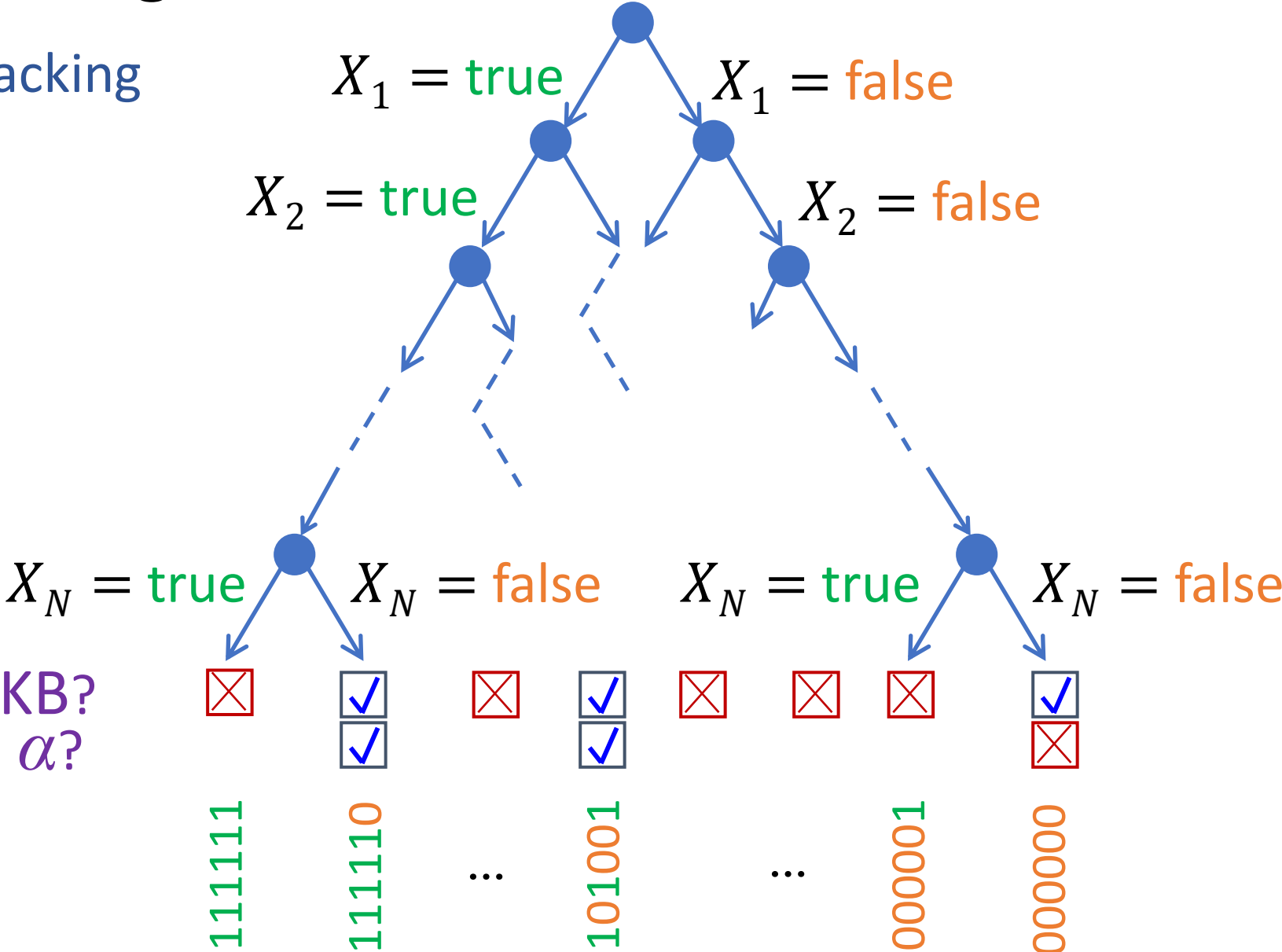
```
function TT-CHECK-ALL(KB,  $\alpha$ , symbols, model)  Returns true or false
  if empty?(symbols) then
    if PL-TRUE?(KB, model) then
      return PL-TRUE?( $\alpha$ , model)
    else
      return true
  else
     $X_i \leftarrow$  first(symbols)
    rest  $\leftarrow$  rest(symbols)
    return and ( TT-CHECK-ALL(KB,  $\alpha$ , rest, model  $\cup$  { $X_i =$  true})
                TT-CHECK-ALL(KB,  $\alpha$ , rest, model  $\cup$  { $X_i =$  false}) )
```

# Simple Model Checking

Same recursion as backtracking

$O(2^N)$  time, linear space

Can we do better?



