

# Knowledge representation in logic

The state space search algorithms covered earlier had a relatively general formulation, but required the problem to be represented in a specific format. This format included the definition of the state space, the set of state transition operators, and a heuristic state evaluation function.

Generally, the structure and format of knowledge representation are highly important and affect the efficiency — or even the ability — of searching for the solution.

There exist a number of paradigms for knowledge representation in artificial intelligence. These knowledge representation paradigms usually come with associated with them algorithms for **reasoning**, i.e. making subsequent findings ultimately leading to determining the solution to the problem.

One of the most powerful and popular knowledge representation schemes is the language of mathematical logic.

## Why is mathematical logic a good representation language in artificial intelligence?

On the one hand, it is close to the way people think about the world and express their thought in natural language. People even view their way of thinking as “logical”. The categories by which people think and speak include such constructs as: objects and relations between them, simple and complex assertions, sentences, connectives, conditionals, and even quantifiers.

On the other hand, the mathematical logic offers a precise apparatus for reasoning, based on theorem proving. People, likewise, use logical reasoning in their thinking, so mathematical logic seems to be a good representation platform for the knowledge base of an intelligent agent, whose way of expressing facts and reasoning should be similar to the human's.

## Example: the wumpus world

It is useful to have a good testing environment for verifying the methods being developed. This environment needs to be simple enough to allow developing intuitions and quickly discovering properties, but at the same time rich enough to pose some significant demands of the problem solving abilities, and allow to formulate problems of various degree of difficulty.

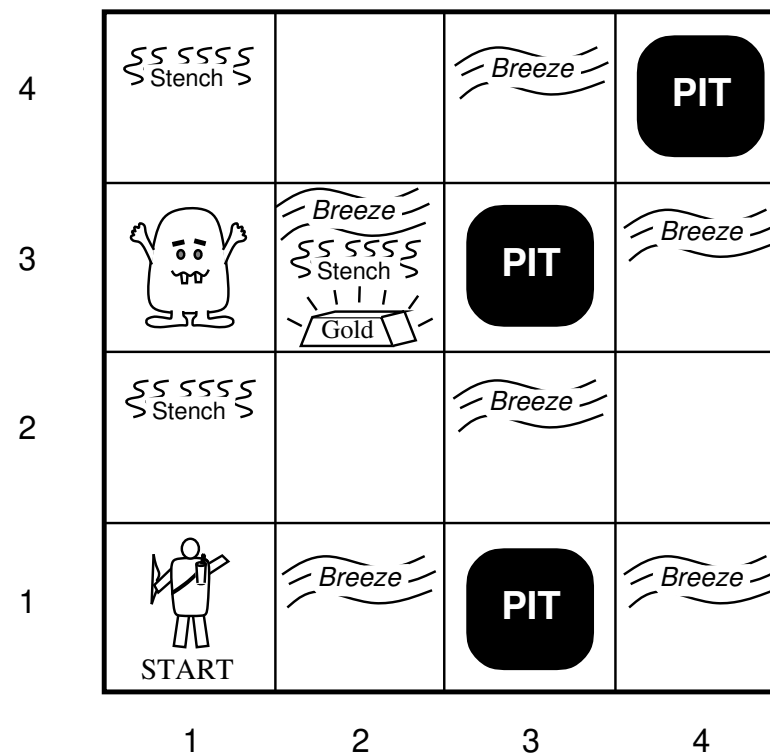
One of such “textbook” testing environment is the **wumpus world**.<sup>1</sup> An intelligent agent moves around this environment in search for gold, which she intends to carry out safely. The agent is however faced with some dangers, such as the *pits*), into which she may fall, and the title *wumpus* monster, which may eat the agent.

The agent may only turn right or left, move forward by one step, shoot a single arrow from a bow (ahead), pick up gold, and leave the environment when she is in the starting position.

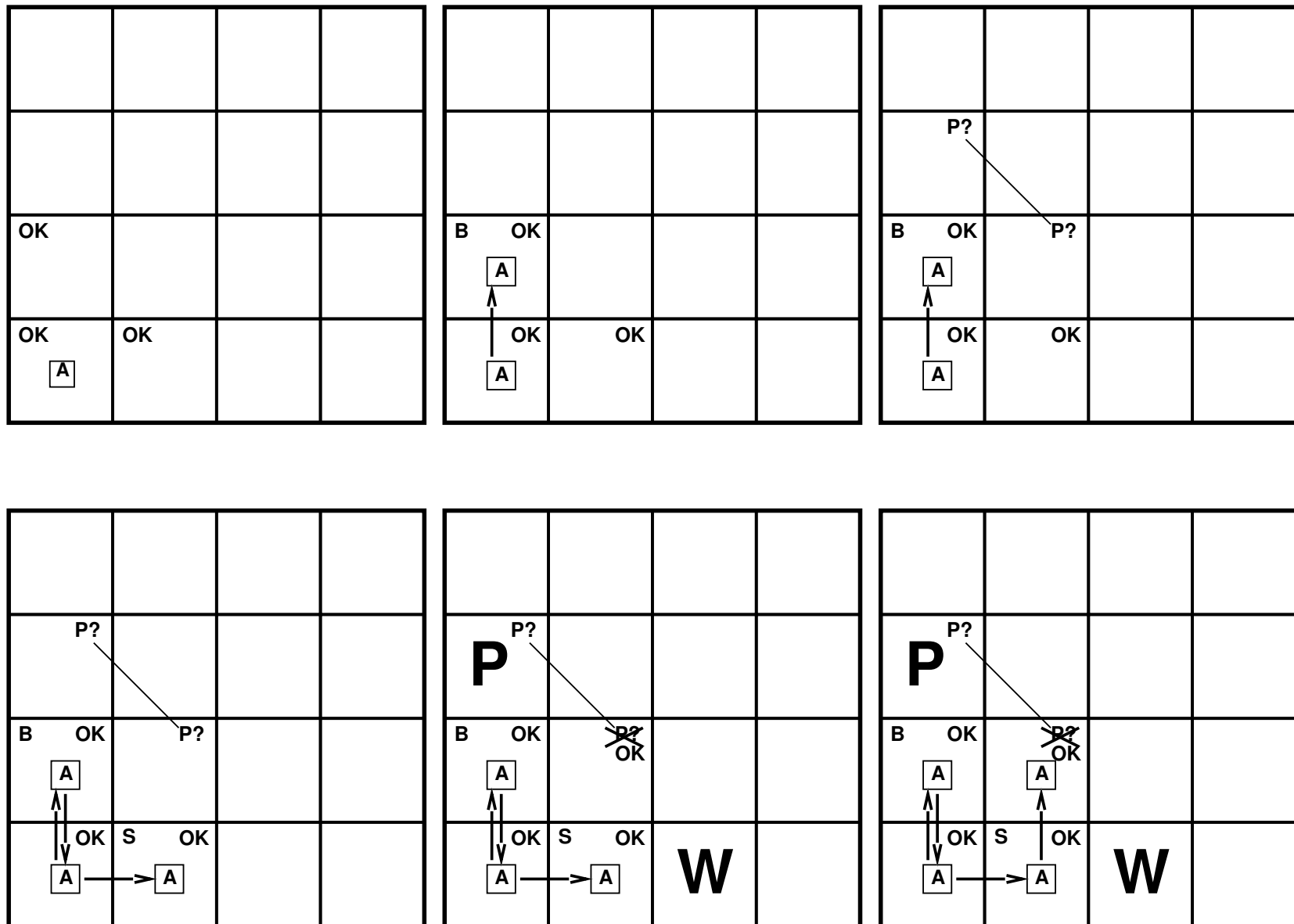
---

<sup>1</sup>The examples and diagrams of the wumpus world presented here are borrowed from the textbook by Russell and Norvig “Artificial Intelligence A Modern Approach” and the materials provided on Stuart Russell’s Web page.

