



Artificial Intelligence

Computer Science, CS541 - A

Jonggi Hong

Announcements

- HW #1 is due today (11:59 pm).



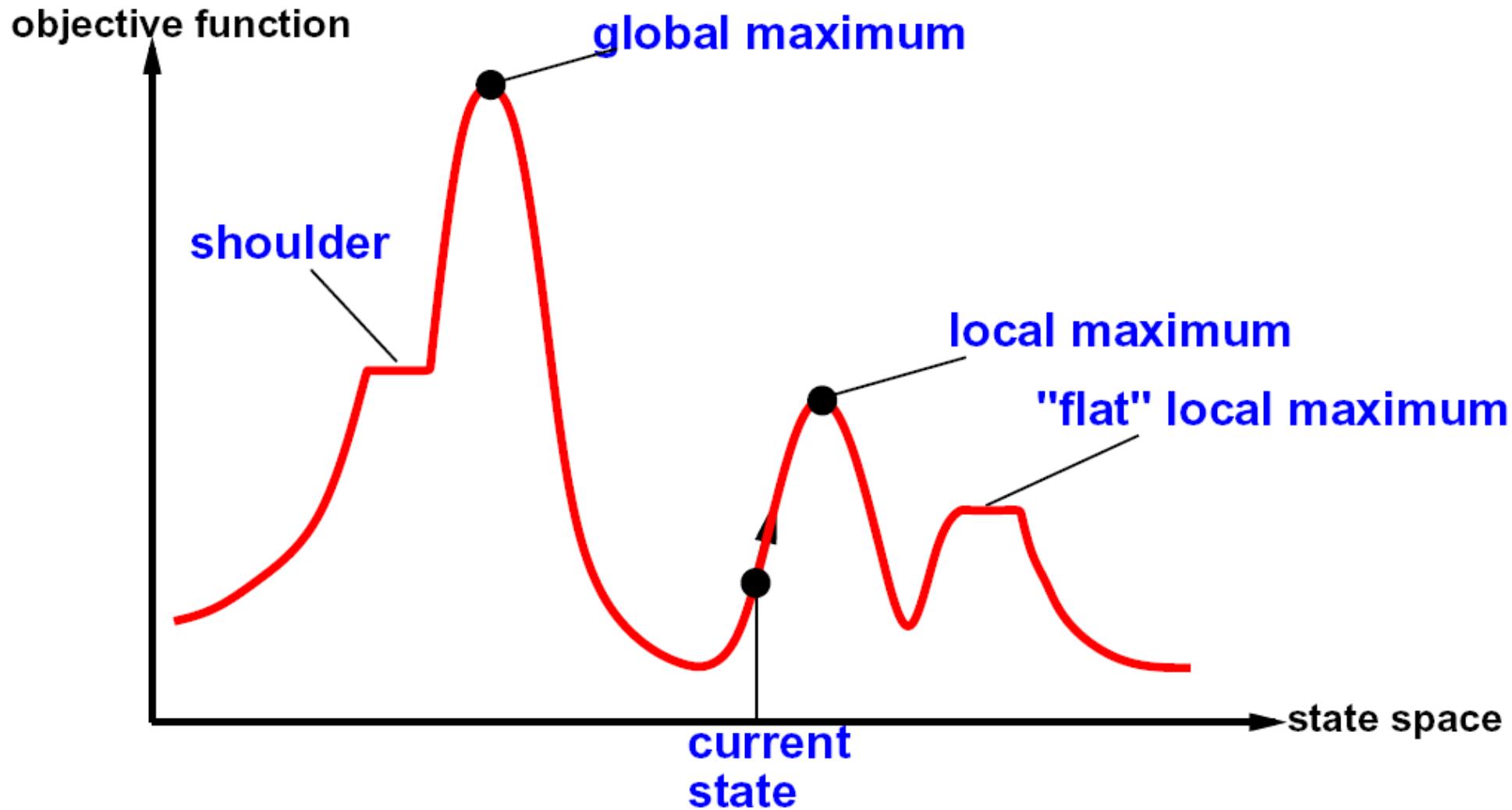
Recap

Local search and adversarial search



Local Search

Hill climbing diagram



Simulated annealing

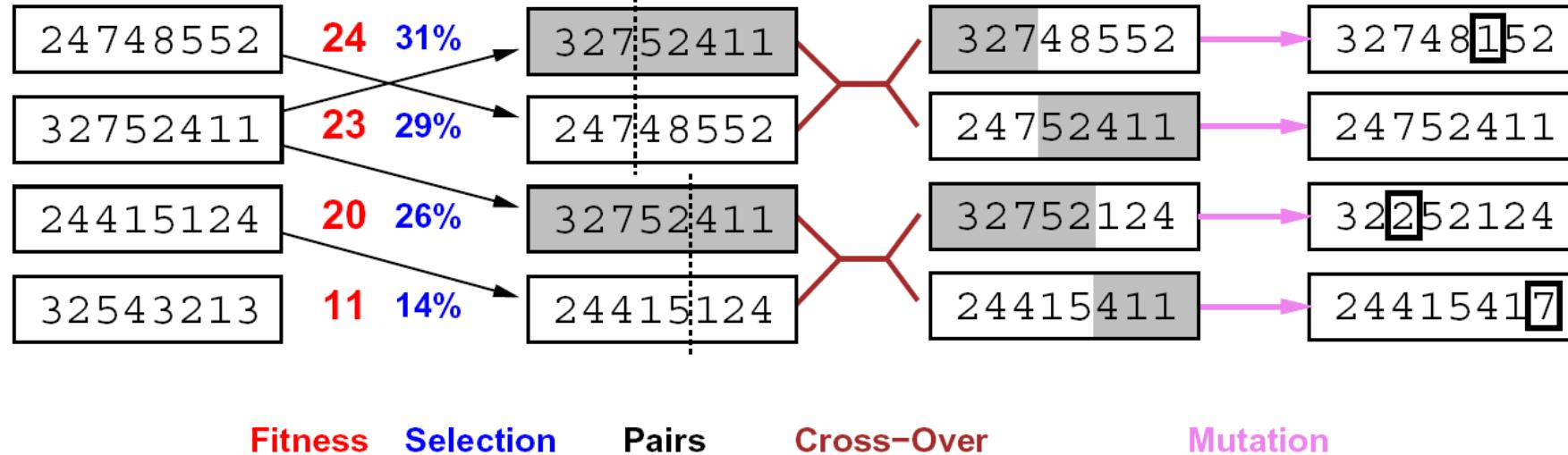
- Idea: Escape local maxima by allowing downhill moves
 - But make them rarer as time goes on

```
function SIMULATED-ANNEALING(problem, schedule) returns a solution state
  inputs: problem, a problem
          schedule, a mapping from time to “temperature”
  local variables: current, a node
                    next, a node
                    T, a “temperature” controlling prob. of downward steps
  current  $\leftarrow$  MAKE-NODE(INITIAL-STATE[problem])
  for t  $\leftarrow$  1 to  $\infty$  do
    T  $\leftarrow$  schedule[t]
    if T = 0 then return current
    next  $\leftarrow$  a randomly selected successor of current
     $\Delta E \leftarrow$  VALUE[next] – VALUE[current]
    if  $\Delta E > 0$  then current  $\leftarrow$  next
    else current  $\leftarrow$  next only with probability  $e^{\Delta E/T}$ 
```

Typical genetic algorithm

1. Initialize: Population $P \leftarrow N$ random individuals
2. Evaluate: For each x in P , compute $\text{fitness}(x)$
3. Loop
 - For $i=1$ to N
 - **Select** 2 parents each with probability proportional to fitness scores
 - **Crossover** the 2 parents to produce a new samples (child)
 - With some small probability **mutate** child
 - Add child to population
 - Until some child is fit enough or you get bored
4. Return best child in the population according to fitness function

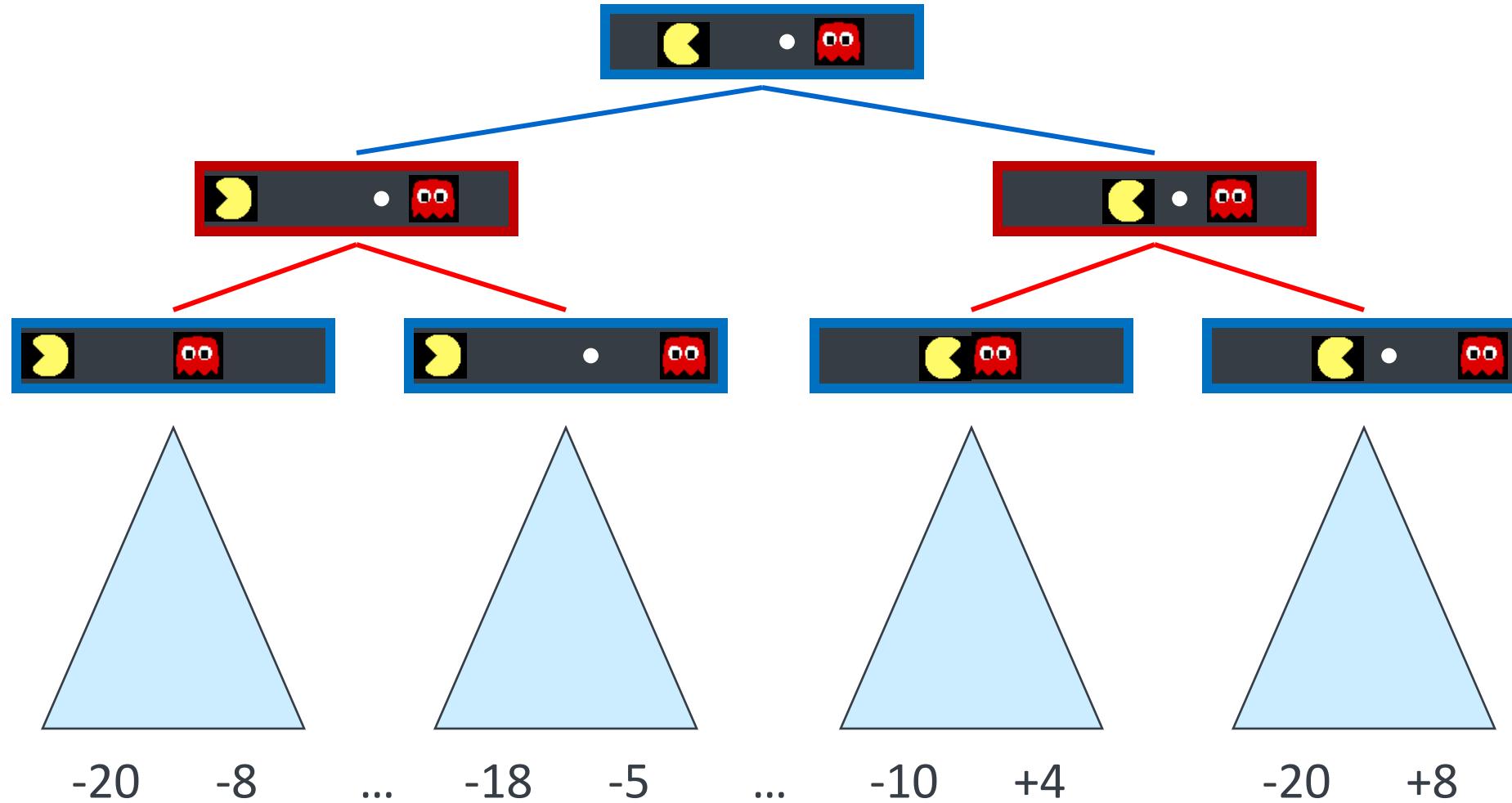
Genetic algorithms



- Genetic algorithms use a natural selection metaphor
 - Keep best N hypotheses at each step (selection) based on a fitness function
 - Also have pairwise crossover operators, with optional mutation to give variety
- Possibly the most misunderstood, misapplied (and even maligned) technique around

Adversarial Search

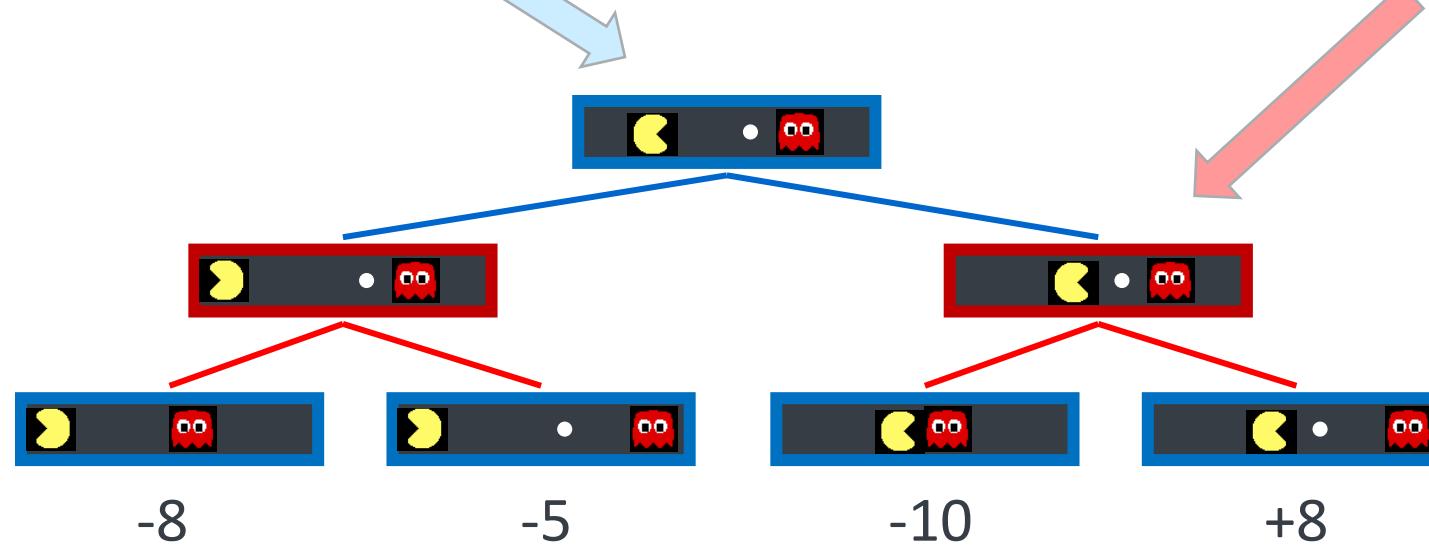
Adversarial game trees



Minimax values

States under agent's control:

$$V(s) = \max_{s' \in \text{successors}(s)} V(s')$$



States under opponent's control:

$$V(s') = \min_{s \in \text{successors}(s')} V(s)$$

Terminal States:

$$V(s) = \text{known}$$

Minimax implementation (dispatch)

```
def value(state):
```

 if the state is a terminal state: return the state's utility

 if the next agent is MAX: return max-value(state)

 if the next agent is MIN: return min-value(state)