# Inference: Proofs

A proof is a *demonstration* of entailment between  $\alpha$  and  $\beta$

## Method 1: *model-checking*

- For every possible world, if  $\alpha$  is true make sure that  $\beta$  is true too
- OK for propositional logic (finitely many worlds); not easy for first-order logic

## Method 2: *theorem-proving*

- Search for a sequence of proof steps (applications of *inference rules*) leading from  $\alpha$  to  $\beta$
- E.g., from  $P \wedge (P \Rightarrow Q)$ , infer  $Q$  by *Modus Ponens*

## Properties

- *Sound* algorithm: everything it claims to prove is in fact entailed
- *Complete* algorithm: every sentence that is entailed can be proved

# Simple Theorem Proving: Forward Chaining

Forward chaining applies Modus Ponens to generate new facts:

- Given  $X_1 \wedge X_2 \wedge \dots \wedge X_n \Rightarrow Y$  and  $X_1, X_2, \dots, X_n$
- Infer Y

$\mathcal{L} \Rightarrow \mathcal{R}$

Forward chaining keeps applying this rule, adding new facts, until nothing more can be added

# Forward Chaining Algorithm

function PL-FC-ENTAILS?(KB, q) Returns true or false

**KB CLAUSES**

$P \Rightarrow Q$	L M N Q
<u><math>L \wedge M \Rightarrow P</math></u>	
<u><math>B \wedge L \Rightarrow M</math></u>	
$A \wedge P \Rightarrow L$	
<u><math>A \wedge B \Rightarrow L</math></u>	
$A$	
$B$	

KB = ?



# Properties

## Forward Chaining is:

- Sound and complete for definite-clause KBs
- Complexity: linear time 😊

Resolution is another theorem-proving algorithm that is:

- Sound and complete for **any PL** KBs!
- Complexity: exponential time 😞

## Vocab Reminder

### Literal

- Atomic sentence:  
T, F, Symbol,  $\neg$ Symbol

### Clause

- Disjunction of literals:  
 $A \vee B \vee \neg C$

### Definite clause

- Disjunction of literals, *exactly one* is positive  
 $\neg A \vee B \vee \neg C$

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

# Inference Rules

Modus Ponens

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

Forward Chaining Notation Alert!

Unit Resolution

$$\frac{(a \vee b) \wedge (\neg b \vee c)}{a \vee c}$$

if  $b = F$   
→  $a$  must be  $T$   
else  $b = T$   
→  $b$  is  $F$   
 $c$  must be  $T$

General Resolution

$$\frac{a_1 \vee \dots \vee a_m \vee b, \neg b \vee c_1 \vee \dots \vee c_n}{a_1 \vee \dots \vee a_m \vee c_1 \vee \dots \vee c_n}$$

# Resolution

## Algorithm Overview

function PL-RESOLUTION?(KB,  $\alpha$ ) returns true or false

$$KB \models \alpha$$

$$\text{SAT}(KB \wedge \neg \alpha) \rightarrow \text{No}$$

We want to prove that KB entails  $\alpha$

In other words, we want to prove that we cannot satisfy (KB and **not**  $\alpha$ )

1. Start with a set of CNF clauses, including the KB as well as  $\neg \alpha$
2. Keep resolving pairs of clauses until

→ A. You resolve the empty clause

Contradiction found!

KB  $\wedge$   $\neg \alpha$  cannot be satisfied

Return true, KB entails  $\alpha$

B. No new clauses added

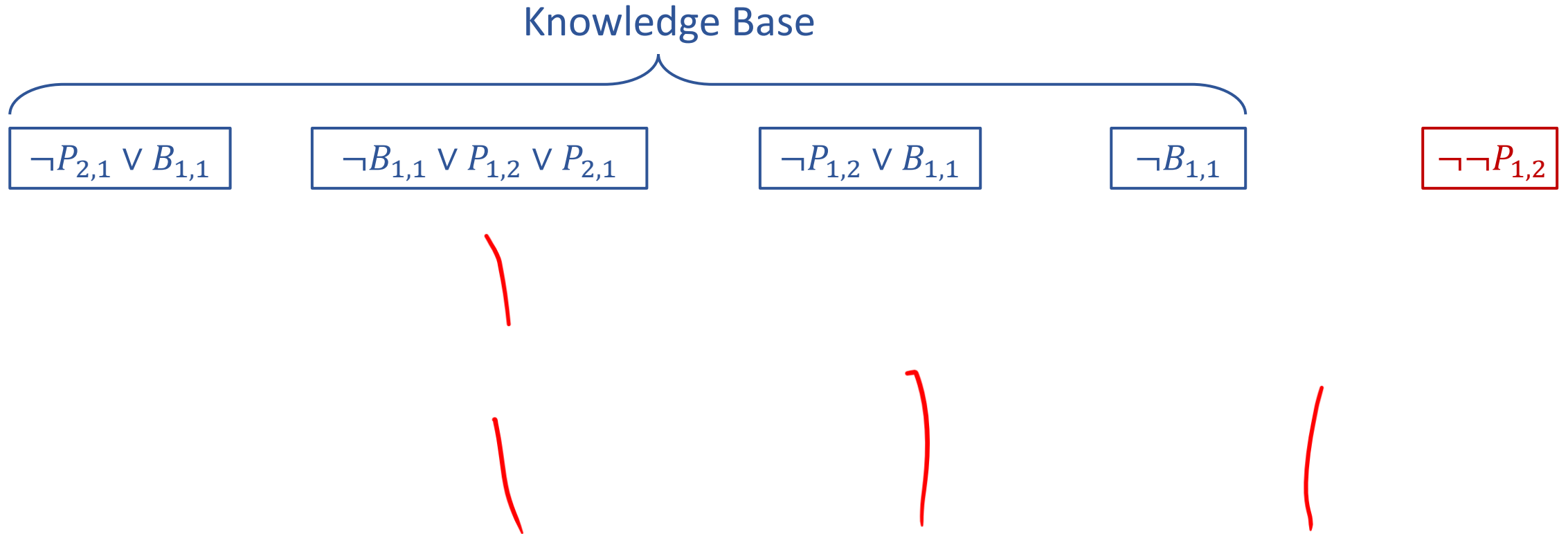
Return false, KB does not entail  $\alpha$