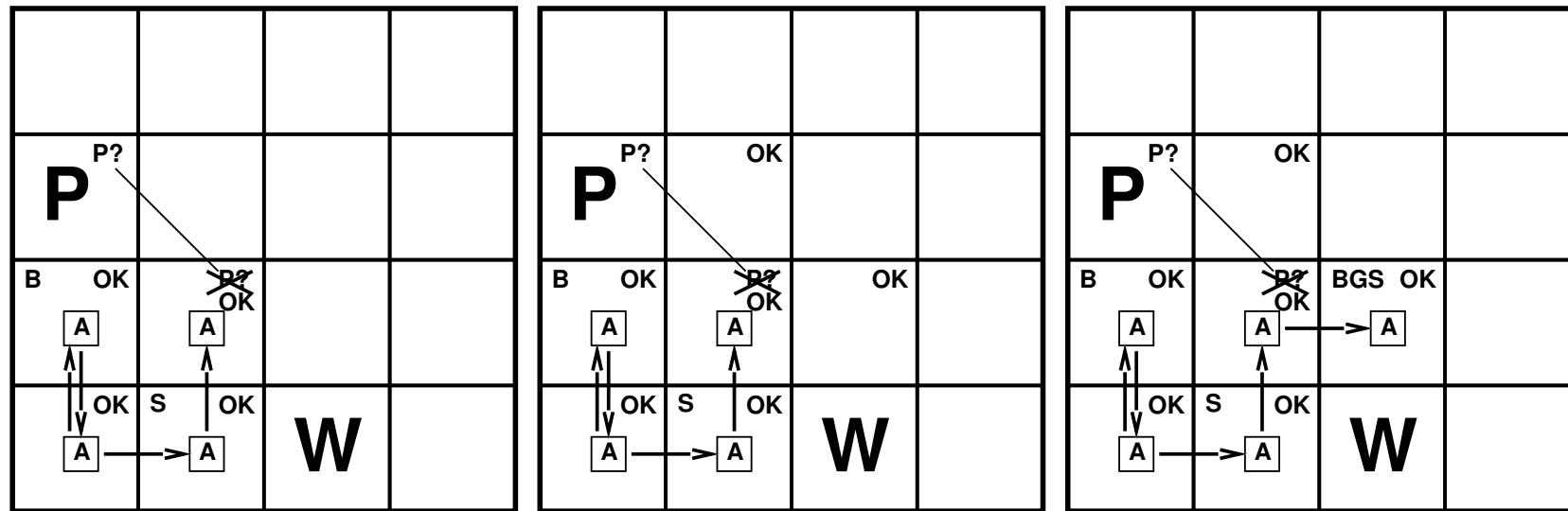
The agent receives some information about her environment (the data obtained by the agent by her perception are called the **percepts**). She can smell the wumpus *stench* and feel the *breeze* from the pits, but only in the fields directly neighboring the wumpus or the pits. She can also detect the *gold*, but only when she enters the field it is in. She cannot determine her absolute position (*à la* GPS), but she can remember her position and covered trail. She can sense the walls only by trying to enter them, which results in getting bumped back.



# Example: moving around the wumpus world

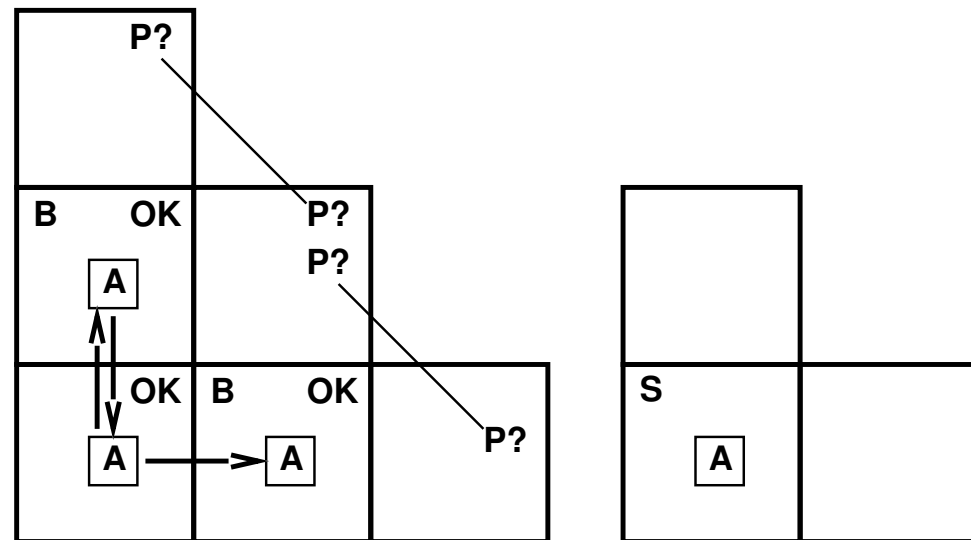


# Example: moving around the wumpus world (contd.)



However, it is not always possible to act so efficiently in the wumpus world by using only logical reasoning.

In some cases the only solution is to “shoot”, ie. blindly select a move, and analyze the outcome. Provided that we survive!!



# Propositional logic: syntax and wffs

**Propositional logic** is a very simple logical language. It allows writing **atomic formulas** based on **propositional symbols**. By writing a logical formula we state some fact. Examples of atomic formulas:  $P, Q, R, WumpusAt\_1\_5, HaveGold$ .

We can also write **complex formulas** which are constructed from other formulas using the **logical connectives**:  $\neg$  (negation),  $\wedge$  (conjunction),  $\vee$  (alternative),  $\Rightarrow$  (implication), and  $\Leftrightarrow$  (biconditional).

Complex formulas can be made up of complex formulas using parentheses, or without parentheses, where not ambiguous. These rules of how legal language expressions, called **well-formed formulas** or wffs, may be created, jointly form the **syntax** of the language.

Examples of wffs:

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

$$\neg \neg P$$

$$(AgentAt\_1\_1 \wedge PitAt\_1\_2) \Rightarrow Breeze$$

$$HaveGold \vee \neg HaveGold$$

$$HaveGold \wedge \neg HaveGold$$

Examples of non-wffs:

$$P \wedge \wedge Q$$

$$P \neg Q$$

$$P(WumpusAt\_1\_5)$$

Explain why these are not wffs.

# Propositional logic: semantics

The syntax defines the language. In the case of propositional logic it consists of: a set of propositional symbols (these can be arbitrary), the set of logical connectives (these are only the five we introduced), and the rules of their use.

The syntax does not concern the meaning of the formulas. This is the role of the **semantics** of the language. The semantics assigns a meaning to each of the propositional symbols. Having defined a meaning of any formula we can start talking about whether it is true or false. Which is the ultimate goal of the logical representation.

Note that we have already written some of the propositional symbols in such a way as to suggest what we intend them to mean: *AgentAt\_1\_1*, *PitAt\_1\_2*, *HaveGold*. Other symbols are just generic; they can be assigned any meaning, abstract or very specific: *P*, *Q*, *R*.

However, the language of logic must be flexible and very general — we can never infer from the mere notation of a propositional symbol what it actually means.



# Propositional logic: semantics — possible worlds, interpretations

If we associate the formula (written with a single propositional symbol) *AgentAt\_1\_1* with the meaning that the wumpus world agent is currently at position (1,1) then it still does not give us a way to verify whether the formula states the truth or not. It is entirely possible, that in one specific instance of the game it is true, while in many other instances it is false.

The semantics resolves this by associating each atomic formula with a **possible world**, which is a concrete configuration of the problem domain being described, where all the objects being described are in precisely defined states. This is done by way of an **interpretation function** which associates each atomic formula (or propositional symbol) with a specific meaning with respect to such possible world, and thus defines the truth value of such formula.

It is important that the interpretation function is completely defined, i.e. each propositional symbol present in the language is associated with some aspect of the possible world, and the corresponding atomic formula can be unambiguously interpreted as having value 1 or 0.

# Propositional logic: semantics — interpretations, models

The interpretation using the possible world on the left assigns the formula *AgentAt\_1\_1* the truth value 1 (or True), while another one using the possible world on the right assigns the same formula the truth value 0 (or False):

			P
		W	
	P		
A			

		P	
A			
	P		W

Possible worlds are also more precisely referred to as **models**.

Note that the locations of all the objects (agent, wumpus, pits), if described by the propositional symbols of the language, must be specified by each model, whether the wumpus world agent knows these locations or not.

With these models it is not possible, for example, to have another object  $F$  whose position might be described by formulas such as  $FAt\_2\_2$ . If this was the case then the above configurations would not be models for such a problem domain, as they do not reflect the location of the  $F$  object.

# Propositional logic: semantics — formula satisfaction

Given a specific propositional formula, atomic or complex, some models assign it the truth value 1, while others assign it the truth value 0. Note that there is no other option. All the objects described by the propositional symbols are present in the model, and all their properties are reflected there.

We will say that a model  $m$  **satisfies** a formula  $f$  if it assigns it the truth value 1. We will also say that a model  $m$ , which satisfies a formula  $f$ , is **a model of this formula**.

Note the different meaning of the word: model. Any model (for a specific problem domain) can either satisfy a given formula, or falsify it. But if it satisfies it, then it is **the model of this formula**.

For the sake of definition, if a formula is satisfied by all models, then it is called a **tautology**. An example of a tautology is  $P \vee \neg P$ . Its truth value does not depend on the model — it must be assigned the truth value 1 with any model.

Conversely, if a formula cannot be satisfied by any model, then it is called **unsatisfiable**. An example of an unsatisfiable formula may be  $P \wedge \neg P$ . Its truth value also does not depend on the model — it is a constant 0.

# Propositional logic: semantics — complex formula satisfaction

Let's stop for a moment to note an important detail. Primarily, a model defines the truth values of all the atomic formulas (propositional symbols). Once that is done, all other formulas (complex) have their truth values determined by the semantics of the specific logical connectives.

For example, assume a model  $m$  assigns the symbol  $AtAgent\_1\_1$  the value 1. Then consider the formula  $\neg AtAgent\_1\_1$ . We are not free to choose any truth value we like for it; we are obliged to assume its value is 0. The same applies to any formulas using  $\wedge, \vee$ , etc. The truth values of any formulas containing them are defined by their truth tables.

At the same time, we would like to make sure that this model assigns 0 to all the formulas (atomic) of the kind:  $AtAgent\_1\_2, AtAgent\_2\_1, AtAgent\_2\_2, \dots$ . But the propositional logic's rules do not enforce this. The preceding formulas are not associated in any way, so a model can assign them any truth values. Having any of them equal to 1 would break the rules of the wumpus world (since there should be only one agent, and can be only in one place at a time), but from the logical point of view nothing would be wrong.