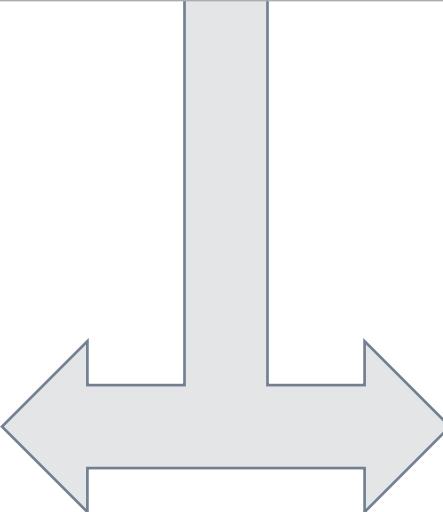
```
def max-value(state):
```

 initialize $v = -\infty$

 for each successor of state:

$v = \max(v, \text{value}(\text{successor}))$

 return v



```
def min-value(state):
```

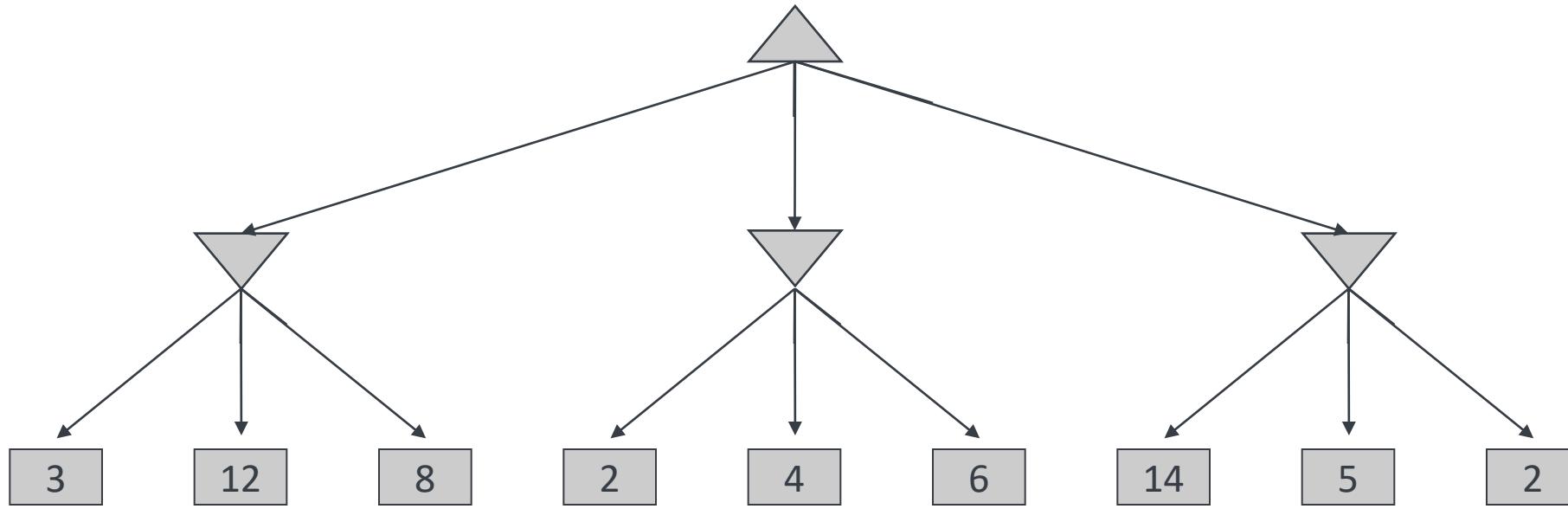
 initialize $v = +\infty$

 for each successor of state:

$v = \min(v, \text{value}(\text{successor}))$

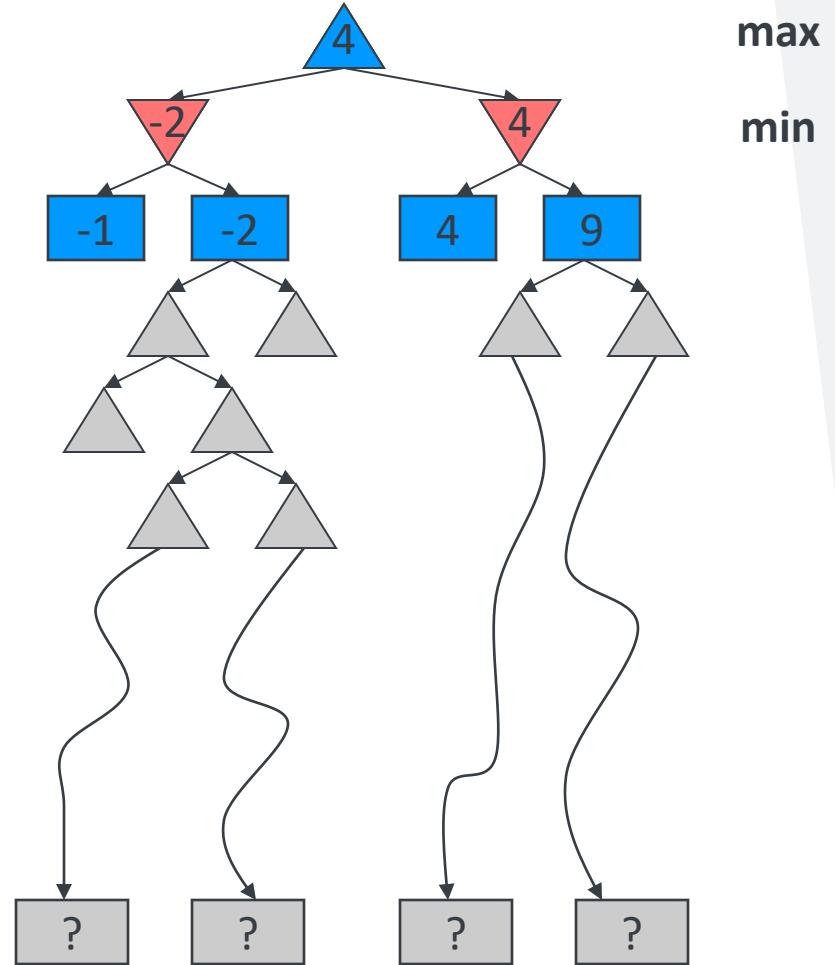
 return v

Minimax example



Resource limits

- Problem: In realistic games, cannot search to leaves!
- Solution: Depth-limited search
 - Instead, search only to a limited depth in the tree
 - Replace terminal utilities with an **evaluation function** for non-terminal positions
- Example:
 - Suppose we have 100 seconds, can explore 10K nodes / sec
 - So can check 1M nodes per move
 - $\alpha\text{-}\beta$ reaches about depth 8
- Guarantee of optimal play is gone
- More plies makes a BIG difference
- Use iterative deepening for an anytime algorithm

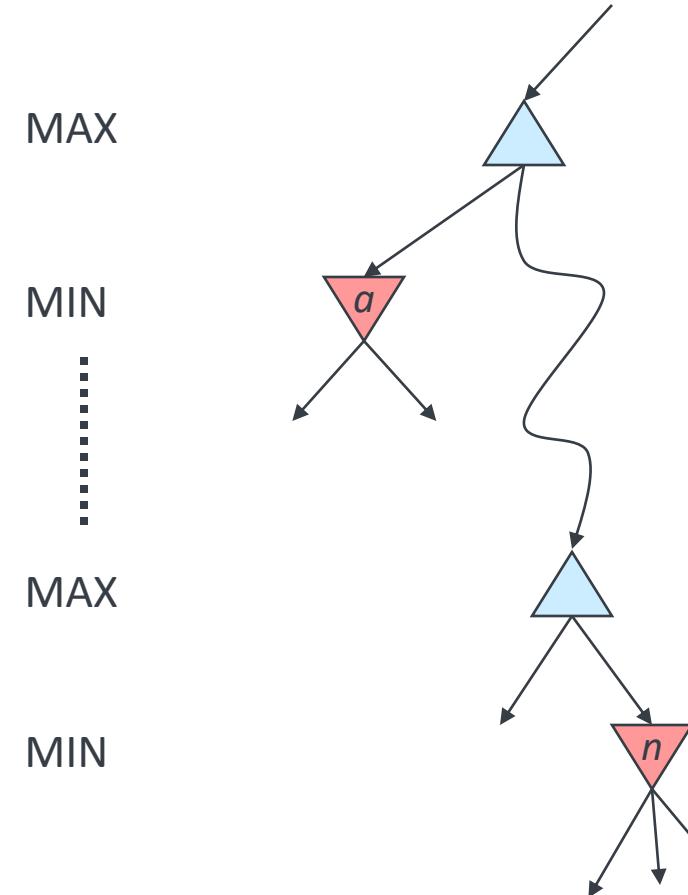


Alpha-beta pruning

- General configuration (MIN version)

- We're computing the MIN-VALUE at some node n
- We're looping over n 's children
- n 's estimate of the children's min is dropping
- Who cares about n 's value? MAX
- Let a be the best value that MAX can get at any choice point along the current path from the root
- If n becomes worse than a , MAX will avoid it, so we can stop considering n 's other children (it's already bad enough that it won't be played)

- MAX version is symmetric



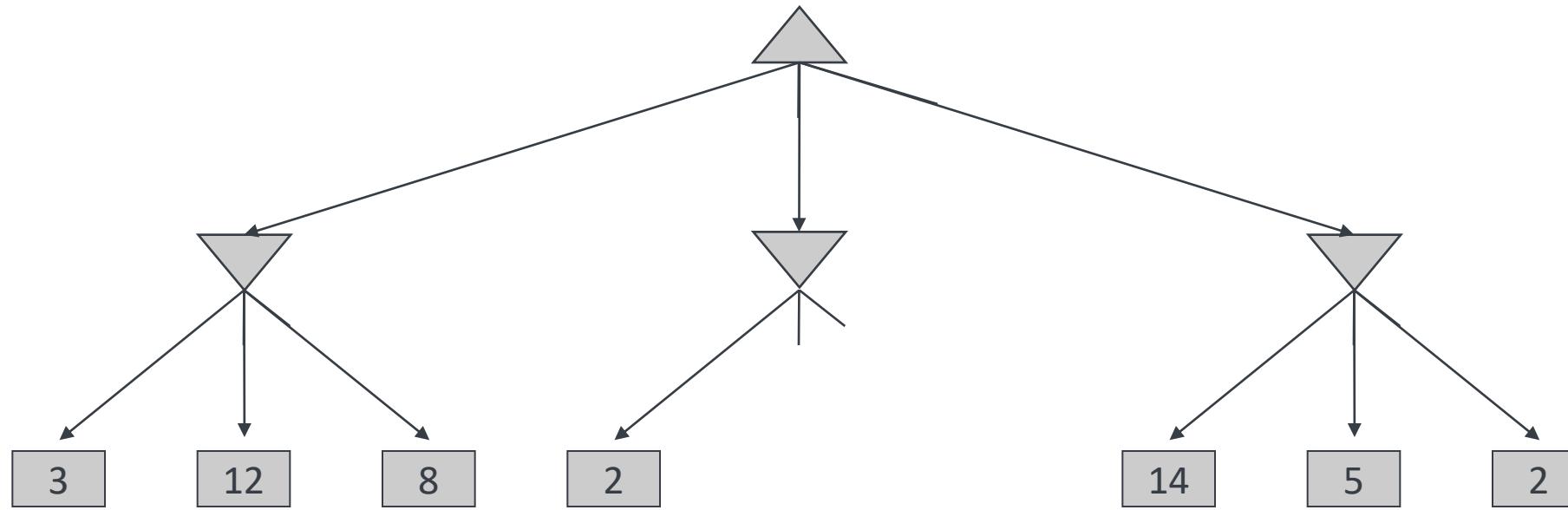
Alpha-beta implementation

α : MAX's best option on path to root
 β : MIN's best option on path to root

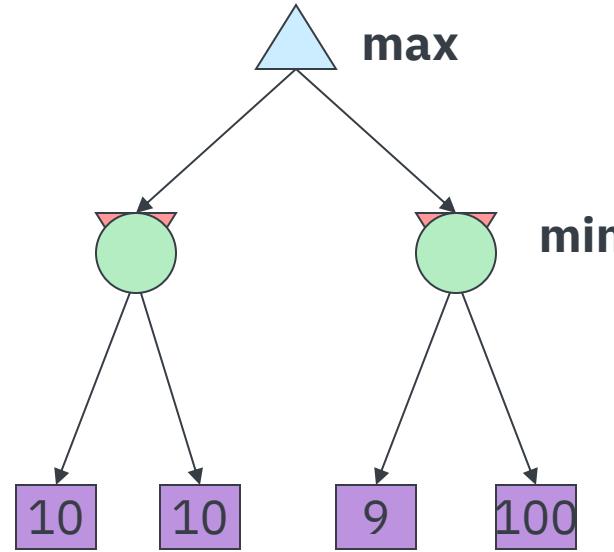
```
def max-value(state,  $\alpha$ ,  $\beta$ ):  
    initialize v = - $\infty$   
    for each successor of state:  
        v = max(v, value(successor,  $\alpha$ ,  $\beta$ ))  
        if v  $\geq \beta$  return v  
         $\alpha$  = max( $\alpha$ , v)  
    return v
```

```
def min-value(state ,  $\alpha$ ,  $\beta$ ):  
    initialize v = + $\infty$   
    for each successor of state:  
        v = min(v, value(successor,  $\alpha$ ,  $\beta$ ))  
        if v  $\leq \alpha$  return v  
         $\beta$  = min( $\beta$ , v)  
    return v
```