# Resolution

Example trying to prove  $\neg P_{1,2}$

General Resolution

$$\frac{a_1 \vee \dots \vee a_m \vee b, \quad \neg b \vee c_1 \vee \dots \vee c_n}{a_1 \vee \dots \vee a_m \vee c_1 \vee \dots \vee c_n}$$

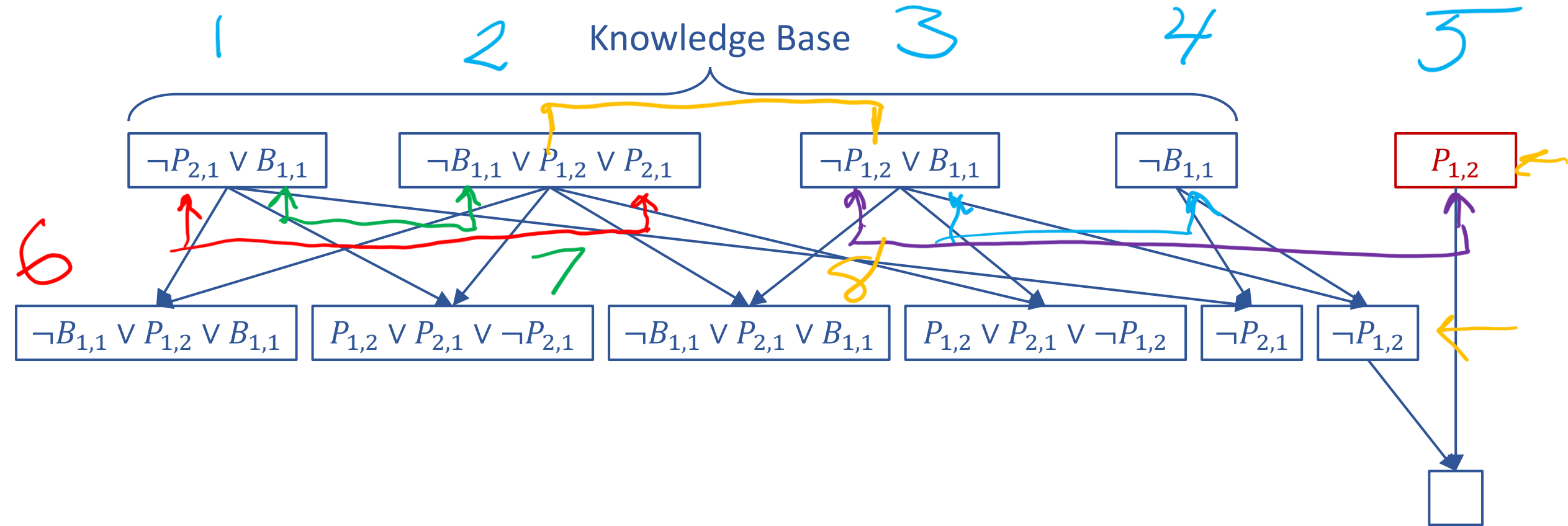


# Resolution

Example trying to prove  $\neg P_{1,2}$

General Resolution

$$\frac{a_1 \vee \dots \vee a_m \vee b, \quad \neg b \vee c_1 \vee \dots \vee c_n}{a_1 \vee \dots \vee a_m \vee c_1 \vee \dots \vee c_n}$$