# Propositional logic: semantics — the interpretation functions

In principle, an interpretation function assigns an atomic formula some model (which may be a model of this formula, or not). There are a lot of models (possible worlds) which can be considered. But from the point of view of establishing truth values of formulas, the only important thing is which propositional symbols they satisfy, and which they do not.

Since in a specific problem domain there might be only some number of propositional symbols (determined by the number of objects and their significant properties), there are only  $2^N$  ( $N$  - the number of the propositional symbols) types of models which really count: those which satisfy a specific symbol, and those which falsify it.

For this reason, for the propositional logic, we often dismiss the great variety of possible worlds, and reduce the set of models to the set of different 0/1  $N$ -tuples associating truth values to all the propositional symbols:

	<i>AgentAt_1_1</i>	<i>WumpusAt_1_5</i>	...
$m_1$	0	0	...
$m_2$	0	1	...
$m_3$	1	0	...
$m_4$	1	1	...
...	...	...	...

# Propositional logic: semantics — sets of models

Given the previous generalization, we can alternatively start looking at logical formulas as a compact way of representing sets of models, namely those which are models of a specific formula.

For example, the formula *AgentAt\_1\_1* can be thought of as representing all the models where the agent is located in position (1,1).

# Propositional logic: some laws of logic

Since for a given (complex) formula, we will be primarily interested in determining its truth value, it could be advantageous to note some transformations which can be done on logical formulas while preserving their truth value. In the following formulas the equivalence symbol  $\equiv$  is used to denote that one side can be replaced by the other in the process of truth evaluation of a formula.

## Associativity:

$$\begin{aligned}p \wedge (q \wedge r) &\equiv (p \wedge q) \wedge r \\p \vee (q \vee r) &\equiv (p \vee q) \vee r\end{aligned}$$

Because of the associativity of both conjunction and alternative, we can write multiple consecutive applications of either of these connectives without parentheses. This is because the truth value of the formula with multiple applications of either of these connectives does not depend on the order in which the connectives are interpreted:

$$\begin{aligned}p_1 \wedge (p_2 \wedge (p_3 \wedge (...))) &\equiv p_1 \wedge p_2 \wedge p_3 \wedge ... \\p_1 \vee (p_2 \vee (p_3 \vee (...))) &\equiv p_1 \vee p_2 \vee p_3 \vee ...\end{aligned}$$

## Distributivity:

$$p \wedge (q \vee r) \equiv (p \wedge q) \vee (p \wedge r)$$

$$p \vee (q \wedge r) \equiv (p \vee q) \wedge (p \vee r)$$

## de Morgan laws:

$$\neg(p \wedge q) \equiv (\neg p) \vee (\neg q)$$

$$\neg(p \vee q) \equiv (\neg p) \wedge (\neg q)$$

## Other useful identities:

$$p \Rightarrow q \equiv \neg p \vee q$$



# Logical reasoning — entailment

A **knowledge base (KB)** is a set of formulas representing a conjunction of all its member formulas. The set of models of such a set is an intersection of the sets of models of the member formulas.

In Artificial Intelligence a knowledge base is typically the database of all facts that an AI agent possesses. **An AI agent's knowledge is naturally a conjunction of many facts.** A typical activity of such an agent is trying to answer a question, of whether another fact, represented by a formula  $f$  holds with such a knowledge base.

It may happen, that assuming all formulas of some set KB are true, another formula  $f$  must also be true, under all possible interpretations.

In such case we say that KB **entails**  $f$ , written  $\text{KB} \models f$ .

Alternatively, we can say the  $f$  **logically follows** from the set of formulas KB.

Examples:  $\{P, Q\} \models P \wedge Q$      $\{P \wedge Q\} \models P$      $\{P \vee Q, \neg P\} \models Q$      $\{P \Rightarrow Q, P\} \models Q$

**In propositional logic, one way of determining entailment is by truth tables. In such a table we enumerate all models (truth assignments for propositional symbols) and verify that all those models which satisfy the full KB, also satisfy  $f$ .**

# Short review

For the following examples answer whether the specified entailment holds.

1.  $\{P \vee Q\} \models P \wedge Q$
2.  $\{P \wedge Q\} \models P \vee Q$
3.  $\{P, Q\} \models P \Rightarrow Q$
4.  $\{P, Q\} \models \neg P \vee Q$
5.  $\{P \Rightarrow Q, \neg Q\} \models \neg P$
6.  $\{P \Rightarrow Q, \neg P\} \models \neg Q$
7.  $\{P \Rightarrow Q, \neg P\} \models Q$
8.  $\{P \Rightarrow Q, Q\} \models P$
9.  $\{P \Rightarrow Q, Q \Rightarrow R\} \models R$

# Logical reasoning — modus ponens

In some cases we can apply a process on logical formulas called an **inference**.

Example:

It is raining.	(Raining)
If it is raining, then the road is wet.	(Raining $\Rightarrow$ RoadWet)
Conclusion: the road is wet.	(RoadWet)

This is an example of an **inference rule** called **modus ponens**:

For any propositional symbols  $p$  and  $q$ :

$$\frac{p, \quad p \Rightarrow q}{q}$$

or more generally:

For any  $p_1, \dots, p_k, q$ :

$$\frac{p_1, \quad \dots, \quad p_k, \quad (p_1 \wedge \dots \wedge p_k) \Rightarrow q}{q}$$

# Logical reasoning — inference rules

There can be other inference rules, generally written according to a scheme:

$$\frac{f_1, \dots, f_k \text{ (premises)}}{g \text{ (conclusion)}}$$

We can use inference rules in a reasoning process, by applying them successively on formulas from the knowledge base, until a desired conclusion is obtained. This is the process of **inferencing**, and any formula  $f$  obtained in this process is said to be **derived** from KB, written  $\text{KB} \vdash f$ .

The **forward inferencing** algorithm:

```
repeat until no change to KB:
  foreach inference rule  $\frac{f_1, \dots, f_k}{g}$ 
    if  $f_1, \dots, f_k \in \text{KB} \wedge g \notin \text{KB}$ 
      add  $g$  to KB
```

If  $f$  gets eventually added to KB then  $\text{KB} \vdash f$ .

# Derivation versus entailment of formulas

Note the inferencing process operates strictly in the syntax domain. It operates on the formulas as they are written, and does not refer to models or truth checking.

A question therefore arises:

How does  $KB \models f$  relate to  $KB \vdash f$  ?

Are these equivalent?

Can any entailed formula be derived?

## Short review

For the following examples answer whether the specified derivability holds.

First assume that the only inference rule is the modus ponens:  $\frac{\phi, \phi \Rightarrow \psi}{\psi}$

In addition to the inference steps you can use the logic equivalence laws such as those presented on pages 15 through 16 to convert formulas to the desired equivalent form, both in the KB set and in the target formula on the right-hand side.

1.  $\{P, Q\} \vdash P \wedge Q$
2.  $\{P \wedge Q\} \vdash P$
3.  $\{P \wedge Q\} \vdash P \vee Q$
4.  $\{P \Rightarrow Q, \neg Q\} \vdash \neg P$

Now assume that in addition to modus ponens you can also use the following inference rules:  $\frac{\phi \wedge \psi}{\phi}$  (conjunction elimination),  $\frac{\phi, \psi}{\phi \wedge \psi}$  (conjunction introduction), and  $\frac{\phi}{\phi \vee \psi}$  (disjunction introduction).

First answer again the above examples, and then also the following:

5.  $\{P, Q\} \vdash P \Rightarrow Q$
6.  $\{\neg P, Q\} \vdash P \Rightarrow Q$
7.  $\{P \Rightarrow Q, P \vee Q\} \vdash Q$
8.  $\{P \Rightarrow Q, P \wedge R\} \vdash Q \wedge R$

# Inference rules: soundness and completeness of inference rules

An inference rule is **sound** if it only allows to derive from any KB formulas which are entailed by that KB. (But it may not derive ALL such formulas.)

An inference rule is **complete** if it allows to derive from any KB ALL the formulas which are entailed by that KB. (But it may also derive other formulas which are not entailed by the KB.)

It follows from the above, that if we had an inference rule which was both sound and complete, then we could use the inferencing process, instead of checking all the models (truth tables).

But it is not easy to find such an inference rule. For example, modus ponens is sound, but is not complete. To see this, consider that for the previous example:

$$KB = \{ \text{Raining}, \text{Raining} \Rightarrow \text{RoadWet} \}$$

we have been able to derive RoadWet using modus ponens, but this KB also entails  $(\text{Raining} \wedge \text{Raining} \Rightarrow \text{RoadWet})$ , which the modus ponens rule cannot derive.

# Inference rules: soundness and completeness for sets of rules

Suppose we have more than one inference rule. How does soundness/completeness work in this case?

It should be clear, that however many inference rules we consider, we want each one of them to be (individually) sound. A single unsound rule would allow the inferencing algorithm to introduce contradictory conclusions to the knowledge base, regardless of what other rules can offer.

Example: consider two hypothetical inference rules:  $\frac{p, q}{p \vee q}$ ,  $\frac{p, q}{\neg(p \vee q)}$

The first is sound, and the second is not.

If both are present in the system, the inferencing algorithm will have to derive both  $(p \vee q)$  and  $\neg(p \vee q)$ . No matter how many sound rules are present, one unsound rule can spoil the whole system by allowing to derive a false formula.