Minimax pruning



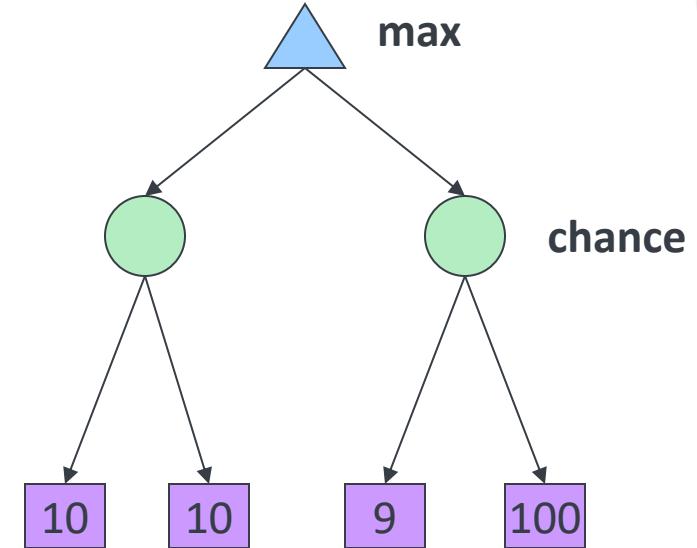
Worst-case vs. Average case



Idea: Uncertain outcomes controlled by chance, not an adversary!

Expectimax search

- Why wouldn't we know what the result of an action will be?
 - Explicit randomness: rolling dice
 - Unpredictable opponents: the ghosts respond randomly
 - Actions can fail: when moving a robot, wheels might slip
- Values should now reflect average-case (expectimax) outcomes, not worst-case (minimax) outcomes
- **Expectimax search:** compute the average score under optimal play
 - Max nodes as in minimax search
 - Chance nodes are like min nodes but the outcome is uncertain
 - Calculate their **expected utilities**
 - *I.e.* take weighted average (expectation) of children



Expectimax pseudocode

```
def value(state):
```

 if the state is a terminal state: return the state's utility

 if the next agent is MAX: return max-value(state)

 if the next agent is EXP: return exp-value(state)

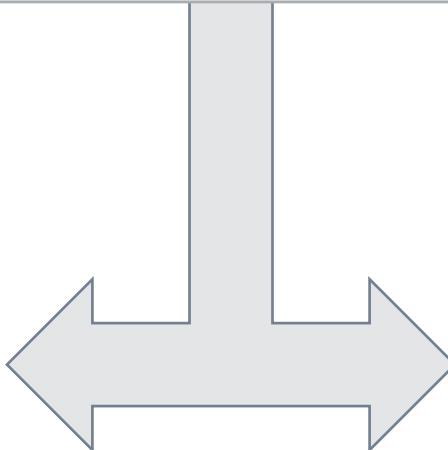
```
def max-value(state):
```

 initialize $v = -\infty$

 for each successor of state:

$v = \max(v, \text{value}(\text{successor}))$

 return v



```
def exp-value(state):
```

 initialize $v = 0$

 for each successor of state:

$p = \text{probability}(\text{successor})$

$v += p * \text{value}(\text{successor})$

 return v



Adversarial Search (cont.)

Chapter 5



Preferences

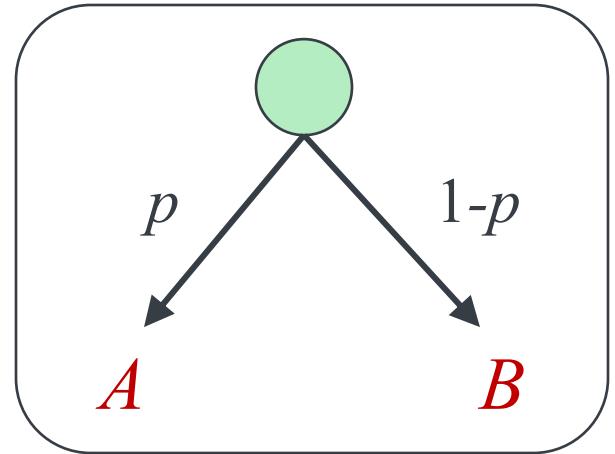
- An agent must have preferences among:
 - Prizes: A , B , etc.
 - Lotteries: situations with uncertain prizes

$$L = [p, A; (1 - p), B]$$

A Prize



A Lottery



- Notation:
 - Preference: $A \succ B$
 - Indifference: $A \sim B$

Rationality

Rational preferences

- We want some constraints on preferences before we call them rational, such as:

Axiom of Transitivity: $(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$

- For example: an agent with **intransitive preferences** can be induced to give away all of its money
 - If $B > C$, then an agent with C would pay (say) 1 cent to get B
 - If $A > B$, then an agent with B would pay (say) 1 cent to get A
 - If $C > A$, then an agent with A would pay (say) 1 cent to get C

Rational preferences

The axioms of rationality

Orderability

$$(A \succ B) \vee (B \succ A) \vee (A \sim B)$$

Transitivity

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

Continuity

$$A \succ B \succ C \Rightarrow \exists p [p, A; 1 - p, C] \sim B$$

Substitutability

$$A \sim B \Rightarrow [p, A; 1 - p, C] \sim [p, B; 1 - p, C]$$

Monotonicity

$$A \succ B \Rightarrow$$

$$(p \geq q \Leftrightarrow [p, A; 1 - p, B] \succeq [q, A; 1 - q, B])$$

Theorem: Rational preferences imply behavior describable as maximization of expected utility.

MEU principle

- Theorem [Ramsey, 1931; von Neumann & Morgenstern, 1944]

- Given any preferences satisfying these constraints, there exists a real-valued function U such that:

$$U(A) \geq U(B) \Leftrightarrow A \succeq B$$

$$U([p_1, S_1; \dots; p_n, S_n]) = \sum_i p_i U(S_i)$$

- *I.e.* values assigned by U preserve preferences of both prizes and lotteries!

- Maximum expected utility (MEU) principle:

- Choose the action that maximizes expected utility
 - Note: an agent can be entirely rational (consistent with MEU) without ever representing or manipulating utilities and probabilities
 - E.g., a lookup table for perfect tic-tac-toe, a reflex vacuum cleaner

Human utilities

Utility scales

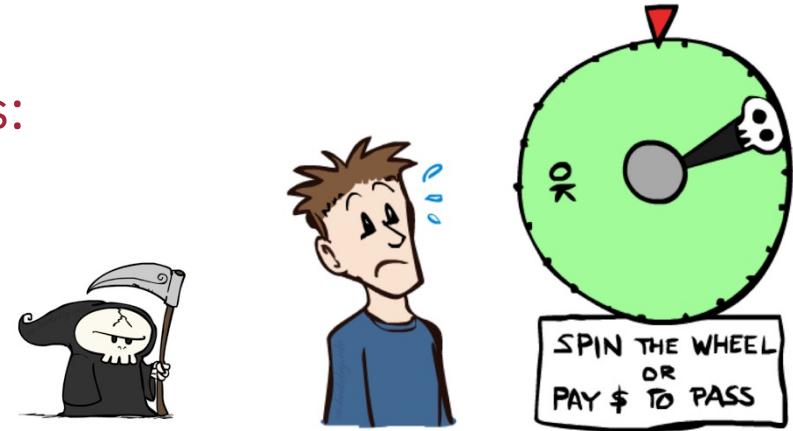
- **Normalized utilities:** $u_+ = 1.0, u_- = 0.0$
- **Micromorts:** one-millionth chance of death, useful for paying to reduce product risks, etc.
- **QALYs:** quality-adjusted life years, useful for medical decisions involving substantial risk
- Note: behavior is invariant under positive linear transformation

$$U'(x) = k_1 U(x) + k_2 \quad \text{where } k_1 > 0$$

- With deterministic prizes only (no lottery choices), only ordinal utility can be determined, *i.e.*, total order on prizes

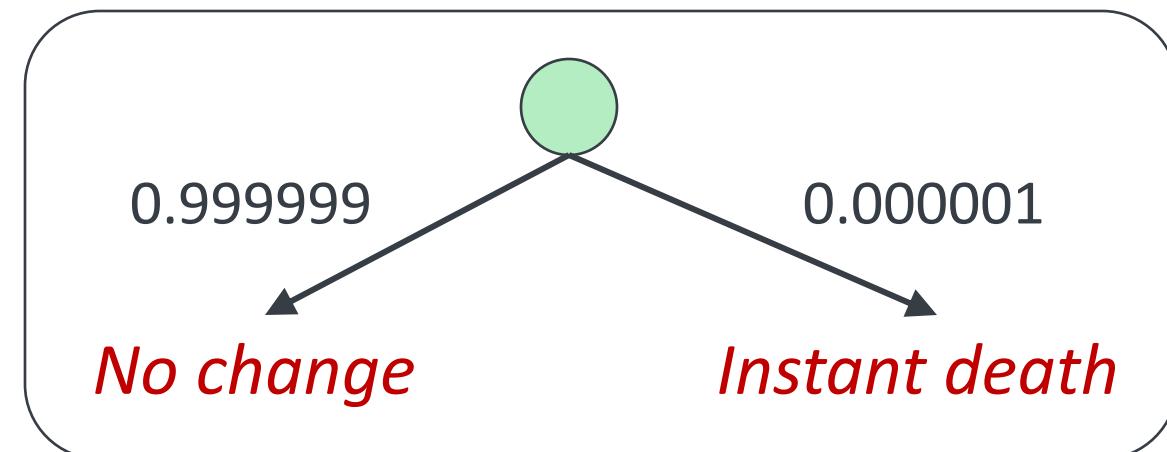
Human utilities

- Utilities map states to real numbers. Which numbers?
- Standard approach to assessment (elicitation) of human utilities:
 - Compare a prize A to a **standard lottery** L_p between
 - “best possible prize” u_+ with probability p
 - “worst possible catastrophe” u_- with probability $1-p$
 - Adjust lottery probability p until indifference: $A \sim L_p$
 - Resulting p is a utility in $[0,1]$



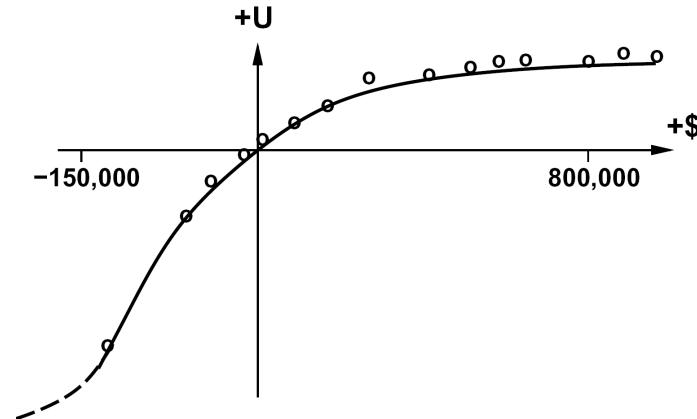
Pay \$30

~



Money

- Money does not behave as a utility function, but we can talk about the utility of having money (or being in debt)
- Given a lottery $L = [p, \$X; (1-p), \$Y]$
 - The **expected monetary value** $EMV(L)$ is $p*X + (1-p)*Y$
 - $U(L) = p*U(\$X) + (1-p)*U(\$Y)$
 - Typically, $U(L) < U(EMV(L))$
 - In this sense, people are **risk-averse**
 - When deep in debt, people are **risk-prone**



Example: Insurance

- Consider the lottery [0.5, \$1000; 0.5, \$0]

- What is its **expected monetary value**? (\$500)
- What is its **certainty equivalent**?
 - Monetary value acceptable in lieu of lottery
 - \$400 for most people
- Difference of \$100 is the **insurance premium**
 - There's an insurance industry because people will pay to reduce their risk
 - If everyone were risk-neutral, no insurance needed!
- It's win-win: you'd rather have the \$400 and the insurance company would rather have the lottery (their utility curve is flat and they have many lotteries)

Which one do you prefer? - 1

A: [0.8, \$4k;
0.2, \$0]

B: [1.0, \$3k;
0.0, \$0]

Which one do you prefer? -2

C: [0.2, \$4k;
0.8, \$0]

D: [0.25, \$3k;
0.75, \$0]

Example: human rationality?

- Famous example of Allais (1953)

- A: [0.8, \$4k; 0.2, \$0]
 - B: [1.0, \$3k; 0.0, \$0]
 - C: [0.2, \$4k; 0.8, \$0]
 - D: [0.25, \$3k; 0.75, \$0]

- Most people prefer B > A, C > D

- But if $U(\$0) = 0$, then

- $B > A \Rightarrow U(\$3k) > 0.8 U(\$4k)$
 - $C > D \Rightarrow 0.8 U(\$4k) > U(\$3k)$



Logic

Chapter 5



Language

- Natural language (informal)
 - English: Two divides even numbers.
 - German: Zwei dividieren geraden zahlen.
- Programming language (formal)
 - Python: `def even(x): return x % 2 == 0`
 - C++: `bool even(int x) { return x % 2 == 0; }`
- Logical languages (formal)
 - First-order logic: $\forall x. \text{Even}(x) \rightarrow \text{Divides}(x, 2)$

Two Goals of a Logic Language

- Represent knowledge about the world
- Reason with that knowledge

What is Logic?