# Resolution

function PL-RESOLUTION?(KB,  $\alpha$ ) returns true or false

clauses  $\leftarrow$  the set of clauses in the CNF representation of  $KB \wedge \neg\alpha$

new  $\leftarrow \{ \}$

loop do

for each pair of clauses  $C_i, C_j$  in clauses do

resolvents  $\leftarrow$  PL-RESOLVE( $C_i, C_j$ )

if resolvents contains the empty clause then

return true 

new  $\leftarrow$  new  $\cup$  resolvents

if new  $\subseteq$  clauses then

return false 

clauses  $\leftarrow$  clauses  $\cup$  new

# Properties

## Forward Chaining is:

- Sound and complete for definite-clause KBs
- Complexity: linear time 😊

## Resolution is another theorem-proving algorithm that is:

- Sound and complete for any PL KBs!
- Complexity: exponential time 😞

## Vocab Reminder

### Literal

- Atomic sentence:  
T, F, Symbol,  $\neg$ Symbol

### Clause

- Disjunction of literals:  
 $A \vee B \vee \neg C$

### Definite clause

- Disjunction of literals, *exactly one* is positive  
 $\neg A \vee B \vee \neg C$

# Outline

## Logical Agent Algorithms

- Vocab
- PL\_TRUE
- Entailment
  - Model checking: Truth table entailment
  - Theorem proving:
    - (Forward chaining), resolution
- Satisfiability: DPLL
- Planning with logic



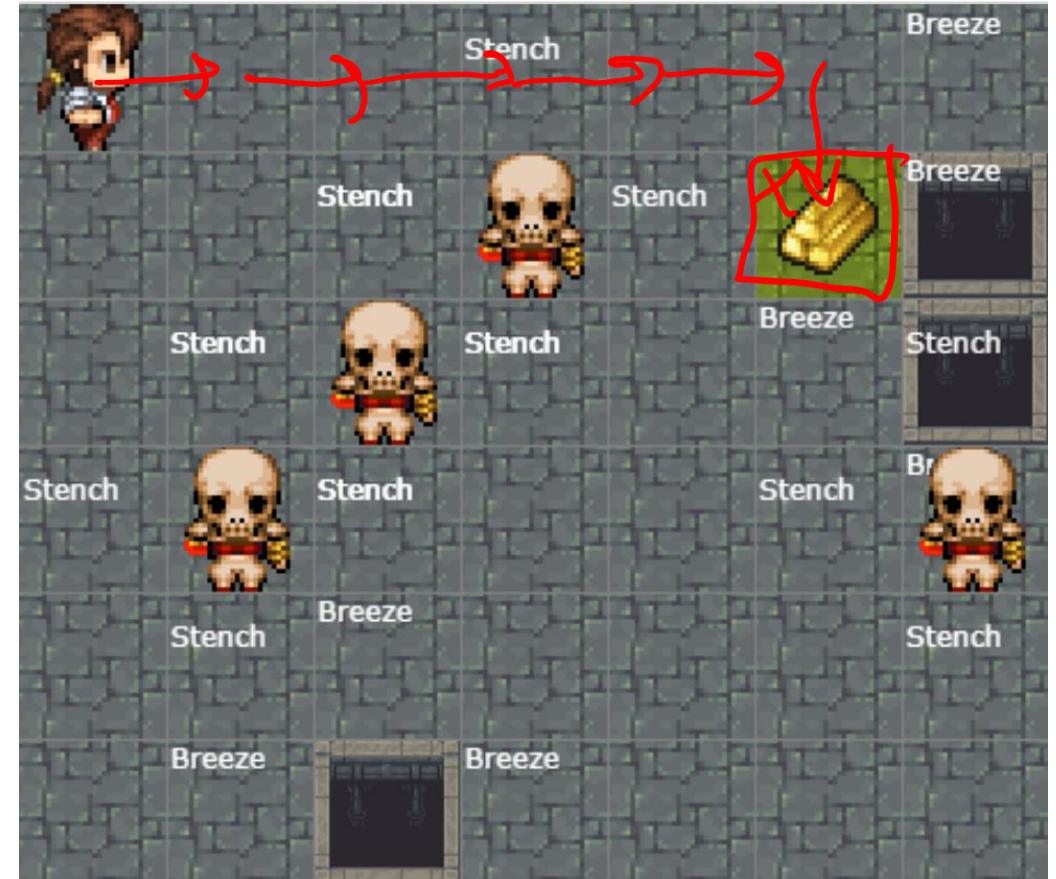
# Satisfiability and Entailment

A sentence is *satisfiable* if it is true in at least one world (e.g. CSPs!)

$$KB \wedge (G_x \wedge A_x)$$



$$\begin{array}{l} E_1 : T \\ E_2 : T \\ \vdots \\ E_t : T \\ S_{t+1} : T \end{array}$$



# Satisfiability and Entailment

A sentence is *satisfiable* if it is true in at least one world (cf CSPs!)

Suppose we have a hyper-efficient SAT solver; how can we use it to test entailment?