# Inference rules: soundness and completeness for sets of rules

## (2)

Completeness is a different story. It may be the case for some set of inference rules  $\{IR_1, IR_2, \dots, IR_n\}$ , that the inferencing system would be able for any KB to derive all the formulas which are entailed by this KB. Thus a set of inference rules can be complete as a whole, even though any or all of these rules could individually not be complete.

So one way to fix the inferencing system to have derivation be equivalent to entailment, would be to find a set of sound inference rules, which together would be complete.

It is possible, but there is also another solution.



# The Conjunctive Normal Form

Some terminology:

We will call a **literal** any atomic formula, or a negation thereof.

We will call a **clause** a formula which is an alternative of literals.

A formula, which is a conjunction of clauses will be said to be in the **Conjunctive Normal Form (CNF)**.

In short, we can say, that a CNF formula is a conjunction of clauses.

Examples:  $[(P \vee Q \vee \neg R) \wedge (P \vee \neg Q) \wedge R], (P \wedge Q), (P \vee Q), \neg P, P$

Non-examples:  $(P \Rightarrow Q) \wedge (Q \Rightarrow P), \neg(P \wedge Q)$

A fact: any propositional formula can be converted to an equivalent formula in the CNF form.

CNF formulas are useful, because they allow to be inferenced over, using resolution.

And, what is equally important, this process can be made completely automatic.

# Inference rules: resolution

By using a well-known and useful identity:

$$(p \Rightarrow q) \equiv (\neg p \vee q)$$

we can rewrite the modus ponens to a different form:

$$\frac{p, \neg p \vee q}{q}$$

Note the red terms cancel out, in a sense. This actually makes sense, because if we know  $p$  to be true, then the  $\neg p$  is certainly false, and can be dropped from the alternative, leaving only  $q$  as a new conclusion.

This observation can be generalized to the following **resolution inference rule**:

$$\frac{p_1 \vee \dots \vee p_n \vee q, \neg q \vee r_1 \vee \dots \vee r_m}{p_1 \vee \dots \vee p_n \vee r_1 \vee \dots \vee r_m}$$

Resolution can be the used as the only inference rule in a sound and complete theorem proving system.

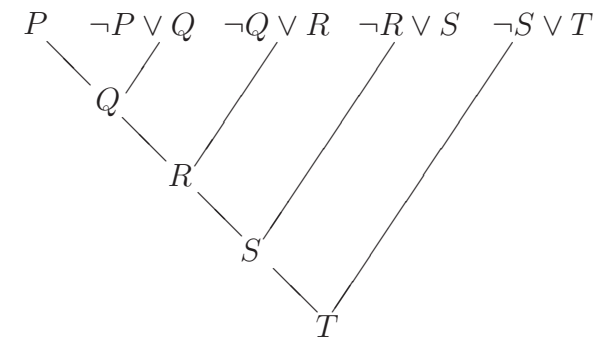
# Logical reasoning using resolution — examples

Let us consider some common patterns of logical reasoning. Suppose that whenever we know  $P$  then also  $Q$  is true, and whenever  $Q$  then also  $R$ , whenever  $R$  then  $S$ , and whenever  $S$  then also  $T$ . And suppose we also know  $P$ . Then we should be able to infer all these fact as the result of a chain of applications of the modus ponens inference rule. Let's see how it works with the CNF form and resolution.

Original facts:  $P, P \Rightarrow Q, Q \Rightarrow R, R \Rightarrow S, S \Rightarrow T$

The same facts in the CNF form, and the graphical representation of a chain of resolution inference steps:

(Note that in the graphical form the chain of the inference steps forms a shape of a tree. This is typical.)



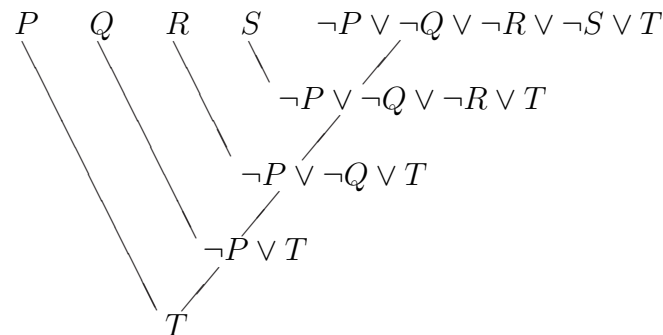
In the above inference scheme, all steps performed on the CNF clauses could equally well be performed with the modus ponens rule on the original implication formulas. But it does not always work this way.

Now let us consider a different reasoning pattern. Suppose several facts:  $P, Q, R, S$  are all known to be true. Further suppose, that these facts, when all combined, imply  $T$ .

Original facts:  $P, Q, R, S, (P \wedge Q \wedge R \wedge S) \Rightarrow T$ .

Facts in the CNF form:  $P, Q, R, S, (\neg P \vee \neg Q \vee \neg R \vee \neg S \vee T)$ .

The resolution tree:



This time, the inference of the final formula  $T$  could not be obtained using modus ponens. For this, we would first have to obtain the formula  $P \wedge Q \wedge R \wedge S$ , which is entailed by the set of the original facts, but cannot be derived using modus ponens. But neither could this be done with resolution. It succeeded to derive  $T$  from the CNF clauses of the original facts, but likewise would not be able to derive the plain conjunction of the original known facts, because the resolution can only produce results by merging two clauses with canceled conflicting literals.

In order to be able to derive formulas like the above conjunction, we need to use the resolution in a special way.

# Empty clauses

We can talk of a single literal as of a unary clause, i.e. an alternative of just this one literal. Moreover, we allow empty clauses, which are treated as alternatives of zero literals. This can be explained using a functional notation for the alternative, thanks to associativity.

$$\begin{aligned} p_1 \vee p_2 \vee \dots \vee p_n &\equiv \vee(p_1, p_2, \dots, p_n) \\ p \vee q \vee r &\equiv \vee(p, q, r) \\ p \vee q &\equiv \vee(p, q) \\ p &\equiv \vee(p) \\ \square &\equiv \vee() \end{aligned}$$

While the truth value of any nonempty clause depends on the truth values of its components, the empty clause must have a constant logical interpretation. By a simple generalization of the definition of the logical values of the alternative we can obtain that **the empty clause is a false (unsatisfiable) formula**. Since we will need to use the empty clause in logical notation, the symbol  $\square$  is used to denote it.

The empty clause can be treated as a neutral element for the alternative connective:

$$\square \vee p \equiv p \equiv p \vee \square$$

# Resolution-based refutation reasoning

A sound and complete theorem proving system can be set up using resolution by using **refutation reasoning**. Whenever we want to obtain:

$$KB \vdash f$$

we add the negation of the theorem formula  $\neg f$  to the  $KB$  set and — hoping that the resulting set of formulas  $KB \cup \{\neg f\}$  is now inconsistent (false) — try to derive an unsatisfiable (false) formula. Assuming that the original  $KB$  is satisfiable, the only source of unsatisfiability can be the added  $(\neg f)$  formula, which proves the original  $f$ .

Since a resolution-based inferencing system works with clauses, the result which is sought in this refutation process is an empty clause. If it can be obtained, the proof is complete. If it is not possible to obtain the empty clause, the theorem we tried to prove must be false (in the propositional logic).

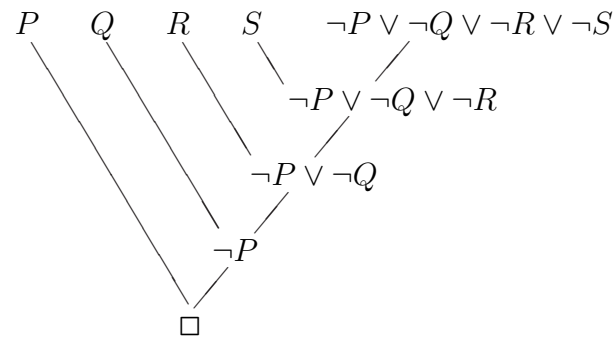
Note however, that a failure to derive the empty clause is not by itself a proof of the falsity of the theorem, just as a failure to find a proof does not mean that one does not exist. But if a search for an empty clause is organized in such way, that it is complete, for example, by assuring that all possible resolution inference steps are tried, then the conclusion of the theorem being false can be correctly made.



# Resolution refutation reasoning — an example

For a simple illustration of a resolution refutation proof let us consider the previous case of having in the database four facts:  $P, Q, R, S$ , and trying to derive their conjunction:  $P \wedge Q \wedge R \wedge S$ . The derivation:  $\{P, Q, R, S\} \vdash (P \wedge Q \wedge R \wedge S)$  is not possible by either modus ponens or resolution used in a straightforward way.

Trying the resolution refutation, the negation of the theorem turns out to be a single clause:  $\neg P \vee \neg Q \vee \neg R \vee \neg S$ . And the sequence of the steps leading to an empty clause is straightforward:



# Short review

For the following set of formulas, write all possible to obtain resolvents.

If it is not possible to perform any resolution, then give a short explanation.

Compare the computed resolvents with logical consequences you can derive intuitively from the formulas given.

Pay attention to commas, to correctly identify formulas in sets.