- A formal language
 - **Syntax**: what expressions are legal
 - **Semantics**: what legal expressions mean
 - **Proof system**: a way of manipulating syntactic expressions to get other syntactic expressions (which will tell us something new)
- Why proofs?
 - Two kinds of inferences an agent might want to make:
 - Multiple percepts → conclusions about the world
 - Current state & operator → properties of the next state

What is Logic?

- Propositional logic
- First-order logic

Propositional Logic

Propositional Logic Syntax

- Syntax: what you are allowed to write
 - for (thing t=fizz; t==fuzz; t++) { ... }
 - “Colorless green ideas sleep furiously.”
- Sentences (wffs: well formed formulas)
 - true and false are sentences.
 - Propositional variables are sentences: P, Q, R, Z
 - If φ and ψ are sentences, then so are
 - (φ) , $\neg \varphi$, $\varphi \vee \psi$, $\varphi \wedge \psi$, $\varphi \rightarrow \psi$, $\varphi \leftrightarrow \psi$
 - Nothing else is a sentence.

Precedence

highest \neg

\wedge

$A \vee B \wedge C$

$A \vee (B \wedge C)$

\vee

$A \wedge B \rightarrow C \vee D$

$(A \wedge B) \rightarrow (C \vee D)$

\rightarrow

$A \rightarrow B \vee C \leftrightarrow D$

$(A \rightarrow (B \vee C)) \leftrightarrow D$

lowest \leftrightarrow

- Precedence rules enable “shorthand” form of sentences, but formally only the fully parenthesized form is legal.
- Syntactically ambiguous forms allowed in shorthand only when semantically equivalent: $A \wedge B \wedge C$ is equivalent to $(A \wedge B) \wedge C$ and $A \wedge (B \wedge C)$.

Semantics

- Meaning of a sentence is truth value {**t**, **f**}
- **Interpretation (model)** is an assignment of truth values to the propositional variables.
 - holds(φ , i) [Sentence φ is **t** in interpretation i]
 - fails(φ , i) [Sentence φ is **f** in interpretation i]

Semantic Rules

- $\text{holds}(\text{true}, i)$ for all i
- $\text{fails}(\text{false}, i)$ for all i
- $\text{holds}(\neg\varphi, i)$ if and only if $\text{fails}(\varphi, i)$ negation
- $\text{holds}(\varphi \wedge \psi, i)$ iff $\text{holds}(\varphi, i)$ and $\text{holds}(\psi, i)$ conjunction
- $\text{holds}(\varphi \vee \psi, i)$ iff $\text{holds}(\varphi, i)$ or $\text{holds}(\psi, i)$ disjunction
- $\text{holds}(P, i)$ iff $i(P) = t$
- $\text{fails}(P, i)$ iff $i(P) = f$

Some Important Shorthand

- $\varphi \rightarrow \psi \equiv \neg \varphi \vee \psi$ (conditional, implication)
 - Implication: antecedent \rightarrow consequent
- $\varphi \leftrightarrow \psi \equiv (\varphi \rightarrow \psi) \wedge (\psi \rightarrow \varphi)$ (biconditional, equivalence)

P	Q	$\neg P$	$P \wedge Q$	$P \vee Q$	$P \rightarrow Q$	$Q \rightarrow P$	$P \leftrightarrow Q$
f	f	t	f	f	t	t	t
f	t	t	f	t	t	f	f
t	f	f	f	t	f	t	f
t	t	f	t	t	t	t	t

Note that implication is not “causality”, if P is f then $P \rightarrow Q$ is t.

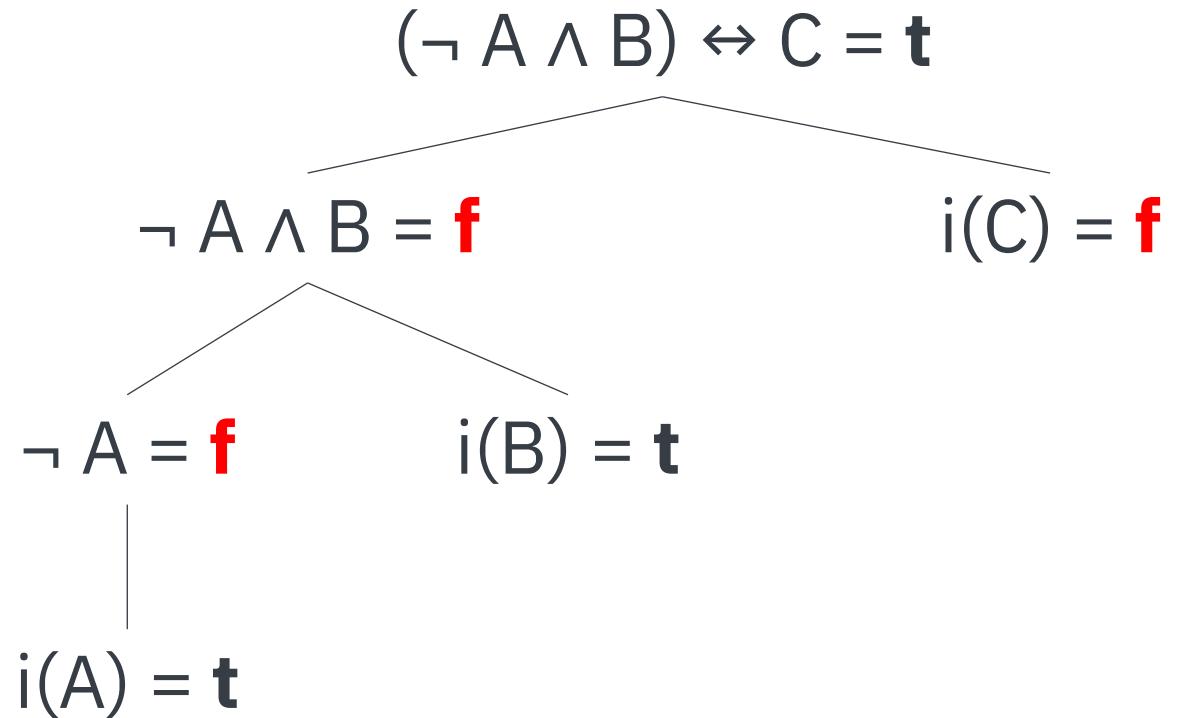
Example: Inverted V-Structure

- Sentence

- $(\neg A \wedge B) \leftrightarrow C$

- Interpretation

- $i(A) = t, i(B) = t, i(C) = f$



Terminology

- A sentence is **valid** iff its truth value is **t** in all interpretations.
 - Valid sentences: true, \neg false, $P \vee \neg P$
- A sentence is **satisfiable** iff its truth value is **t** in at least one interpretation.
 - Satisfiable sentences: P , true, $\neg P$
- A sentence is **unsatisfiable** iff its truth value is **f** in all interpretations.
 - Unsatisfiable sentences: $P \wedge \neg P$, false, \neg true

These are finitely decidable in propositional logic.

Examples

Sentence	Valid?	Interpretation that make sentence's truth value = f
$\text{smoke} \rightarrow \text{smoke}$	yes	
$\text{smoke} \vee \neg \text{smoke}$	yes	
$\text{smoke} \rightarrow \text{fire}$	satisfiable, not valid	$\text{smoke} = \mathbf{t}, \text{fire} = \mathbf{f}$
$(s \rightarrow f) \rightarrow (\neg s \rightarrow \neg f)$	satisfiable, not valid	$s = \mathbf{f}, f = \mathbf{t}$ $s \rightarrow f = \mathbf{t}, \neg s \rightarrow \neg f = \mathbf{f}$
$(s \rightarrow f) \rightarrow (\neg f \rightarrow \neg s)$ <i>(contrapositive)</i>	valid	
$b \vee d \vee (b \rightarrow d)$	valid	
$b \vee d \vee \neg b \vee d$	valid	

Satisfiability

- Related to constraint satisfaction
 - Given a sentence S , try to find an interpretation i such that $\text{holds}(S, i)$.
 - Analogous to find an assignment of values to variables such that the constraints hold.
-
- Brute force method
 - enumerate all interpretations and check
 - Better methods
 - heuristic search
 - constraint propagation
 - stochastic search

Satisfiability Problems

- Scheduling nurses to work in a hospital
 - Propositional variables represent the availability of nurses. (For example, Pat is working on Tuesday at 2.)
 - Constraints on the schedule are represented using logical expressions over the variables.
- Finding bugs in software
 - Propositional variables represent the state of the program.
 - Use logic to describe how the program works and to assert there is a bug.
 - If the sentence is satisfiable, you've found a bug!

A Good Lecture?

- Imagine we knew that
 - If today is sunny, then Tomas will be happy. ($S \rightarrow H$)
 - If Tomas is happy, the lecture will be good. ($H \rightarrow G$)
 - Today is sunny. (S)
- Should we conclude that the lecture will be good?

Checking Interpretations

S	H	G
t	t	t
t	t	f
t	f	t
t	f	f
f	t	t
f	t	f
f	f	f
f	f	f

Checking Interpretations

S	H	G	$S \rightarrow H$	$H \rightarrow G$	S
t	t	t	t	t	t
t	t	f	t	f	t
t	f	t	f	t	t
t	f	f	f	t	t
f	t	t	t	t	f
f	t	f	t	f	f
f	f	f	t	t	f
f	f	f	t	t	f

Knowledge base
(KB)

Checking Interpretations

S	H	G	$S \rightarrow H$	$H \rightarrow G$	S	G
t	t	t	t	t	t	t
t	t	f	t	f	t	f
t	f	t	f	t	t	t
t	f	f	f	t	t	f
f	t	t	t	t	f	t
f	t	f	t	f	f	f
f	f	f	t	t	f	t
f	f	f	t	t	f	f

Good lecture!

Adding a Variable

Another variable: Leslie is happy. (L)

L	S	H	G	S → H	H → G	S	G
t	t	t	t	t	t	t	t
t	t	t	f	t	f	t	f
t	t	f	t	f	t	t	t
t	t	f	f	f	t	t	f
t	f	t	t	t	t	f	t
t	f	t	f	t	f	f	f
t	f	f	f	t	t	f	t
t	f	f	f	t	t	f	f
f	t	t	t	t	t	t	t
f	t	t	f	t	f	t	f
...

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about $A \vee C$?

$A \vee C$ is guaranteed to be true

$A \vee C$ is guaranteed to be false

We don't have enough information to say anything definitive about $A \vee C$

Interpretations

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about $A \vee C$?

A	B	C	$A \vee B$	$\neg B \vee C$	$A \vee C$
t	t	t	t	t	t
t	t	f	t	f	t
t	f	t	t	t	t
t	f	f	t	t	t
f	t	t	t	t	t
f	t	f	t	f	f
f	f	t	f	t	t
f	f	f	f	t	f

$A \vee C$ is guaranteed to be true.

Interpretations

If we know that $A \vee B$ and $\neg B \vee C$ are true, what do we know about A ?

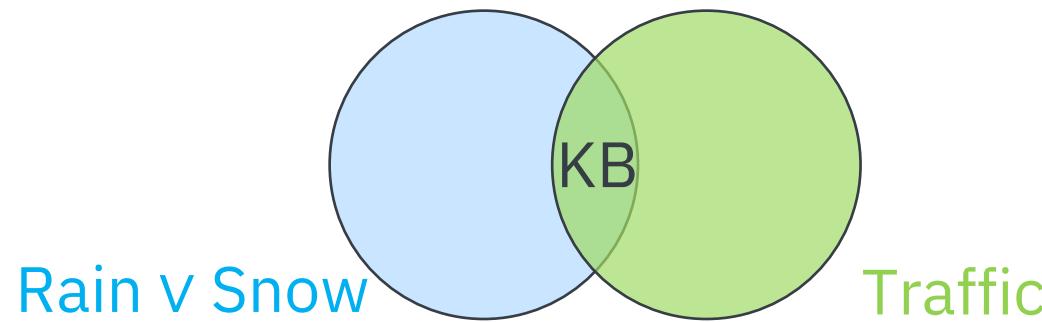
A	B	C	$A \vee B$	$\neg B \vee C$	$A \vee C$
t	t	t	t	t	t
t	t	f	t	f	t
t	f	t	t	t	t
t	f	f	t	t	t
f	t	t	t	t	t
f	t	f	t	f	f
f	f	t	f	t	t
f	f	f	f	t	f

We don't have enough information to say anything definitive about A.

Knowledge Base

- A **knowledge base** (KB) is a set of sentences representing their conjunction/intersection.
- Intuition
 - KB specifies constraints on the world. KB represents the set of interpretations that satisfies those constraints.