- Suppose  $\alpha \models \beta$
- Then  $\alpha \Rightarrow \beta$  is true in all worlds
- Hence  $\neg(\alpha \Rightarrow \beta)$  is false in all worlds
- Hence  $\alpha \wedge \neg\beta$  is false in all worlds, i.e., unsatisfiable

$KB \wedge \neg q$

So, add the negated conclusion to what you know, test for (un)satisfiability; also known as *reductio ad absurdum*

Efficient SAT solvers operate on *conjunctive normal form*



# Efficient SAT solvers

DPLL (Davis-Putnam-Logemann-Loveland) is the core of modern solvers

Essentially a backtracking search over models with some extras:

- *Early termination*: stop if
  - all clauses are satisfied; e.g.,  $(A \vee B) \wedge (A \vee \neg C)$  is satisfied by  $\{A=\text{true}\}$
  - any clause is falsified; e.g.,  $(A \vee B) \wedge (A \vee \neg C)$  is satisfied by  $\{A=\text{false}, B=\text{false}\}$
- *Pure literals*: if all occurrences of a symbol in as-yet-unsatisfied clauses have the same sign, then give the symbol that value
  - E.g.,  $A$  is pure and positive in  $(A \vee B) \wedge (A \vee \neg C) \wedge (C \vee \neg B)$  so set it to **true**
- *Unit clauses*: if a clause is left with a single literal, set symbol to satisfy clause
  - E.g., if  $A=\text{false}$ ,  $(A \vee B) \wedge (A \vee \neg C)$  becomes  $(\text{false} \vee B) \wedge (\text{false} \vee \neg C)$ , i.e.  $(B) \wedge (\neg C)$
  - Satisfying the unit clauses often leads to further propagation, new unit clauses, etc.

# DPLL algorithm

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  $\leftarrow$  FIND-PURE-SYMBOL(symbols, clauses, model)  
if P is non-null then return DPLL(clauses, symbols-P, modelU{P=value})

P, value  $\leftarrow$  FIND-UNIT-CLAUSE(clauses, model)  
if P is non-null then return DPLL(clauses, symbols-P, modelU{P=value})

P  $\leftarrow$  First(symbols)  
rest  $\leftarrow$  Rest(symbols)

return or(DPLL(clauses, rest, modelU{P=true}),  
          DPLL(clauses, rest, modelU{P=false}))