1.  $\{ P \vee Q , \neg P \vee \neg Q \}$
2.  $\{ P \Rightarrow Q , Q \Rightarrow R \}$
3.  $\{ \neg P \Rightarrow Q , Q \Rightarrow R \}$
4.  $\{ P \vee Q \vee R , \neg P \vee Q \vee R \}$
5.  $\{ P \vee Q \vee R , \neg P \vee \neg Q \vee \neg R \}$
6.  $\{ P \vee Q , P \vee \neg Q , \neg P \vee Q \}$
7.  $\{ P \Rightarrow (Q \vee R) , \neg Q \wedge \neg R \}$
8.  $\{ P \Rightarrow Q , R \Rightarrow Q , P \vee R \}$

# First order predicate calculus — terms

The **terms** represent objects in the language of logic and may be: constants (denoting a specific object), variables (can assume the values of various objects), or functions (determine an object from the value of their object argument(s), or map some objects into some others).

Examples of terms:  $A$ ,  $123$ ,  $x$ ,  $f(A)$ ,  $f(g(x))$ ,  $+(x, 1)$

By convention, we will write constant terms in capital letters, and variables in lowercase.

Let us make a note, that the last term in the above examples is an indirect notation of the subsequent value for  $x$ , and not a subtraction. In pure logic there is no arithmetic. We will see the consequences of this often.

# First order predicate calculus — predicates

The **predicates** represent relations over the set of terms. We can treat them as functions assuming the values of true or false (1 or 0), assigning 1 to each vector of  $n$  terms satisfying the relation, and 0 to each vector of  $n$  terms not satisfying the relation.

A predicate symbol written with the set of terms is called an **atomic formula**.

Examples of atomic formulas:  $P$ ,  $Q(A)$ ,  $R(x, f(A))$ ,  $> (x, 10)$

The expression  $> (x, 10)$  is the functional equivalent of  $x > 10$ . In arithmetic we treat such an expression as inequality and we could solve it. But as a logical formula we can only **evaluate** it, meaning determine its truth value. But if a formula contains a variable then often its truth value cannot be determined.

# Representing facts with logical formulas

What is the purpose of the predicate language?

We could use it to write the facts we want to express, like:

$$\begin{aligned} &At(Wumpus, 2, 2) \\ &At(Agent, 1, 1) \\ &At(Gold, 3, 2) \end{aligned}$$

The selection of the set of symbols, by which we intend to describe the objects and relations of some world is called **conceptualization**. For example, an alternative conceptualization for the above facts could be the following:

$$\begin{aligned} &AtWumpus(loc(2, 2)) \\ &AtAgent(loc(1, 1)) \\ &AtGold(loc(3, 2)) \end{aligned}$$

These two conceptualizations are similar, but have different properties. For example, in the latter the wumpus, agent and gold are not mentioned directly. In general, the accepted conceptualization has influence on the ease or even the ability to express different facts about the problem domain.

## Representing facts with logical formulas (contd.)

A problem with the conceptualization of the wumpus world is the description of the presence and location of the pits. We could give the pits full citizenship and identity:

$$At(Pit_4, 3, 3)$$

In this way it would be easy to describe the “bird’s eye view” of the wumpus world, by giving different pits some names (constant terms). But from the point of view of the agent acting in the wumpus world this conceptualization is very uncomfortable. It would be hard to describe the world as it is gradually learned, when at first the agent does not even know the total number of pits. The presence of a pit at some location would have to be described by a variable:

$$At(x, 3, 3)$$

Unfortunately, this description does not indicate that  $x$  is a pit so this requires further descriptions. A comfortable alternative is to view the pits as anonymous, and only denote the presence or absence of pits at specific locations:

$$PitAt(3, 3)$$

$$NoPitAt(1, 1)$$

# Logical connectives and complex formulas

**Complex formulas** can be constructed from atomic formulas using the **logical connectives**:  $\neg, \wedge, \vee, \Rightarrow, \Leftrightarrow$ . As a special case, an atomic formula or a negation of an atomic formula is called a **literal**.

Examples of complex formulas (the first one is a single literal):

$$\neg At(Wumpus, 1, 1)$$

$$PitAt(2, 1) \vee PitAt(1, 2)$$

$$[At(Agent, 1, 1) \wedge PitAt(2, 1)] \Rightarrow Percept(Breeze)$$

Let us observe that the last formula is of a different nature. The first two could be a fragment of a world description obtained or constructed by the intelligent agent during her activity in the wumpus world. But the last one expresses one of the laws of this world. The agent knows this law and can hold such a formula in her knowledge base.

The facts generally true in a problem domain are called the **axioms of the world**. The facts describing a specific instance of the problem are called **incidental**.

# Quantifiers

The complex formulas can also be built using the **quantifiers**:  $\forall, \exists$ , which **bind** variables in formulas. The general scheme for the formula with a quantifier is:

$$\forall x P(x)$$

A variable not bound by a quantifier in a formula is called **free**. The formula:

$$\exists y Q(x, y)$$

contains two variables, one free ( $x$ ) and one bound by a quantifier ( $y$ ).

A **sentence**, or a **closed** formula is a formula without free variables.

Examples:

$$\exists x, y \text{ At}(\text{Gold}, x, y)$$

$$\forall x, y [\text{At}(\text{Wumpus}, x, y) \wedge \text{At}(\text{Agent}, x, y)] \Rightarrow \text{AgentDead}$$

$$\forall x, y [\text{At}(\text{Wumpus}, x, y) \wedge \text{At}(\text{Agent}, -(x, 1), y)] \Rightarrow \text{Percept}(\text{Stench})$$

Let's note that the  $-(x, 1)$  is an indirect notation of the column left of  $x$ , and not a subtraction. There is no subtracting in logic.



# Short review

1. Work out a complete representation for the wumpus world in the first order predicate calculus. That is: introduce term symbols (constants and functions), and predicate symbols necessary to describe problem instances for the domain.

Note: we do not consider the process of searching for the solution, analyzing alternative moves and their consequences, describing sequences of steps etc. We only seek a scheme for describing static snapshots of a problem instance.

2. Using the representation worked out in the previous question, describe a problem instance given on page 4.
3. Try to write the axioms for the wumpus world, that is, the general rules for this domain.



# Rewriting logical formulas as sets of clauses

A variable-free formula can be converted to a set of clauses, also called the **prenex** form, where all quantifiers are written in front of the formula:

- (i) rename the variables bound by quantifiers to be unique,
- (ii) replace all logical connectives with only conjunctions and alternatives,
- (iii) move negations to inside the formulas (to predicate symbols),
- (iv) extract the quantifiers outside the formula,
- (v) convert the formula to CNF,
- (vi) replace all existential quantifiers with Skolem functions.

The first five steps are logically equivalent transformations (as long as the right order of the extracted quantifiers is maintained in step (iv)). The (vi) step, called **skolemization**, converts all the formulas of the form:

$$\forall x_1 \forall x_2 \dots \forall x_n \exists y \Phi(x_1, x_2, \dots, x_n, y)$$

with:

$$\forall x_1 \forall x_2 \dots \forall x_n \Phi(x_1, x_2, \dots, x_n, f_y(x_1, x_2, \dots, x_n))$$

where  $f_y$  is a newly introduced functional symbol called the **Skolem function**. In case there are no universal quantifiers  $\forall$  this will be a **Skolem constant**.

# Skolem's theorem

The last step in the algorithm for the conversion of formulas into the prenex form is not a logically equivalent transformation. That means, that for the original formula  $\Phi$  and the resulting prenex formula  $\Phi'$ , in general  $\Phi \not\equiv \Phi'$ .

However, the following property, called the **Skolem theorem** holds: for a closed formula  $\Phi$ , if  $\Phi'$  is its prenex form, then  $\Phi$  is satisfiable if and only if  $\Phi'$  is satisfiable.

Therefore, while we cannot in general use the derived prenex form  $\Phi'$  for any logical reasoning instead of  $\Phi$ , we can use it for proving satisfiability (or unsatisfiability).

There exists an extremely useful inferencing scheme, using formulas in prenex form, often written as sets (or lists) of clauses, with clauses written as sets (or lists) of literals.

# Short review

Convert to prenex form the following first order predicate calculus formulas:

1.  $(\exists x P(x)) \wedge (\exists x Q(x))$
2.  $\exists x \forall y \exists z [P(x) \Rightarrow Q(y, z)]$
3.  $\exists x \forall y [P(x, y) \Rightarrow Q(A, x)]$
4.  $\forall x \exists y [P(x, y) \Rightarrow Q(y, f(y))]$
5.  $\forall x [(P(x) \Rightarrow Q(x)) \wedge (P(x) \Rightarrow R(x))]$
6.  $\forall x [(P(x) \wedge Q(x)) \vee (R(x) \wedge S(x))]$



# Substituting variables in formulas

We shall consider transformations of formulas consisting in replacing variable occurrences with other expressions (terms). Since the variables in prenex formulas are implicitly bound with universal quantifiers, replacing variables with other terms means taking specific cases of the formula.

We will call a **substitution** a set of mappings indicating terms to be substituted for specific variables. The term may not contain the variable it is to replace. An example of a substitution:  $s = \{x \mapsto A, y \mapsto f(z)\}$ .

**Applying a substitution** works by syntactically replacing all the occurrences of a given variable within a formula with its associated term. All replacements are done simultaneously, so e.g. by applying the substitution  $s = \{x \mapsto y, y \mapsto A\}$  to the term  $f(x, y)$  the result will be the term  $f(y, A)$ .