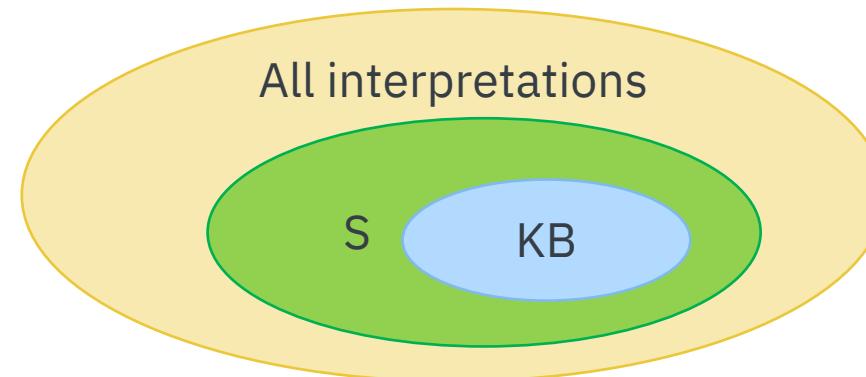
Let $\text{KB} = \{\text{Rain} \vee \text{Snow}, \text{Traffic}\}$



Entailment

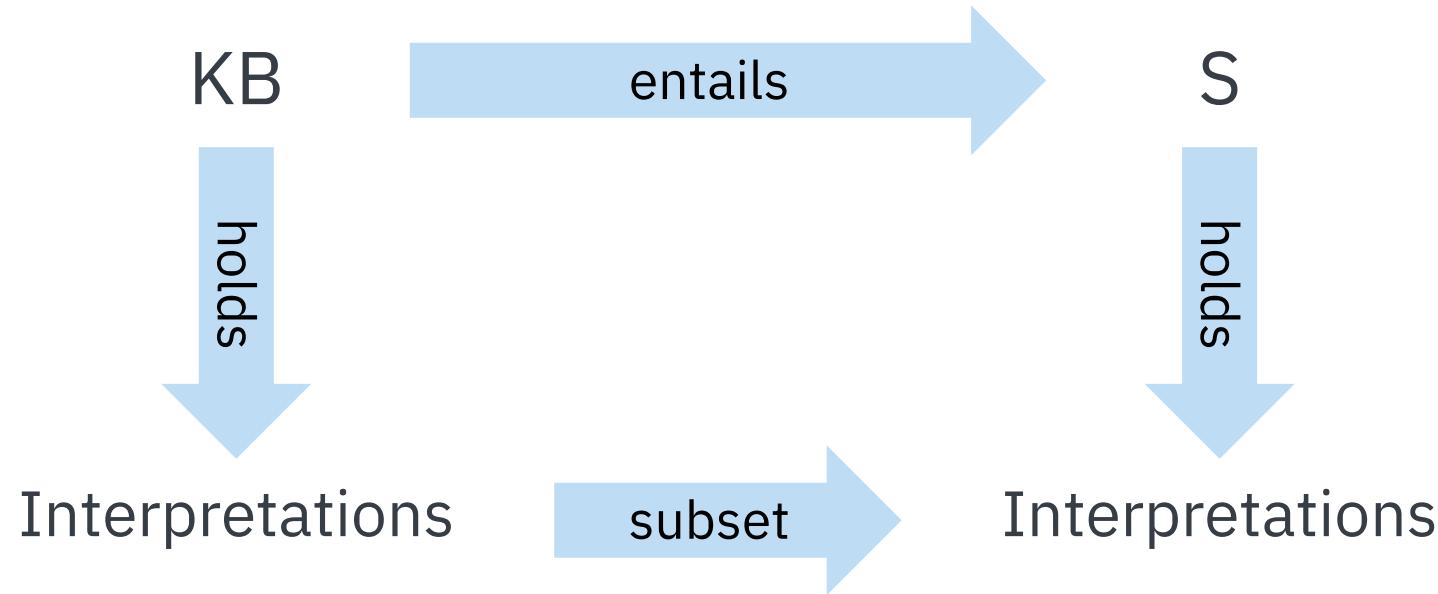
A knowledge base (KB) **entails** a sentence S iff every interpretation that makes KB true also makes S true.

- Intuition
 - The sentence adds no information/constraints (it was already known)
- Entailment
 - KB entails S (written $\text{KB} \models S$)
- Example
 - $\text{Rain} \wedge \text{Snow} \models \text{Snow}$



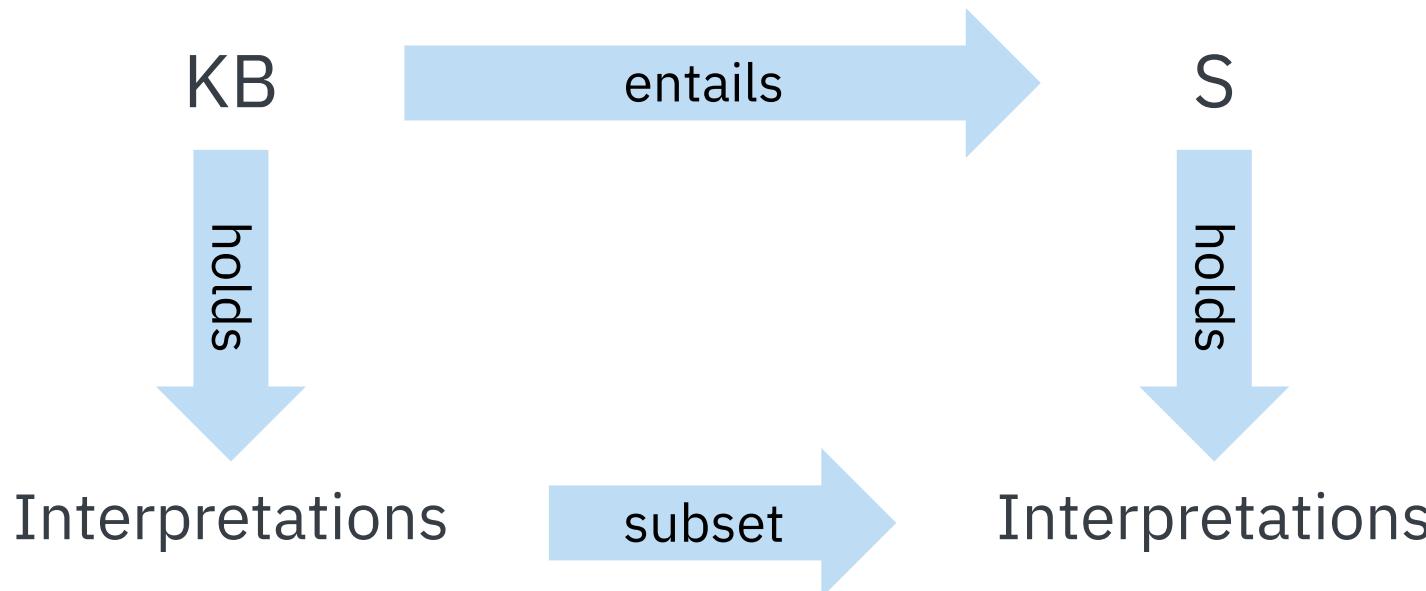
Entailment

A knowledge base (KB) **entails** a sentence S iff every interpretation that makes KB true also makes S true.



Computing Entailment

1. Enumerate all interpretations.
2. Select those in which all elements of KB are true.
3. Check to see if S is true in all of those interpretations.

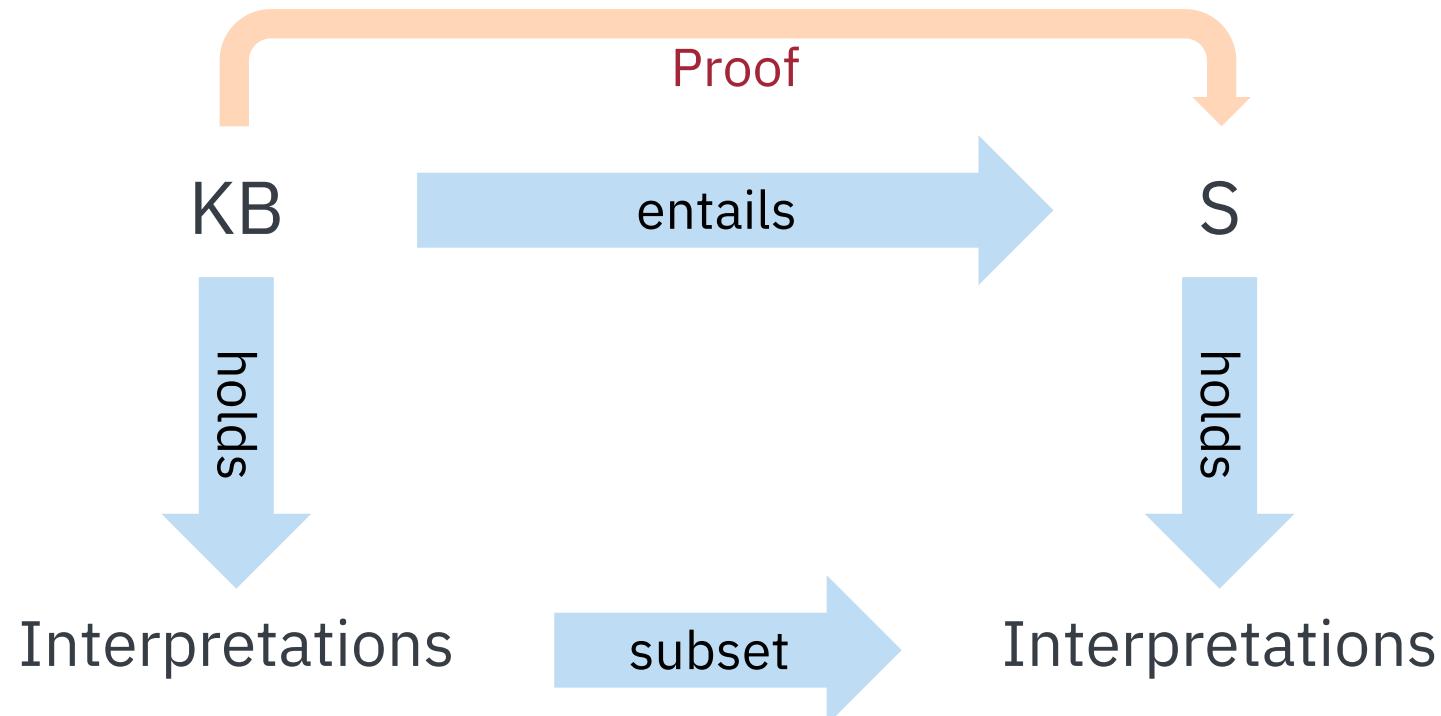


Usually, there are too many interpretations.

$2^{10} = 1024$
interpretations with
10 variables.

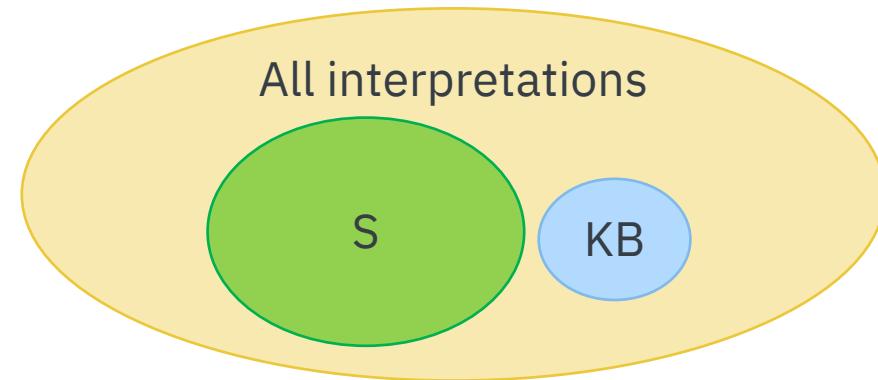
Entailment and Proof

A **proof** is a way to test whether a KB entails a sentence, without enumerating all possible interpretations.

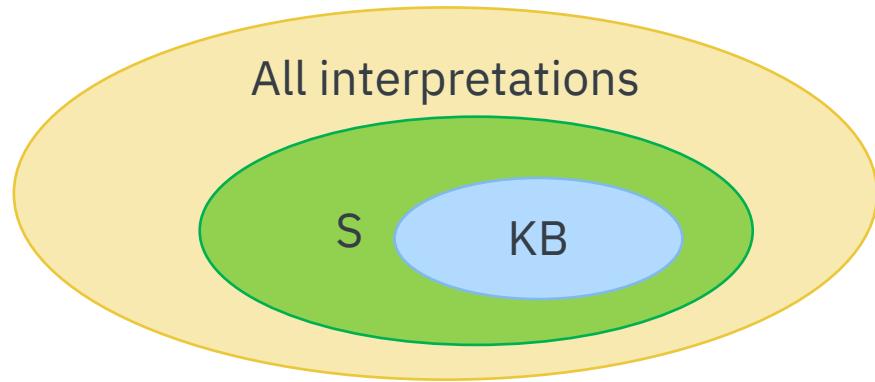


Contradiction

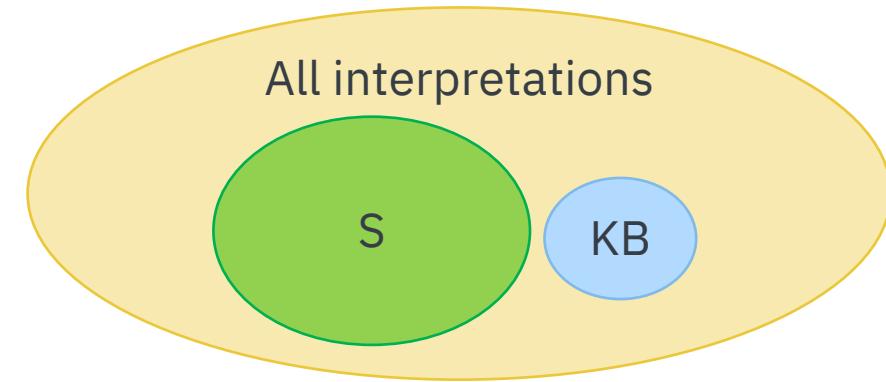
- Intuition
 - The sentence contradicts what we know (captured in KB)
- Contradiction
 - KB contradicts f iff $M(\text{KB}) \cap M(S) = \emptyset$
- Example
 - Rain \wedge Snow contradicts \neg Snow



Entailment and Contradiction



Entailment



Contradiction

KB contradicts S iff KB entails $\neg S$.

Proof

- Proof is a sequence of sentences.
- First ones are premises (KB).
- Then, you can write down on the next line the result of applying an inference rule to previous lines.
- When S is on a line, you have proved S from KB.