# Outline

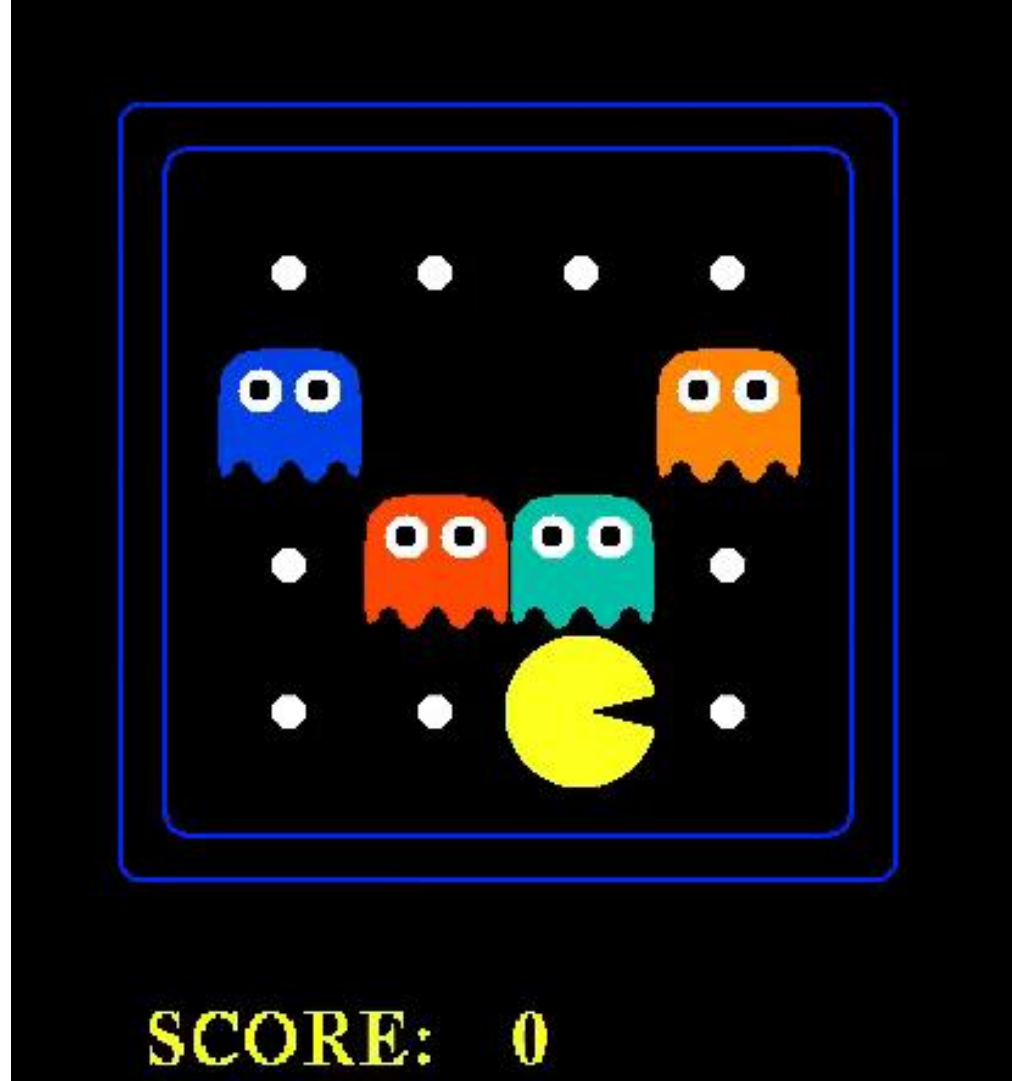
## Logical Agent Algorithms

- Vocab
- PL\_TRUE
- Entailment
  - Model checking: Truth table entailment
  - Theorem proving:
    - (Forward chaining), resolution
- Satisfiability: DPLL
- Planning with logic

# Planning as Satisfiability

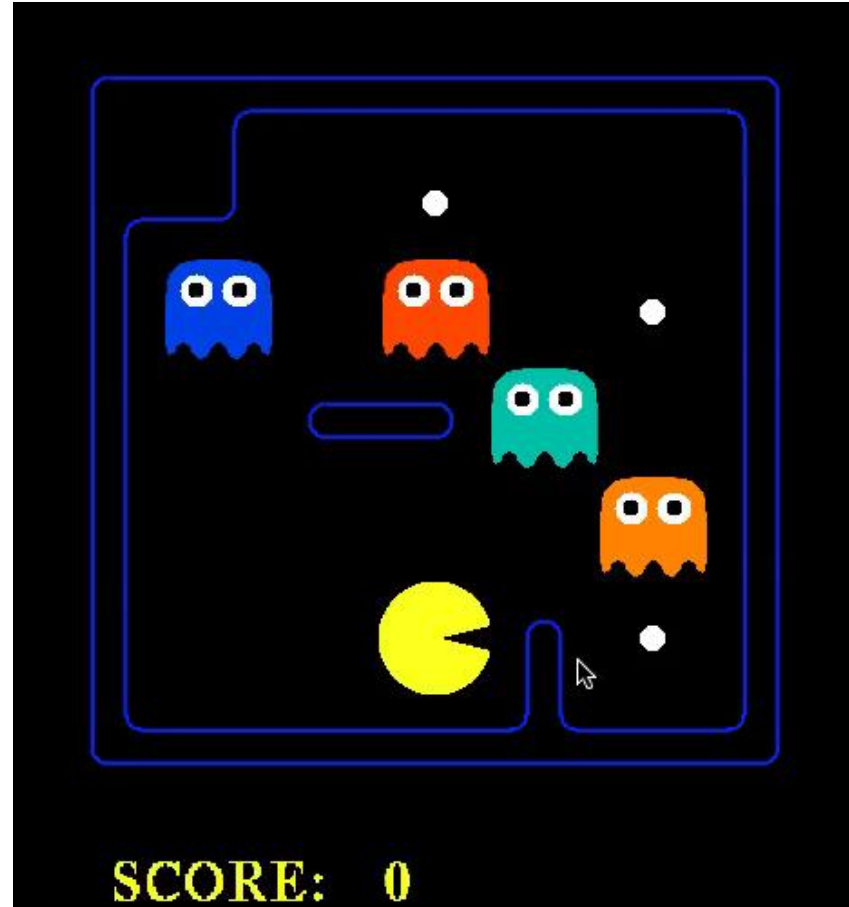
Given a hyper-efficient SAT solver, can we use it to make plans?

Yes, for fully observable, deterministic case: planning problem is solvable iff there is some satisfying assignment for actions etc.









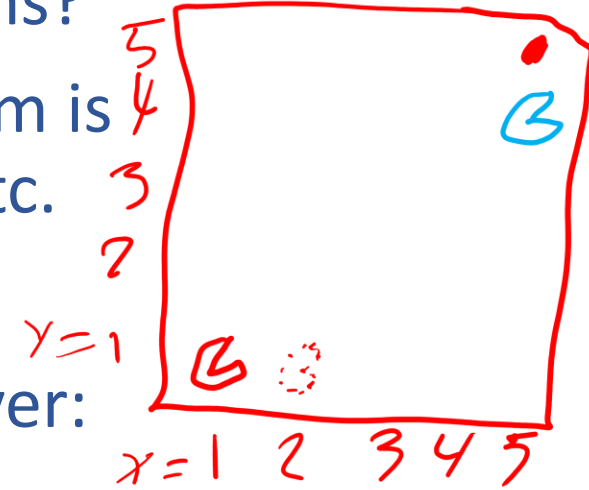
# Planning as Satisfiability

Given a hyper-efficient SAT solver, can we use it to make plans?