Note that this way it does not matter in which order the variables are substituted, even though a substitution is a set (unordered).

A **composition** of substitutions  $s_1$  and  $s_2$  (written as:  $s_1s_2$ ) is called a substitution obtained by applying the substitution  $s_2$  on terms from  $s_1$ , and appending to the resulting set all the pairs from  $s_2$  with variables not in  $s_1$ .

$$\begin{aligned}\Phi s_1s_2 &= (\Phi s_1)s_2 \\ s_1(s_2s_3) &= (s_1s_2)s_3\end{aligned}$$



# Unification

**Unification** is the procedure of finding a substitution of terms to variables in a set of formulas, to reduce it to a singleton set (or to logically equivalent formulas, see explanation below).

A **unifier** of a set of formulas is a substitution reducing it to a singleton set. A set of formulas is **unifiable** if there exists a unifier for it.

For example, the set  $\{P(x), P(A)\}$  is unifiable, and its unifier is  $s = \{x \mapsto A\}$ .

Likewise, the set  $\{P(x), P(y), P(A)\}$  is unifiable, and its unifier is  $s = \{x \mapsto A, y \mapsto A\}$ .

The set  $\{P(A), P(B)\}$  is not unifiable, and neither is  $\{P(A), Q(x)\}$ .

## Unification (contd.)

While unification is a general procedure, here we will compute it only on sets of clauses. Consider the following example clause sets:

$$\Phi = \{P \vee Q(x), P \vee Q(A), P \vee Q(y)\}$$

$$\Psi = \{P \vee Q(x), P \vee Q(A), P \vee Q(f(y))\}$$

$$\Omega = \{P \vee Q(x), P \vee Q(A) \vee Q(y)\}$$

The set  $\Phi$  is unifiable, its unifier is:  $s = \{x \mapsto A, y \mapsto A\}$ , and the unified set is the singleton set:  $\Phi_s = \{P \wedge Q(A)\}$ .

The set  $\Psi$  is not unifiable.

The set  $\Omega$  is a more complex case. By applying a purely **syntactic unification**, it is not unifiable, since after applying the substitution the formulas are not the same.

However, by applying a **semantic unification**, the set is unifiable, since the formulas after applying the substitution are logically equivalent. We will allow semantic unification using associativity and commutativity of the alternative.

# Most general unifier (mgu)

The **most general unifier (mgu)** of a unifiable set of formulas is the simplest (minimal) unifier for that set.

For a unifiable set of formulas there always exists its mgu, and any other unifier for this set can be obtained by composing the mgu with some additional substitution. The **unification algorithm** computes the mgu of a set of formulas.

## Short review

For the following set of clauses answer if each set is unifiable.

If so, then write its unifier. Try to give both the mgu, and another unifier, which is not mgu. If the set is not unifiable, then explain why.

Pay attention to commas, to correctly identify formulas in sets.

1.  $\{P(x) , P(f(x))\}$

2.  $\{P(x, y) , P(y, x)\}$

3.  $\{P(x, y) , P(y, f(x))\}$

4.  $\{P(x, y) , P(y, f(y))\}$

5.  $\{P(x, y) , P(y, z) , P(z, A)\}$

# Resolution — the general case

Resolution in the general case: if for two clauses (sets of literals):  $\{L_i\}$  and  $\{M_i\}$  there exist respective subsets of literals:  $\{l_i\}$  and  $\{m_i\}$ , called the **collision literals** such, that the set:  $\{l_i\} \cup \{\neg m_i\}$  is unifiable and  $s$  is its mgu, then their resolvent is the set:  $[\{L_i\} - \{l_i\}]s \cup [\{M_i\} - \{m_i\}]s$ .

There can exist different resolvents for given clauses, by different selection of collision literals. For example, consider the following clauses:

$$P[x, f(A)] \vee P[x, f(y)] \vee Q(y) \quad \text{and} \quad \neg P[z, f(A)] \vee \neg Q(z)$$

By choosing  $\{l_i\} = \{P[x, f(A)]\}$  and  $\{m_i\} = \{\neg P[z, f(A)]\}$  we obtain the resolvent:

$$P[z, f(y)] \vee \neg Q(z) \vee Q(y)$$

But by choosing  $\{l_i\} = \{P[x, f(A)], P[x, f(y)]\}$  and  $\{m_i\} = \{\neg P[z, f(A)]\}$  we obtain:

$$Q(A) \vee \neg Q(z)$$

# Short review

For the following set of clauses, write all possible to obtain resolvents.

For each resolvent, note which clauses it was derived from, and what substitution was used. If it is not possible to compute a resolution, then give a short explanation.

Pay attention to commas, to correctly identify formulas in sets.

1.  $\{\neg P(x) \vee Q(x) , P(A)\}$
2.  $\{\neg P(x) \vee Q(x) , \neg Q(x)\}$
3.  $\{\neg P(x) \vee Q(x) , P(f(x)) , \neg Q(x)\}$

# Resolution as an inference rule

Resolution is a sound inference rule, since a clause obtained from a pair of clauses by resolution is their logical consequence. It is, however, not complete, i.e. we cannot derive by resolution just any conclusion  $\varphi$  of a given formula  $\Delta$ , such that  $\Delta \vdash \varphi$ .

For example, for  $\Delta = \{P, Q\}$  we cannot derive by resolution the formula  $P \vee Q$  or  $P \wedge Q$ , and for  $\Delta = \{\forall x R(x)\}$  we cannot derive the formula  $\exists x R(x)$ .

However, for an unsatisfiable set of clauses, it is always possible with resolution to derive a null clause (denoted by  $\square$ ), representing contradiction. So it is possible to use the resolution in the **refutation** proof procedure, by negating the thesis and trying to derive the empty clause. Using this procedure, the validity of any theorem can be verified. Therefore, resolution is said to be **refutation complete**.

Consider the above examples. For:  $\Delta = \{P, Q\}$  negating the formula  $P \vee Q$  gives the clauses  $\neg P$  and  $\neg Q$  and each of them immediately gives the empty clause with the corresponding clause from  $\Delta$ . The negation of  $P \wedge Q$  is the clause  $\neg P \vee \neg Q$  and the empty clause can be derived in two resolution steps. For  $\Delta = \{\forall x R(x)\}$  the negation of  $\exists x R(x)$  is  $\neg R(y)$ , which unifies with the clause  $R(x)$  derived from  $\Delta$  and derives the empty clause in one resolution step.

# Theorem proving based on resolution

The basic reasoning scheme based on resolution, when we have a set of axioms  $\Delta$  and want to derive from it the formula  $\varphi$ , is the following. We make a union of the sets of clauses obtained from  $\Delta$  and  $\neg\varphi$ , and we try to derive falsity (the empty clause) from it, generating subsequent resolvents from the selected pairs of clauses. At each step we add the newly obtained resolvent to the main set of clauses, and repeat the procedure.

The main result from the mathematical logic being used here is the following two facts. If resolution is executed on a set of clauses obtained from an unsatisfiable formula, with some systematic algorithm of generating resolvents, then we will obtain the empty clause at some point. And the other way around, if the empty clause can be generated from a set of clauses obtained from some formula, then this set of clauses, but also the original formula, are both unsatisfiable. This applies as well to the clauses created by skolemization, so is a confirmation of the correctness of the whole procedure.



# Theorem proving: an example

We know that:

- |                                   |   |
|-----------------------------------|---|
| 1. Whoever can read is literate.  | $(\forall x)[R(x) \Rightarrow L(x)]$      |
| 2. Dolphins are not literate.     | $(\forall x)[D(x) \Rightarrow \neg L(x)]$ |
| 3. Some dolphins are intelligent. | $(\exists x)[D(x) \wedge I(x)]$           |

We need to prove the statement:

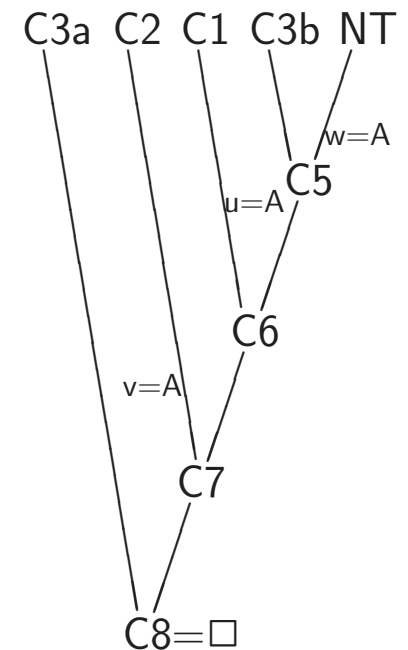
4. Some who are intelligent cannot read.  $(\exists x)[I(x) \wedge \neg R(x)]$

After converting the statements to the prenex CNF form we obtain the clauses:

- |      |                            |                                  |
|------|----------------------------|----------------------------------|
| C1:  | $\neg R(u) \vee L(u)$      | from the first axiom             |
| C2:  | $\neg D(v) \vee \neg L(v)$ | from the second axiom            |
| C3a: | $D(A)$                     | from the third axiom, p.1        |
| C3b: | $I(A)$                     | from the third axiom, p.2        |
| NT:  | $\neg I(w) \vee R(w)$      | from the negation of the theorem |

From the subsequent resolution steps we obtain:

- |     |             |                                 |
|-----|-------------|---------------------------------|
| C5: | $R(A)$      | resolvent of clauses C3b and NT |
| C6: | $L(A)$      | resolvent of clauses C5 and C1  |
| C7: | $\neg D(A)$ | resolvent of clauses C6 and C2  |
| C8: | $\square$   | resolvent of clauses C7 and C3a |





# Theorem proving: an example from mathematics

Let us consider the following example from mathematics.<sup>2</sup> We would like to prove that the intersection of two sets is contained in either one of them. We start from writing the axioms which are necessary for this reasoning. In this case they are the definitions of the set intersection and inclusion.

$$\begin{aligned}\forall x \forall s \forall t \quad (x \in s \wedge x \in t) &\Leftrightarrow x \in s \cap t \\ \forall s \forall t \quad (\forall x \, x \in s \Rightarrow x \in t) &\Leftrightarrow s \subseteq t\end{aligned}$$

Theorem to be proved:

$$\forall s \forall t \quad s \cap t \subseteq s$$

---

<sup>2</sup>This example is borrowed from the book “Logical Foundations of Artificial Intelligence” by Genesereth and Nilsson.