Proof

- Example of making an inference
 - It is raining. (Rain)
 - If it is raining, then it is wet. ($\text{Rain} \rightarrow \text{Wet}$)
 - Therefore, it is wet. (Wet)

$$\begin{array}{c} \text{Rain} \\ \hline \text{Rain} \rightarrow \text{Wet} \\ \hline \text{Wet} \end{array}$$

(Modus ponens rule)

Entailment and Proof

■ Semantics

- Interpretation defines entailed/true sentences
- $\text{KB} \vDash S$ (entailment)

■ Syntax

- Inference rules derive sentences
- $\text{KB} \vdash S$ (proof)

Proof

- If inference rules are **sound**, then any S you can prove from KB is entailed by KB .
 - $\text{KB} \vdash S \rightarrow \text{KB} \models S$
- If inference rules are **complete**, then any S that is entailed by KB can be proved from KB .
 - $\text{KB} \models S \rightarrow \text{KB} \vdash S$

Soundness and Completeness

- Knowledge base

- Rain
 - Rain \rightarrow Wet

Is the inference rule sound?

yes

- Sentence

- Wet

Is the inference rule complete?

yes

- Inference rule

- $$\frac{\alpha \rightarrow \beta}{\frac{\alpha}{\beta}}$$

Soundness and Completeness

- Knowledge base

- Rain
- Rain \vee Snow \rightarrow Wet

Is the inference rule sound?

yes

- Sentence

- Wet

Is the inference rule complete?

no

- Inference rule

- $$\frac{\alpha \rightarrow \beta}{\frac{\alpha}{\beta}} \quad \}$$

- Semantically: $\text{KB} \models S$
- Syntactically: $\text{KB} \not\models S$

Fixing Completeness

- Option 1

- Restrict the allowed set of sentences.
- Propositional logic → propositional logic with only Horn clauses.

- Option 2

- Use more powerful inference rules



STEVENS
INSTITUTE OF TECHNOLOGY
1870

15 min. break



Natural Deduction

Natural Deduction

- Example of making an inference

- It is raining. (Rain)
- If it is raining, then it is wet. ($\text{Rain} \rightarrow \text{Wet}$)
- Therefore, it is wet. (Wet)

$$\frac{\begin{array}{c} \text{Rain} \\[1ex] \text{Rain} \rightarrow \text{Wet} \end{array}}{\text{Wet}}$$

(Modus ponens rule)

Natural Deduction

■ Some inference rules

$$\frac{\alpha \rightarrow \beta}{\frac{\alpha}{\beta}}$$

Modus
ponens

$$\frac{\alpha \rightarrow \beta}{\frac{\neg \beta}{\neg \alpha}}$$

Modus
tolens

$$\frac{\alpha}{\frac{\beta}{\alpha \wedge \beta}}$$

And-
introduction

$$\frac{\alpha \wedge \beta}{\alpha}$$

And-
elimination

Natural Deduction Example

Step	Formula	Derivation
1	$P \wedge Q$	Given
2	$P \rightarrow R$	Given
3	$(Q \wedge R) \rightarrow \mathbf{S}$	Given

Natural Deduction Example

Step	Formula	Derivation
1	$P \wedge Q$	Given
2	$P \rightarrow R$	Given
3	$(Q \wedge R) \rightarrow S$	Given
4	P	1 And-Elim

Natural Deduction Example

Step	Formula	Derivation
1	$P \wedge Q$	Given
2	$P \rightarrow R$	Given
3	$(Q \wedge R) \rightarrow S$	Given
4	P	1 And-Elim
5	R	4, 2 Modus Ponens

Natural Deduction Example

Step	Formula	Derivation
1	$P \wedge Q$	Given
2	$P \rightarrow R$	Given
3	$(Q \wedge R) \rightarrow S$	Given
4	P	1 And-Elim
5	R	4, 2 Modus Ponens
6	Q	1 And-Elim

Natural Deduction Example

Step	Formula	Derivation
1	$P \wedge Q$	Given
2	$P \rightarrow R$	Given
3	$(Q \wedge R) \rightarrow S$	Given
4	P	1 And-Elim
5	R	4, 2 Modus Ponens
6	Q	1 And-Elim
7	$Q \wedge R$	5, 6 And-Intro

Natural Deduction Example

Step	Formula	Derivation
1	$P \wedge Q$	Given
2	$P \rightarrow R$	Given
3	$(Q \wedge R) \rightarrow S$	Given
4	P	1 And-Elim
5	R	4, 2 Modus Ponens
6	Q	1 And-Elim
7	$Q \wedge R$	5, 6 And-Intro
8	S	7, 3 Modus Ponens

Natural Deduction Algorithm

Input: set of inference rules Rules

Repeat until no changes to KB:

 Choose a set of formulas f_1, \dots, f_k in KB

 If matching rule $\frac{f_1, \dots, f_k}{S}$ exists:

 Add S to KB

KB derives/proves S ($KB \vdash S$) iff S eventually gets added to KB.

Proof Systems

- There are many natural deduction (ND) systems.
 - They are typically “proof checkers”.
 - The ND systems can be complete and sound with propositional logic.
- Natural deduction uses lots of inference rules which introduces a large branching factor in the search for a proof.
- In general, you need to do “proof by cases” which introduces even more branching.

Prove R	
1	$P \vee Q$
2	$Q \rightarrow R$
3	$P \rightarrow R$

Propositional Resolution

Definite Clauses

- A **definite clause**

$$(p_1 \wedge p_2 \wedge \dots \wedge p_k) \rightarrow q$$

where $p_1 \wedge p_2 \wedge \dots \wedge p_k, q$ are positive propositional symbols

- Intuition

If hold, $p_1 \wedge p_2 \wedge \dots \wedge p_k$ then q holds.

- Examples of definite clauses

- $(\text{Rain} \wedge \text{Snow}) \rightarrow \text{Traffic}$
- Traffic

- Examples of non-definite clauses

- $\neg \text{Traffic}$
- $(\text{Rain} \wedge \text{Snow}) \rightarrow (\text{Traffic} \vee \text{Peaceful})$

Horn Clauses

- A **Horn clause** is either

- a definite clause: $(p_1 \wedge p_2 \wedge \dots \wedge p_k) \rightarrow q$
- a goal clause: $(p_1 \wedge p_2 \wedge \dots \wedge p_k) \rightarrow \mathbf{f}$

- Examples

- $(\text{Rain} \wedge \text{Snow}) \rightarrow \text{Traffic}$ (definite)
- $\text{Traffic} \wedge \text{Accident} \rightarrow \mathbf{f}$ (goal)
 - equivalent: $\neg(\text{Traffic} \wedge \text{Accident})$

Completeness of Modus Ponens

$$\frac{\alpha \rightarrow \beta}{\frac{\alpha}{\beta}}$$

Modus
ponens

■ Claim

Modus ponens is complete with respect to Horn clauses.

- Suppose KB contains only Horn clauses and p is an entailed propositional symbol.
- Then applying modus ponens will derive p .
- $\text{KB} \vDash S \rightarrow \text{KB} \vdash S$

Forward Chaining

$$P \Rightarrow Q$$

$$L \wedge M \Rightarrow P$$

$$B \wedge L \Rightarrow M$$

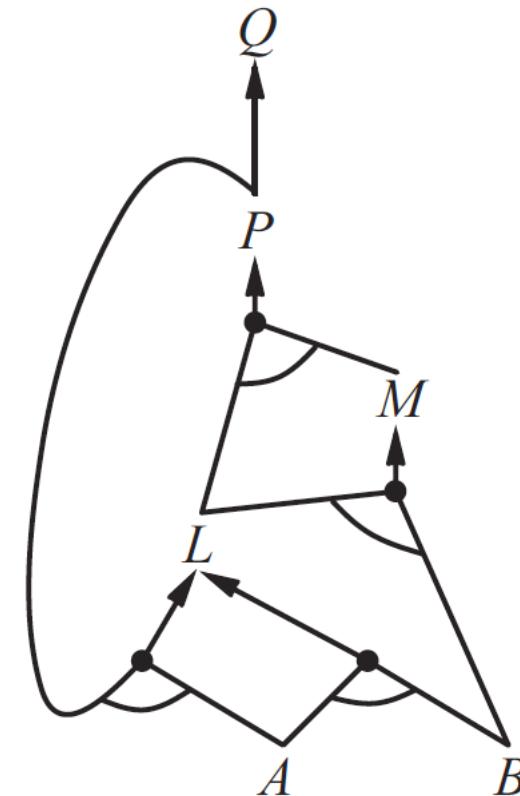
$$A \wedge P \Rightarrow L$$

$$A \wedge B \Rightarrow L$$

$$A$$

$$B$$

Set of Horn clauses



AND-OR graph

- links with arc: conjunction
- links without arc: disjunction

Forward Chaining

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 \leftarrow a table, where *count*[*c*] is the number of symbols in *c*'s premise

inferred \leftarrow a table, where *inferred*[*s*] is initially *false* for all symbols

agenda \leftarrow a queue of symbols, initially symbols known to be true in *KB*

while *agenda* is not empty **do**

p \leftarrow POP(*agenda*)

if *p* = *q* **then return** *true*

if *inferred*[*p*] = *false* **then**

inferred[*p*] \leftarrow *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 *agenda*

return *false*

Horn Clauses and Disjunction

- Horn clauses can be converted to disjunctive clauses

Horn clauses

$$A \rightarrow C$$

$$A \wedge B \rightarrow C$$

Disjunctive clauses

$$\neg A \vee C$$

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

- Literal

- Either p or $\neg p$, where p is a propositional symbol

- Clause

- Disjunction of literals

Propositional Resolution

- Resolution rule

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

- A single inference rule is a sound and complete proof system.
- Requires all sentences to be converted to conjunctive normal form.

Conjunctive Normal Form

- Conjunctive normal form (CNF) formulas

Example: $(A \vee B \vee \neg C) \wedge (B \vee D) \wedge (\neg A) \wedge (B \vee C)$

$(A \vee B \vee \neg C)$ is a **clause**.

A clause is a disjunction of literals.

$A, B, \neg C$ are **literals**.

A literal is a variable or a negation of a variable.

- Each clause is a requirement that must be satisfied and can be satisfied in multiple ways.
- Every sentence in propositional logic can be written in CNF.

Converting to CNF

1. Eliminate arrows (implications) using definitions.
2. Drive in negations using De Morgan's Laws.

$$\begin{aligned}\neg(\varphi \vee \psi) &\equiv \neg \varphi \wedge \neg \psi \\ \neg(\varphi \wedge \psi) &\equiv \neg \varphi \vee \neg \psi\end{aligned}$$

3. Distribute **or** over **and**.

$$A \vee (B \wedge C) \equiv (A \vee B) \wedge (A \vee C)$$

Every sentence can be converted to CNF, but it may grow exponentially in size.

CNF Conversion Example

$$(A \vee B) \rightarrow (C \rightarrow D)$$

1. Eliminate arrows (implications) using definitions.

$$\neg(A \vee B) \vee (\neg C \vee D)$$

2. Drive in negations using De Morgan's Laws.

$$(\neg A \wedge \neg B) \vee (\neg C \vee D)$$

3. Distribute **or** over **and**.

$$(\neg A \vee \neg C \vee D) \wedge (\neg B \vee \neg C \vee D)$$

Propositional Resolution

- Resolution rule

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

- Resolution refutation

- Convert all sentences to CNF.
- Negate the desired conclusion (converted to CNF).
- Apply resolution rule until either
 - Derive false (a contradiction)
 - Can't apply any more

Propositional Resolution

- Resolution refutation is sound and complete.
 - If we derive a contradiction, then the conclusion follows from the axioms.
 - If we can't apply any more, then the conclusion cannot be proved from the axioms.

Propositional Resolution Example

Prove R	
1	$P \vee Q$
2	$Q \rightarrow R$
3	$P \rightarrow R$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion

Propositional Resolution Example

Prove R	
1	$P \vee Q$
2	$Q \rightarrow R$
3	$P \rightarrow R$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion
5	$Q \vee R$	1, 2

Propositional Resolution Example

Prove R	
1	$P \vee Q$
2	$Q \rightarrow R$
3	$P \rightarrow R$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion
5	$Q \vee R$	1, 2
6	$\neg P$	2, 4

Propositional Resolution Example

Prove R	
1	$P \vee Q$
2	$Q \rightarrow R$
3	$P \rightarrow R$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion
5	$Q \vee R$	1, 2
6	$\neg P$	2, 4
7	$\neg Q$	3, 4

Propositional Resolution Example

Prove R	
1	$P \vee Q$
2	$Q \rightarrow R$
3	$P \rightarrow R$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion
5	$Q \vee R$	1, 2
6	$\neg P$	2, 4
7	$\neg Q$	3, 4
8	R	5, 7

Propositional Resolution Example

$$\begin{array}{c} \text{Prove } R \\ \hline \begin{array}{c|c} & P \vee Q \\ \hline 1 & P \vee Q \\ \hline 2 & Q \rightarrow R \\ \hline 3 & P \rightarrow R \\ \hline \end{array} \\ \begin{array}{c} \text{false} \vee R \\ \hline \neg R \vee \text{false} \\ \hline \text{false} \vee \text{false} \end{array} \end{array}$$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion
5	$Q \vee R$	1, 2
6	$\neg P$	2, 4
7	$\neg Q$	3, 4
8	R	5, 7
9	\bullet	4, 8

Propositional Resolution Example

$$\begin{array}{c} \text{Prove } R \\ \hline \begin{array}{c|c} & P \vee Q \\ \hline 1 & P \vee Q \\ \hline 2 & Q \rightarrow R \\ \hline 3 & P \rightarrow R \\ \hline \end{array} \\ \\ \begin{array}{c} \text{false } \vee R \\ \hline \neg R \vee \text{false} \\ \hline \text{false } \vee \text{false} \end{array} \end{array}$$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$\neg P \vee R$	Given
3	$\neg Q \vee R$	Given
4	$\neg R$	Negated conclusion
5	$Q \vee R$	1, 2
6	$\neg P$	2, 4
7	$\neg Q$	3, 4
8	R	5, 7
9	\bullet	4, 8

The Power of False

Prove Z	
1	P
2	$\neg P$

Step	Formula	Derivation
1	P	Given
2	$\neg P$	Given
3	$\neg Z$	Negated conclusion

The Power of False

Prove Z	
1	P
2	¬ P

Step	Formula	Derivation
1	P	Given
2	¬ P	Given
3	¬ Z	Negated conclusion
4	•	1, 2

- Note that $(P \wedge \neg P) \rightarrow Z$ is valid.
- Any conclusion follows from a contradiction – and so strict logic systems are very brittle.