Yes, for fully observable, deterministic case: planning problem is solvable iff there is some satisfying assignment for actions etc.

For  $T = 1$  to infinity, set up the KB as follows and run SAT solver:

- Initial state, domain constraints
- Transition model sentences up to time  $T$
- Goal is true at time  $T$
- **Precondition axioms:**  $At_{1,1,0} \wedge N_0 \Rightarrow \neg Wall_{1,2}$  etc.
- **Action exclusion axioms:**  $\neg(N_0 \wedge W_0) \wedge \neg(N_0 \wedge S_0) \wedge \dots$  etc.



$$\left. \begin{array}{l} P_{1,1,0} \\ P_{1,1,0} \wedge East_0 \Rightarrow P_{2,1,1} \\ \dots \end{array} \right\}$$

$$P_{5,4,0} = T \quad P_{1,1,0} = T \quad North_0 = T \quad P_{5,4,0} \wedge North_0 \Rightarrow P_{5,5,1}$$

# Initial State

The agent may know its initial location:

- $At_{1,1_0}$

Or, it may not:

- $At_{1,1_0} \vee At_{1,2_0} \vee At_{1,3_0} \vee \dots \vee At_{3,3_0}$

We also need a *domain constraint* – cannot be in two places at once!

- $\neg(AT_{1,1_0} \wedge At_{1,2_0}) \wedge \neg(AT_{1,1_0} \wedge At_{1,3_0}) \wedge \dots$
- $\neg(AT_{1,1_1} \wedge At_{1,2_1}) \wedge \neg(AT_{1,1_1} \wedge At_{1,3_1}) \wedge \dots$
- ...

# Fluents and Effect Axioms

A *fluent* is a state variable that changes over time

How does each *state variable* or *fluent* at each time gets its value?

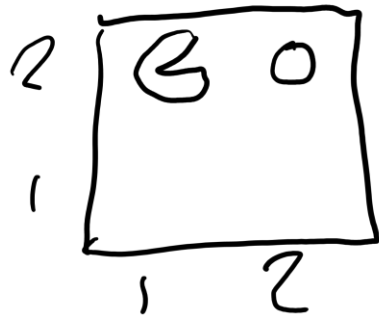
Fluents for PL Pacman are  $\text{Pacman}_{x,y,t}$ , e.g.,  $\text{Pacman}_{3,3,17}$  KB

Model

$$P_{1,2,0} = T$$

$$E_0 = F$$

$$P_{2,2,1} = T$$



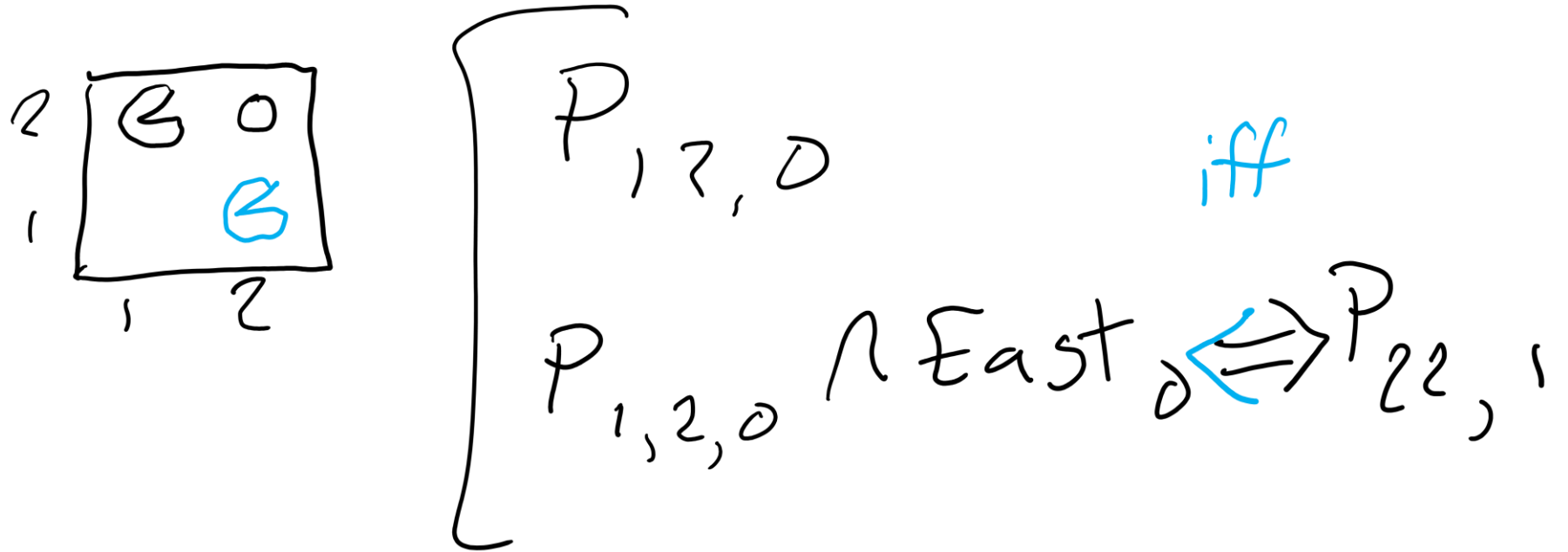
$$\left[ \begin{array}{l} P_{1,2,0} \\ P_{1,2,0} \wedge \text{East}_0 \Rightarrow P_{2,2,1} \\ T \wedge F \Rightarrow T \end{array} \right]$$