Propositional Resolution Example 2

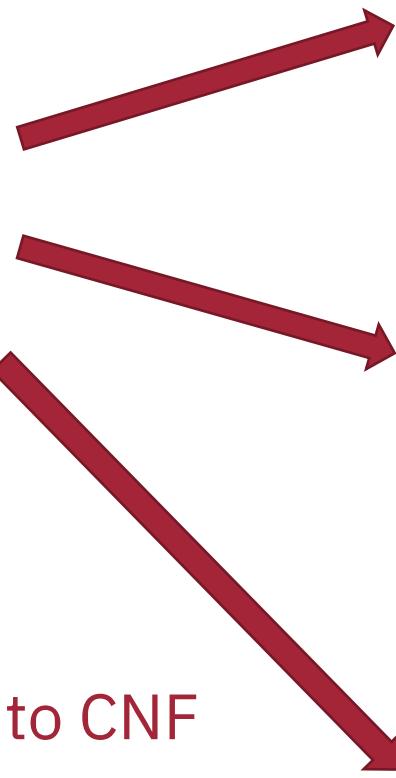
Prove R

1	$(P \rightarrow Q) \rightarrow Q$
2	$(P \rightarrow P) \rightarrow R$
3	$(R \rightarrow S) \rightarrow \neg(S \rightarrow Q)$

Propositional Resolution Example 2

Prove R

1	$(P \rightarrow Q) \rightarrow Q$
2	$(P \rightarrow P) \rightarrow R$
3	$(R \rightarrow S) \rightarrow \neg(S \rightarrow Q)$



$$\begin{aligned} & \neg(\neg P \vee Q) \vee Q \\ & (P \wedge \neg Q) \vee Q \\ & (P \vee Q) \wedge (\neg Q \vee Q) \\ & (P \vee Q) \end{aligned}$$

$$\begin{aligned} & \neg(\neg P \vee P) \vee R \\ & (P \wedge \neg P) \vee R \\ & (P \vee R) \wedge (\neg P \vee R) \end{aligned}$$

$$\begin{aligned} & \neg(\neg R \vee S) \vee \neg(\neg S \vee Q) \\ & (R \wedge \neg S) \vee (S \wedge \neg Q) \\ & (R \vee S) \wedge (\neg S \vee S) \wedge (R \vee \neg Q) \wedge (\neg S \vee \neg Q) \\ & (R \vee S) \wedge (R \vee \neg Q) \wedge (\neg S \vee \neg Q) \end{aligned}$$

Convert to CNF

Propositional Resolution Example 2

Prove R	
1	$(P \rightarrow Q) \rightarrow Q$
2	$(P \rightarrow P) \rightarrow R$
3	$(R \rightarrow S) \rightarrow \neg(S \rightarrow Q)$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$P \vee R$	Given
3	$\neg P \vee R$	Given
4	$R \vee S$	Given
5	$R \vee \neg Q$	Given
6	$\neg S \vee \neg Q$	Given
7	$\neg R$	Negated conclusion

Propositional Resolution Example 2

Prove R	
1	$(P \rightarrow Q) \rightarrow Q$
2	$(P \rightarrow P) \rightarrow R$
3	$(R \rightarrow S) \rightarrow \neg(S \rightarrow Q)$

Step	Formula	Derivation
1	$P \vee Q$	Given
2	$P \vee R$	Given
3	$\neg P \vee R$	Given
4	$R \vee S$	Given
5	$R \vee \neg Q$	Given
6	$\neg S \vee \neg Q$	Given
7	$\neg R$	Negated conclusion
8	S	4, 7
9	$\neg Q$	6, 8
10	P	1, 9
11	R	3, 10
12	\bullet	7, 11

Proof Strategies

- Unit preference

Prefer a resolution step **involving a unit clause** (clause with one literal)

- It produces a shorter clause. Our goal is to produce a zero-length clause!

- Set of support

Choose a resolution **involving the negated goal** or any clause derived **from the negated goal**.

- We are trying to produce a contradiction that follows from the negated goal.
These are relevant clauses.
- If a contradiction exists, it can be found with the set-of-support strategy.

First-Order Logic

Limitations of Propositional Logic

- Propositional logic only deals with fact statements that may or may not be true of the world.
 - Ex) $R = \text{"It is raining"}$
 - You cannot have variables that stands for books or tables.
- Examples
 - When you paint a block with green paint, it becomes green.
 - You would need a statement about every single block. You cannot make a general statement about all blocks.
 - When you sterilize a jar, all the bacteria are dead.
 - In First-Order Logic, we can talk about all bacteria without naming them explicitly.

First-Order Logic (FOL)

- **Variables** refer to things in the world.
- You can **quantify** over the variables.
 - You can make sentences talking about all of them or some of them without having to name them explicitly.

FOL Syntax

- Term
 - Constant symbols
 - Fred, Japan, Bacterium39
 - Variables
 - x, y, a
 - Function symbol applied to one or more terms
 - $f(x)$, $f(f(x))$, mother-of(John)

FOL Syntax

■ Sentence

- A predicate symbol applied to zero or more terms
 - $\text{On}(a, b)$, $\text{Sister}(\text{Jane}, \text{Joan})$, $\text{Sister}(\text{mother-of}(\text{John}), \text{Jane})$
- $t_1 = t_2$
- If v is a variable and Φ is a sentence, then $\forall v.\Phi$ and $\exists v.\Phi$ are sentences.
- Closure under sentential operators: $\vee \wedge \neg \rightarrow \leftrightarrow ()$

FOL Interpretations

■ Interpretation I

- U set of objects
 - called “domain of discourse” or “universe”
- Maps constant symbols to elements of U
- Maps predicate symbols to relations on U
 - An n -ary relation is a set of lists of n objects. A binary relation is a set of pairs.
- Maps function symbols to functions on U
 - A function takes one or more arguments and returns a value. A predicate takes one or more arguments and is either true or false.
 - Assume our functions are *total* (*i.e.*, a function that returns a well-defined value for every element of its domain, with no undefined or “partial” values).

Denotation of Terms

- Terms name objects in U
 - $I(Fred)$
 - If Fred is constant, then given.
 - $I(x)$
 - If x is a variable, then undefined.
 - $I(f(t_1, \dots, t_n))$
 - $I(f)(I(t_1), \dots, I(t_n))$

Holds

- When does a sentence hold in an interpretation?

- P is a relation symbol.
- t_1, \dots, t_n are terms.

$\text{holds}(P(t_1, \dots, t_n), I) \text{ iff } \langle I(t_1), \dots, I(t_n) \rangle \in I(P)$

Brother(Jon, Joe)?

- $I(\text{Jon}) = \text{😊}$ [an element of U]
- $I(\text{Joe}) = \text{😎}$ [an element of U]
- $I(\text{Brother}) = \{\langle \text{😊}, \text{😎} \rangle, \langle \text{😊}, \text{😎} \rangle, \langle \dots, \dots \rangle, \dots\}$

Equality

- When are two terms equal?

$\text{holds}(t_1 = t_2, I)$ iff $I(t_1)$ is the same object as $I(t_2)$.

Jon = Jack?

- $I(\text{Jon}) = \text{😊}$ [an element of U]
- $I(\text{Jack}) = \text{😊}$ [an element of U]
- $\text{holds}(\text{Jon} = \text{Jack}, I)$

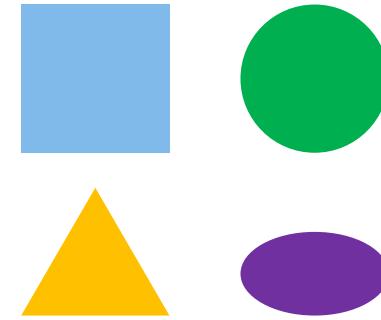
Semantics of Quantifiers

- Extend an interpretation I to bind variable x to element $a \in U$: $I_{x/a}$
 - $\text{holds}(\forall x.\Phi), I$ iff $\text{holds}(\Phi, I_{x/a})$ for all $a \in U$
 - $\text{holds}(\exists x.\Phi), I$ iff $\text{holds}(\Phi, I_{x/a})$ for some $a \in U$
- Quantifier applies to formula to right until an enclosing right parenthesis.

$$(\forall x.P(x) \vee Q(x)) \vee \exists x.R(x) \rightarrow Q(x)$$

FOL Example Domain

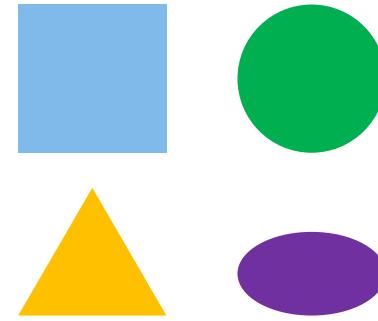
- $U = \{ \text{□}, \text{○}, \text{△}, \text{○} \}$
- Constants
 - Fred
- Preds
 - Above₂, Circle₁, Oval₁, Square₁
- Function
 - hat
- $I(Fred) = \triangle$



The real
world

FOL Example Domain

- $I(Fred) = \triangle$
- $I(Above) = \{<\square, \triangle>, <\bullet, \circ>\}$
- $I(Circle) = \{\bullet\}$
- $I(Oval) = \{\bullet, \circ\}$
- $I(\hat{}) = \{\triangle, \square, \circ, \bullet, \bullet, \bullet\}$
- $I(Square) = \{\square\}$



The real
world

FOL Example Domain

- $I(Fred) = \triangle$ holds(Square(Fred), I)?
- $I(Above) = \{<\square, \triangle>, <\bullet, \circlearrowleft>\}$ holds(Above(Fred, hat(Fred)), I)?
- $I(Circle) = \{<\bullet>\}$ • $I(hat(Fred)) = \square$
- $I(Oval) = \{<\bullet>, <\circlearrowleft>\}$ • holds(Above(\triangle , \square), I)?
- $I(hat) = \{<\triangle, \square>, <\circlearrowleft, \bullet>, <\square, \square>, <\bullet, \bullet>\}$
- $I(Square) = \{<\triangle>\}$

FOL Example Domain

- $I(Fred) = \triangle$
- $I(Above) = \{<\square, \triangle>, <\bullet, \circlearrowleft>\}$
- $I(Circle) = \{<\bullet>\}$
- $I(Oval) = \{<\bullet>, <\circlearrowright>\}$
- $I(\hat{}) = \{<\triangle, \square>, <\circlearrowleft, \bullet>, <\square, \square>, <\bullet, \bullet>\}$
- $I(Square) = \{<\triangle>\}$

$\text{holds}(\exists x. \text{Oval}(x), I)?$

- $\text{holds}(\text{Oval}(x), I_{x/\bullet})?$

$\text{holds}(\forall x. \exists y. \text{Above}(x,y) \vee \text{Above}(y,x), I)?$

- $\text{holds}(\exists y. \text{Above}(x,y) \vee \text{Above}(y,x), I_{x/\triangle})?$
 - $\text{holds}(\text{Above}(x,y) \vee \text{Above}(y,x), I_{x/\triangle, y/\square})?$

- verify for all other values of x.

$\text{holds}(\forall x. \forall y. \text{Above}(x,y) \vee \text{Above}(y,x), I)?$

- $\text{holds}(\text{Above}(x,y) \vee \text{Above}(y,x), I_{x/\square, y/\bullet})?$

Writing FOL

1. Cats are mammals. [Cat¹, Mammal¹]

- $\forall x. \text{Cat}(x) \rightarrow \text{Mammal}(x)$

2. Jane is a tall surveyor. [Tall¹, Surveyor¹, Jane]

- $\text{Tall}(\text{Jane}) \wedge \text{Surveyor}(\text{Jane})$

3. A nephew is a sibling's son. [Nephew², Sibling², Son²]

- $\forall xy. [\text{Nephew}(x,y) \leftrightarrow \exists z. [\text{Sibling}(y,z) \wedge \text{Son}(x,z)]]$

4. A maternal grandmother is a mother's mother. [functions: mgm, mother-of]

- $\forall xy. x=\text{mgm}(y) \rightarrow \exists z. x=\text{mother-of}(z) \wedge z=\text{mother-of}(y)$

Everybody loves somebody [loves2]

$\forall x. \exists y. \text{Loves}(x, y)$

$\exists y. \forall x. \text{Loves}(x, y)$

$\exists y. \exists x. \text{Loves}(x, y)$

$\forall xy. \text{Loves}(x, y)$

Writing More FOL

5. Nobody loves Jane.

- $\forall x. \neg \text{Loves}(x, \text{Jane})$
- $\neg \exists x. \text{Loves}(x, \text{Jane})$

6. Everybody has a father.

- $\forall x. \exists y. \text{Father}(y, x)$

7. Everybody has a father and a mother.

- $\forall x. \exists yz. \text{Father}(y, x) \wedge \text{Mother}(z, x)$

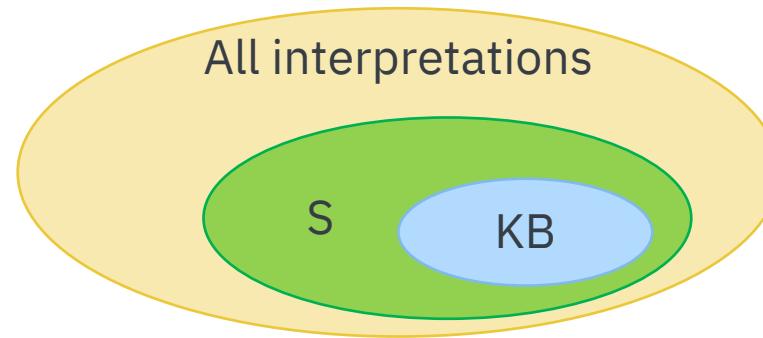
8. Whoever has a father, has a mother.

- $\forall x. [[\exists y. \text{Father}(y, x)] \rightarrow [\exists y. \text{Mother}(y, x)]]$

Entailment in First-Order Logic

KB entails S

- **For every interpretation I, if KB holds in I, then S holds in I.**



- Computing entailment is impossible in general because there are infinitely many possible interpretations.
- Even computing holds is impossible for interpretations with infinite universes.

Intended Interpretations

KB: $(\forall x.\text{Circle}(x) \rightarrow \text{Oval}(x)) \wedge (\forall x.\text{Square}(x) \rightarrow \neg \text{Oval}(x))$

S: $\forall x.\text{Square}(x) \rightarrow \neg \text{Oval}(x)$

- We know $\text{holds}(\text{KB}, I)$.
- We wonder whether $\text{holds}(S, I)$.
- We could ask “Does KB entail S?”
- Or we could just try to check whether $\text{holds}(S, I)$

- $I(\text{Fred}) = \triangle$
- $I(\text{Above}) = \{\langle \square, \triangle \rangle, \langle \bullet, \circ \rangle\}$
- $I(\text{Circle}) = \{\langle \bullet \rangle\}$
- $I(\text{Oval}) = \{\langle \bullet \rangle, \langle \circ \rangle\}$
- $I(\hat{\text{hat}}) = \{\langle \triangle, \square \rangle, \langle \circ, \bullet \rangle, \langle \square, \square \rangle, \langle \bullet, \bullet \rangle\}$
- $I(\text{Square}) = \{\langle \triangle \rangle\}$

An Infinite Interpretation

KB: $(\forall x.\text{Circle}(x) \rightarrow \text{Oval}(x)) \wedge (\forall x.\text{Square}(x) \rightarrow \neg \text{Oval}(x))$

S: $\forall x.\text{Square}(x) \rightarrow \neg \text{Oval}(x)$

- Does KB hold in I_1 ?
 - Yes, but can't verify by enumerating U.
- S also holds in I_1 .
 - Yes, but can't verify by enumerating U.

- $U_1 = \{1, 2, 3, \dots\}$
- $I_1(\text{Circle}) = \{\langle 4, 8, 12, 16, \dots \rangle\}$
- $I_1(\text{Oval}) = \{2, 4, 6, 8, \dots\}$
- $I_1(\text{Square}) = \{1, 3, 5, 7, \dots\}$

An argument for Entailment

KB: $(\forall x. \text{Circle}(x) \rightarrow \text{Oval}(x)) \wedge (\forall x. \text{Square}(x) \rightarrow \neg \text{Oval}(x))$

$S_1: \forall xy. (\text{Circle}(x) \wedge \text{Over}(y) \wedge \neg \text{Circle}(y)) \rightarrow \text{Above}(x,y)$

- $I(\text{Fred}) = \triangle$
- $I(\text{Above}) = \{\langle \square, \triangle \rangle, \langle \bullet, \circ \rangle\}$
- $I(\text{Circle}) = \{\langle \bullet \rangle\}$
- $I(\text{Oval}) = \{\langle \bullet \rangle, \langle \circ \rangle\}$
- $I(\hat{\text{I}}) = \{\langle \triangle, \square \rangle, \langle \circ, \bullet \rangle, \langle \square, \square \rangle, \langle \bullet, \bullet \rangle\}$
- $I(\text{Square}) = \{\langle \triangle \rangle\}$
- $\text{holds}(\text{KB}, I)$
- $\text{holds}(S_1, I)$

- $U_1 = \{1, 2, 3, \dots\}$
- $I_1(\text{Circle}) = \{\langle 4, 8, 12, 16, \dots \rangle\}$
- $I_1(\text{Oval}) = \{2, 4, 6, 8, \dots\}$
- $I_1(\text{Square}) = \{1, 3, 5, 7, \dots\}$
- $I_1(\text{Above}) = >$
- $\text{holds}(\text{KB}, I_1)$
- $\text{fails}(S_1, I_1)$

KB doesn't entail S_1 !

Proof and Entailment

- Entailment captures general notion of “follows from.”
- We can’t evaluate it directly by enumerating interpretations.
- So, we will do proofs.
- In FOL, if S is entailed by KB , then there is a finite proof of S from KB .

Axiomatization

- What if we have a particular interpretation, I , in mind, and want to test whether $\text{holds}(S, I)$?
- Write down a set of sentences, called axioms, that will serve as our KB.
- We would like KB to hold in I , and as few other interpretations as possible.
- No matter what,
 - If $\text{holds}(\text{KB}, I)$ and KB entails S , then $\text{holds}(S, I)$.
- If your axioms are weak, it might be that
 - $\text{holds}(\text{KB}, I)$ and $\text{holds}(S, I)$, but
 - KB doesn't entail S .

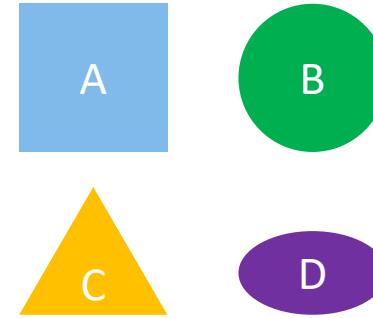
Axiomatization Example

KB₂

- Above(A, C)
- Above(B, D)
- $\forall xy. \text{Above}(x,y) \rightarrow \text{hat}(y) = x$
- $\forall x. (\neg \exists y. \text{Above}(y,x)) \rightarrow \text{hat}(x) = x$

S

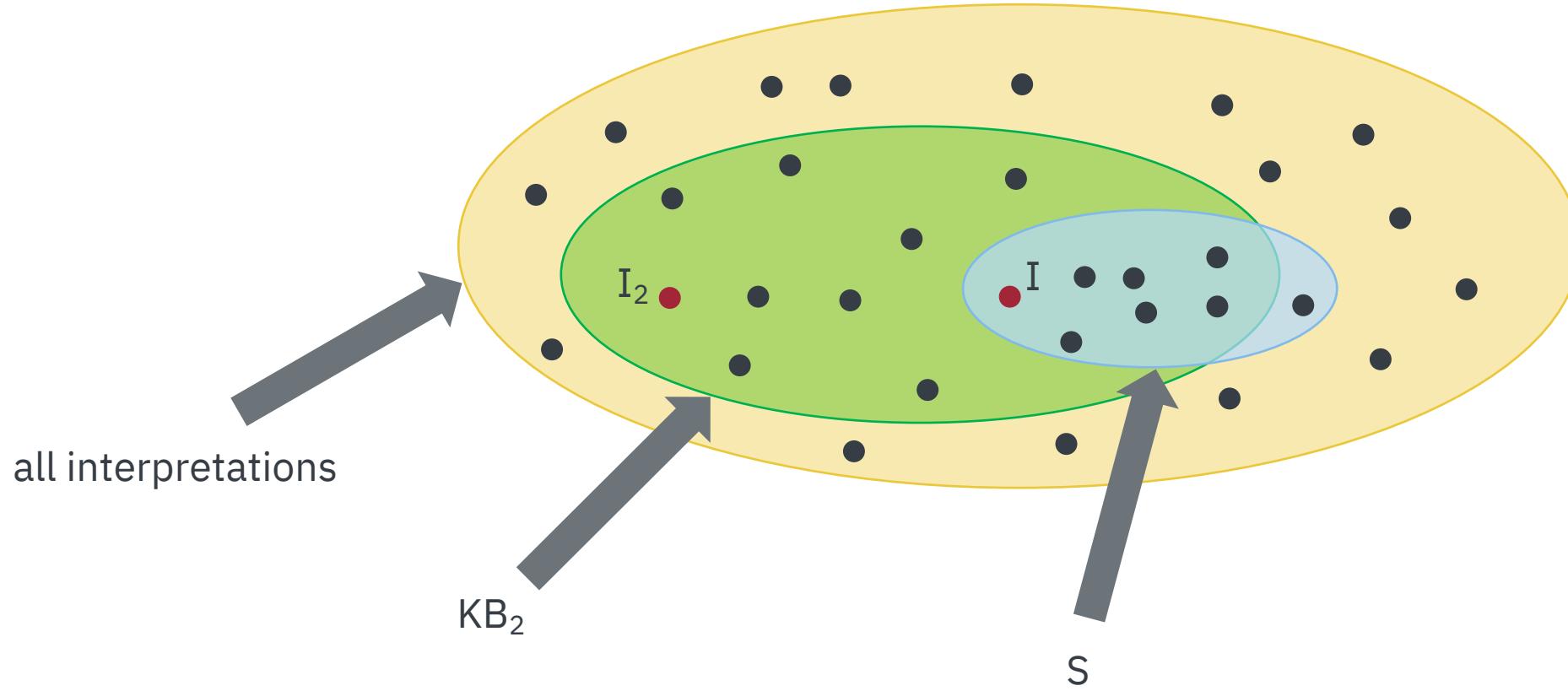
- $\text{hat}(A) = A$



- $I_2(A) = \blacksquare$
- $I_2(B) = \bullet$
- $I_2(C) = \blacktriangle$
- $I_2(D) = \blacklozenge$
- $I_2(\text{Above}) = \{\langle \blacksquare, \blacktriangle \rangle, \langle \bullet, \blacklozenge \rangle, \langle \blacktriangle, \blacksquare \rangle, \langle \blacklozenge, \bullet \rangle\}$
- $I_2(\text{hat}) = \{\langle \blacktriangle, \blacksquare \rangle, \langle \blacklozenge, \bullet \rangle, \langle \bullet, \blacklozenge \rangle, \langle \blacksquare, \blacktriangle \rangle\}$

- holds(KB_2, I_2)
- fails(S, I_2)
- KB_2 doesn't entail S .

KB_2 is a Weakling!



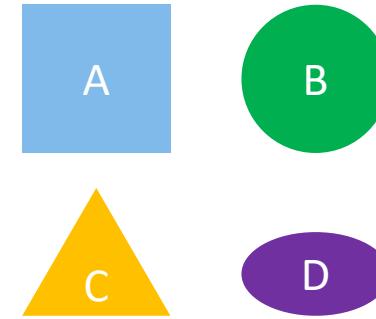
Axiomatization Example: Another Try

KB₃

- Above(A, C)
- Above(B, D)
- $\forall xy. \text{Above}(x,y) \rightarrow \text{hat}(y) = x$
- $\forall x. (\neg \exists y. \text{Above}(y,x)) \rightarrow \text{hat}(x) = x$
- **$\forall xy. \text{Above}(x,y) \rightarrow \neg \text{Above}(y,x)$**

S

- $\text{hat}(A) = A$



- $I_3(A) = \blacksquare$
- $I_3(B) = \bullet$
- $I_3(C) = \blacktriangle$
- $I_3(D) = \blacklozenge$
- $I_3(\text{Above}) = \{\langle \blacksquare, \blacktriangle \rangle, \langle \bullet, \blacklozenge \rangle, \langle \bullet, \blacksquare \rangle\}$
- $I_3(\text{hat}) = \{\langle \blacktriangle, \blacksquare \rangle, \langle \blacklozenge, \bullet \rangle, \langle \bullet, \bullet \rangle, \langle \blacksquare, \bullet \rangle\}$

- fails (KB_3, I_2)
- holds(KB_3, I_3)
- fails(S, I_3)
- KB_3 doesn't entail S .

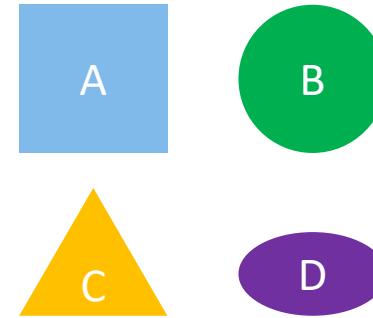
Axiomatization Example: Another Try

KB₄

- Above(A, C)
- Above(B, D)
- $\neg \exists x. \text{Above}(x, A)$
- $\neg \exists x. \text{Above}(x, B)$
- $\forall xy. \text{Above}(x, y) \rightarrow \text{hat}(y) = x$
- $\forall x. (\neg \exists y. \text{Above}(y, x)) \rightarrow \text{hat}(x) = x$

S

- $\text{hat}(A) = A$



- $I_3(A) = \blacksquare$
- $I_3(B) = \bullet$
- $I_3(C) = \blacktriangle$
- $I_3(D) = \blacklozenge$
- $I_3(\text{Above}) = \{\langle \blacksquare, \blacktriangle \rangle, \langle \bullet, \blacklozenge \rangle, \langle \bullet, \blacksquare \rangle\}$
- $I_3(\text{hat}) = \{\langle \blacktriangle, \blacksquare \rangle, \langle \blacklozenge, \bullet \rangle, \langle \bullet, \bullet \rangle, \langle \blacksquare, \bullet \rangle\}$

- fails (KB_4, I_3)
- KB_4 entails S.

We will prove KB_4 entails S.



Questions?

Acknowledgement

Fahiem Bacchus, University of Toronto
Dan Klein, UC Berkeley
Kate Larson, University of Waterloo