# Fluents and Effect Axioms

A *fluent* is a state variable that changes over time

How does each *state variable* or *fluent* at each time gets its value?

Fluents for PL Pacman are  $\text{Pacman}_{x,y,t}$ , e.g.,  $\text{Pacman}_{3,3,17}$



# Fluents and Successor-state Axioms

A *fluent* is a state variable that changes over time

How does each *state variable* or *fluent* at each time gets its value?

Fluents for PL Pacman are  $\text{Pacman}_{x,y,t}$ , e.g.,  $\text{Pacman}_{3,3,17}$

A state variable gets its value according to a *successor-state axiom*

$$\begin{aligned} \blacksquare X_t \Leftrightarrow & [X_{t-1} \wedge \neg(\text{some action}_{t-1} \text{ made it false})] \vee \\ & [\neg X_{t-1} \wedge (\text{some action}_{t-1} \text{ made it true})] \end{aligned}$$

$$P_{22,1} \Leftrightarrow [P_{22,0} \wedge \neg(\text{move P})] \vee [\neg P_{22,0} \wedge (\text{move P})]$$

# Fluents and Successor-state Axioms

A *fluent* is a state variable that changes over time

How does each *state variable* or *fluent* at each time gets its value?

Fluents for PL Pacman are  $\text{Pacman}_{x,y,t}$ , e.g.,  $\text{Pacman}_{3,3,17}$

A state variable gets its value according to a *successor-state axiom*

$$\begin{aligned} \blacksquare X_t \Leftrightarrow & [X_{t-1} \wedge \neg(\text{some action}_{t-1} \text{ made it false})] \vee \\ & [\neg X_{t-1} \wedge (\text{some action}_{t-1} \text{ made it true})] \end{aligned}$$

For Pacman location:

$$\begin{aligned} \blacksquare \text{Pacman}_{3,3,17} \Leftrightarrow & [\text{Pacman}_{3,3,16} \wedge \neg((\neg \text{Wall}_{3,4} \wedge \text{N}_{16}) \vee (\neg \text{Wall}_{4,3} \wedge \text{E}_{16}) \vee \dots)] \\ & \vee [\neg \text{Pacman}_{3,3,16} \wedge ((\text{Pacman}_{3,2,16} \wedge \neg \text{Wall}_{3,3} \wedge \text{N}_{16}) \vee \\ & (\text{Pacman}_{2,3,16} \wedge \neg \text{Wall}_{3,3} \wedge \text{N}_{16}) \vee \dots)] \end{aligned}$$

# Fluents and Successor-state Axioms

A *fluent* is a state variable that changes over time

How does each *state variable* or *fluent* at each time gets its value?

Fluents for PL Pacman are  $\text{Pacman}_{x,y,t}$ , e.g.,  $\text{Pacman}_{3,3,17}$

A state variable gets its value according to a *successor-state axiom*

$$\begin{aligned} \blacksquare X_t \Leftrightarrow & [X_{t-1} \wedge \neg(\text{some action}_{t-1} \text{ made it false})] \vee \\ & [\neg X_{t-1} \wedge (\text{some action}_{t-1} \text{ made it true})] \end{aligned}$$

For Pacman location:

$$\begin{aligned} \blacksquare \text{Pacman}_{3,3,17} \Leftrightarrow & [\text{Pacman}_{3,3,16} \wedge \neg((\neg \text{Wall}_{3,4} \wedge \text{N}_{16}) \vee (\neg \text{Wall}_{4,3} \wedge \text{E}_{16}) \vee \dots)] \\ & \vee [\neg \text{Pacman}_{3,3,16} \wedge ((\text{Pacman}_{3,2,16} \wedge \neg \text{Wall}_{3,3} \wedge \text{N}_{16}) \vee \\ & (\text{Pacman}_{2,3,16} \wedge \neg \text{Wall}_{3,3} \wedge \text{N}_{16}) \vee \dots)] \end{aligned}$$