After converting the formulas to prenex CNF sets of clauses:

1.  $\{x \notin s, x \notin t, x \in s \cap t\}$  from the definition of intersection
2.  $\{x \notin s \cap t, x \in s\}$  from the definition of intersection
3.  $\{x \notin s \cap t, x \in t\}$  from the definition of intersection
4.  $\{F(s, t) \in s, s \subseteq t\}$  from the definition of inclusion
5.  $\{F(s, t) \notin t, s \subseteq t\}$  from the definition of inclusion
6.  $\{A \cap B \not\subseteq A\}$  from the negation of the theorem

Let us note the Skolem functions in clauses 4 and 5, and Skolem constants in clause 6. Below is the proof sequence leading quite directly to the empty clause.

7.  $\{F(A \cap B, A) \in A \cap B\}$  from clauses 4. and 6.
8.  $\{F(A \cap B, A) \notin A\}$  from clauses 5. and 6.
9.  $\{F(A \cap B, A) \in A\}$  from clauses 2. and 7.
10.  $\{\}$  from clauses 8. and 9.

This is it. Theorem proved. But somehow, it is hard to feel the satisfaction which typically accompanies completing a real mathematical proof. Furthermore, in case one wanted to review the proof, or verify it, one has to do some nontrivial extra work, although in this particular case it is still relatively simple.

# Short review

For the following axiom sets  $\Delta$  and theorems  $\varphi$ , try proving  $\Delta \vdash \varphi$  using resolution refutation.

1.  $\Delta = \{\forall x(\text{Likes}(x, \text{Wine}) \Rightarrow \text{Likes}(\text{Dick}, x)), \text{Likes}(\text{Ed}, \text{Wine})\}$   
 $\varphi = \text{Likes}(\text{Dick}, \text{Ed})$
2.  $\Delta = \{\forall x(\text{Likes}(x, \text{Dick}) \Rightarrow \text{Likes}(\text{Dick}, x)), \neg \text{Likes}(\text{wife}(\text{Ed}), \text{Dick})\}$   
 $\varphi = \text{Likes}(\text{Dick}, \text{wife}(\text{Ed}))$
3.  $\Delta = \{\forall x(\text{Likes}(x, \text{Wine}) \Rightarrow \text{Likes}(\text{Dick}, x)), \text{Likes}(\text{Ed}, \text{Wine})\}$   
 $\varphi = (\text{Likes}(\text{Dick}, \text{Ed}) \vee \text{Likes}(\text{Dick}, \text{wife}(\text{Ed})))$
4.  $\Delta = \{\forall x(\text{Likes}(x, \text{Wine}) \Rightarrow \text{Likes}(\text{Dick}, x)), \text{Likes}(\text{Ed}, \text{Wine})\}$   
 $\varphi = (\text{Likes}(\text{Dick}, \text{Ed}) \wedge \text{Likes}(\text{Dick}, \text{wife}(\text{Ed})))$



# Knowledge engineering

The presented formalism of first order predicate logic, along with resolution as a method of theorem proving, constitute a technique for building intelligent agents capable of solving problems presented to them. The construction of such an agent, however, requires an efficient design of a representation, which can be formulated as the following process termed **knowledge engineering**:

## **problem identification**

Define the scope of questions which the agent would have to be able to answer, the type of facts, which she will be able to use, etc. For example, in relation to the wumpus world, we must declare whether the agent should be able to plan activities, or, for example, only create the representation of the state of the world identified by previous actions.

## **knowledge acquisition**

The developer of the agent software (knowledge engineer) may not understand all the nuances of the described world, and must cooperate with experts to obtain all the necessary knowledge.

## **definition of the representation dictionary**

The concepts and objects of the problem domain must be described with logical formulas. It is necessary to define a vocabulary of predicates and terms, i.e. term functions and constants. This stage may prove crucial for the ability to effectively solve problems, e.g. in the wumpus world, would pits better be represented as objects, or properties of locations.

## **encoding general knowledge**

Encode axioms containing general knowledge about the problem domain, the rules governing this world, existing heuristics, etc.

## **encoding specific knowledge**

The statement of a specific problem to be solved by the agent, including the facts about all specific objects, as well as the question to answer, or, more generally, the theorem to prove.

## **submit queries to the reasoning device**

Run the theorem proving procedure on the knowledge base constructed.



## **debug the knowledge base**

Unfortunately, as is also the case with normal programs, rarely designed system will immediately work properly. There may occur such problems as the lack of some key axioms, or axioms imprecisely stated, that permit proving too strong assertions.



# Dealing with equality

One very special relation occurring in logical formulas is the equality (identity) of terms.



Example:

$\Delta = \{=(\text{wife}(\text{Ed}), \text{Meg}), \text{Owns}(\text{wife}(\text{Ed}), \text{alfa-8c})\}.$

Does it mean that Meg owns an Alfa 8C Competizione?

Can we prove it using resolution?

$\text{Owns}(\text{Meg}, \text{alfa-8c})?$

Unfortunately not. The resolution proof procedure does not treat the equality predicate in any special way, and will not take advantage of the term equality information it has. For the resolution proof in the above example to succeed, we would have to formulate an appropriate equality axiom:

$$\forall x, y, z [\text{Owns}(x, y) \wedge =(x, z) \Rightarrow \text{Owns}(z, y)]$$

Using the above axiom, connecting Owns with equality, we can prove that Meg owns the Alfa, as well as any other facts and equalities of the owners. However, to likewise extend the reasoning to the equality of the objects owned, an additional axiom has to be introduced:

$$\forall x, y, z [\text{Owns}(x, y) \wedge =(y, z) \Rightarrow \text{Owns}(x, z)]$$

Worse yet, for the system to properly handle all the facts of the term identities with respect to all relations, similar axioms would have to be written for all the predicate symbols. Unfortunately, in the first order predicate language it is not possible to express this in one general formula, like:

$$\forall P, y, z [P(y) \wedge =(y, z) \Rightarrow P(z)]$$

An alternative solution would be to incorporate the processing of term equality in the theorem proving process. Several solutions exist: a formula reduction rule with respect to term equality, called **demodulation**, a generalized resolution rule called **paramodulation**, and an extended unification procedure which handles equalities.

# Answer extraction from the proof tree

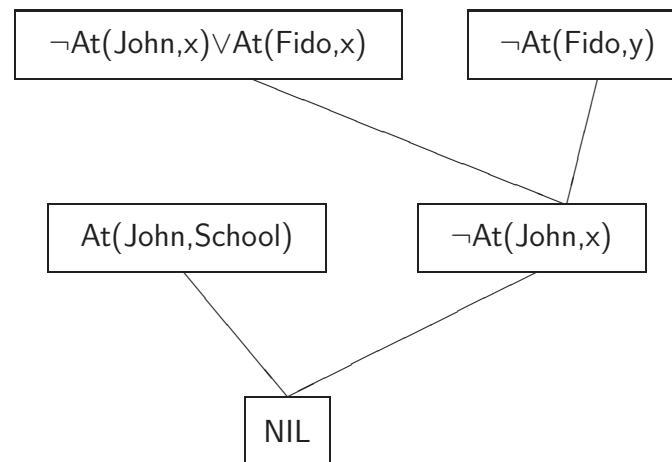
Consider a simple example, we know:

1. Where is John, there is Fido.  $(\forall x)[\text{At}(\text{John}, x) \Rightarrow \text{At}(\text{Fido}, x)]$
2. John is at school.  $\text{At}(\text{John}, \text{School})$

The question we need to answer is:

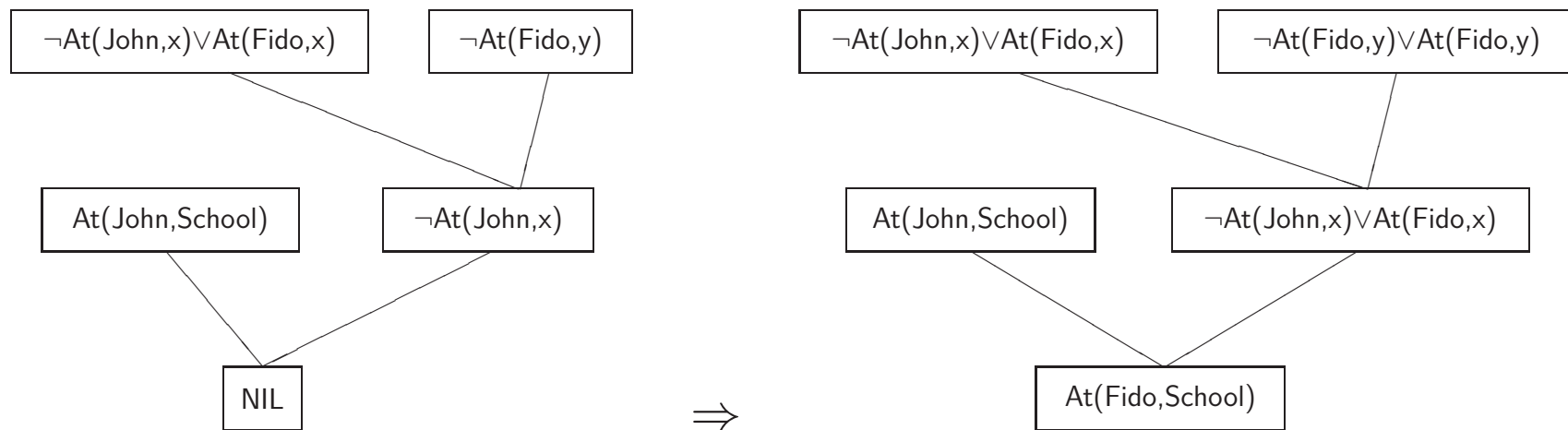
3. Where is Fido?  $(\exists x)[\text{At}(\text{Fido}, x)]$

The logical version of the original question is  $\neg \text{At}(\text{Fido}, x)$ , and the proof is generated easily:



Unfortunately, it does not provide an answer to the original question.

## Answer extraction from the proof tree (contd.)



- The basic procedure converts the refutation proof to a direct proof of the theorem.
- If the theorem contains alternatives (which become conjunctions after negation) then the resulting formula may be complex and hard to interpret.
- If the theorem contains a universal quantifier then after negation it contains Skolem functions or constants, which are carried to the result formula, but can be converted to a universally quantified variable.

# Resolution speedup strategies

In proving theorems using the resolution refutation procedure, we aim to generate the empty clause, indicating a contradiction. To be sure that the empty clause obtains, assuming that this is at all possible, we need to generate the resolvents in some systematic way, for example, using the breadth-first search. But with larger databases, this can lead to generating a large number of conclusions, of which most may not have anything to do with the theorem being proved.

It would be useful to have some **speedup strategies**, which would cut off this search and prevent generating at least some resolvents. These strategies can be **complete**, i.e. such that will always find a solutions (contradiction) if at all possible, or **incomplete** i.e. giving no such a guarantee (but typically much more effective).

## Speedup strategies:

- single literal preference (by itself incomplete, but complete if used as a preference)
- set of support: allow only resolution steps using one clause from a certain set, initially equal to the negated theorem clauses (complete)
- input resolution permits only resolution steps using an original input clause (complete only in some cases, e.g. for Horn clause databases)
- linear resolution (incomplete)
- repetition and subsumption reduction (complete)



# Undecidability of predicate calculus

The predicate calculus seems well suited to expressing facts and reasoning in artificial intelligence systems. However, we need to be aware of some of its fundamental limitations, which constrain its practical applications.

Church's theorem (1936, of the undecidability of predicate calculus): there does not exist a **decision procedure**, which could test the validity of any predicate calculus formula. We say that the predicate calculus is **undecidable**.

This property significantly restricts what can be inferred in the predicate calculus. However, for a number of classes of formulas, there does exist a decision procedure. Furthermore, the predicate calculus has the property of **semidecidability**, which means that there exists a procedure which can determine, in a finite number of steps, that a formula is unsatisfiable, if it is so. However, for satisfiable formulas such procedure may not terminate.

# Incompleteness in predicate calculus

One could think, that the undecidability of predicate calculus can be circumvented by taking advantage of its semidecidability. Attempting to derive a formula  $\varphi$  from an axiom set  $\Delta$ , we start two simultaneous proof procedures:  $\Delta \vdash \varphi$  and  $\Delta \vdash \neg\varphi$ . By semidecidability, we could expect that at least one of these proofs should terminate. Unfortunately, this is not so.

Gödel's theorem (1931, of incompleteness): in predicate calculus one can formulate **incomplete** theories: theories with (closed) formulas, which cannot be proved true or false. What's more, such theories are quite simple and common, e.g. the theory of natural numbers is one such theory.

A theory  $\mathcal{T}$  is called **decidable** if there exists an algorithm that for any closed formula  $\varphi$  can test whether  $\varphi \in \mathcal{T}$ , or  $\varphi \notin \mathcal{T}$ . Incomplete theories are obviously undecidable.

The effect of the Gödel's theorem is, that if, after some number of steps of a proof  $\Delta \vdash \varphi$  (and, perhaps a simultaneous proof  $\Delta \vdash \neg\varphi$ ), there is still no derivation, then we still cannot be certain whether the proof will eventually terminate (or at least one of them will), or that we have an incomplete theory.

# Representing changes

The predicate calculus works well as a representation language for static domains, i.e. such where nothing ever changes, and whatever is true, stays so forever. Unfortunately, the real world is not like this.

For example, if the formula:  $At(John, School)$  correctly describes the current state of the morning of some weekday, then, unfortunately, we must accept that John will go home eventually. If the axioms correctly describe the effects of agents' actions, then the system might be able to derive a new fact:  $At(John, Home)$ . Unfortunately, the fact database will then contain a contradiction, which for a logical system is a disaster. A proof system containing a false formula among its axioms can prove any theorem!

Exercise: assume the axiom set  $\Delta$  contains, among others, two clauses:  $P$  and  $\neg P$ . Give a proof of an arbitrary formula  $Q$ . Hint: prove first that  $P \vee \neg P \vee Q$  is a tautology (sentence always true) for any  $P$  and  $Q$ . This can be accomplished constructing a proof:  $\models (P \vee \neg P \vee Q)$ , this is, proving the thesis with an empty set of axioms. Next add such a tautology to the  $\Delta$  set and proceed to prove  $Q$ .

# Temporal logics

To solve the problem of representation of the changes, a number of special logic theories, called **temporal logics** were created. Ordinary facts expressed in these logics occur at specific points in time. However, the time, its properties, and special inference rules concerning its passage, are built into the theory (instead of being represented explicitly, along with other properties of the world).

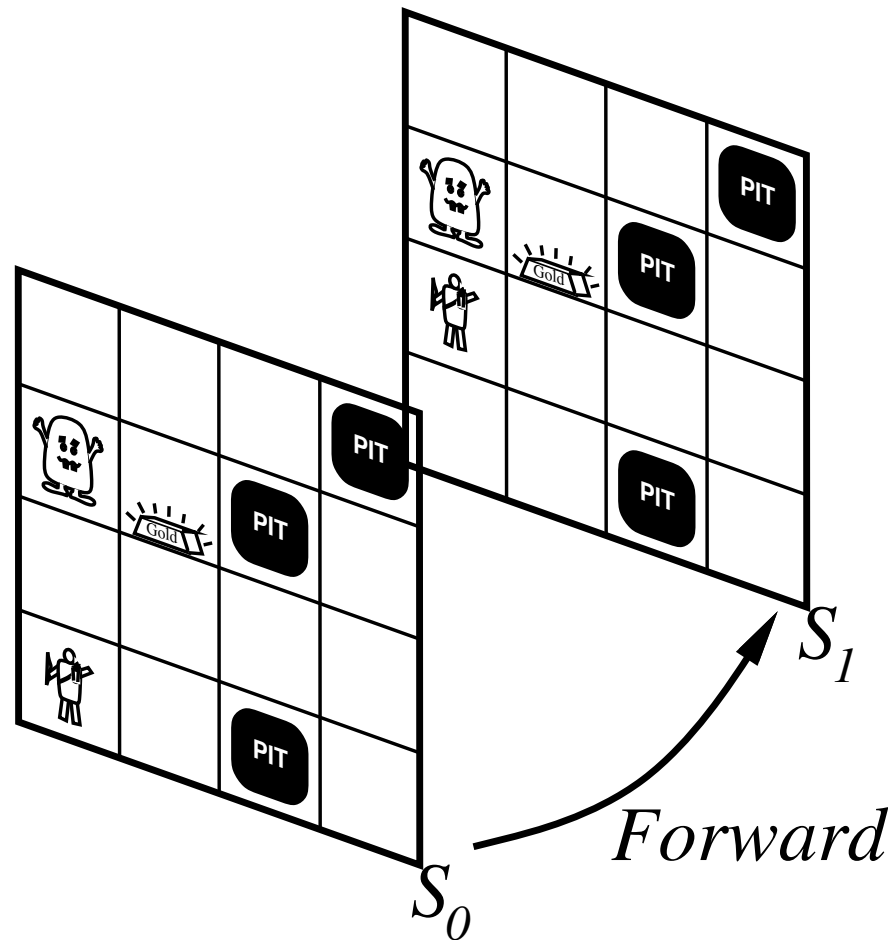
One of the main issues, which these theories treat differently, is the representation of the time itself. Time can be discrete or continuous, may be given in the form of points or intervals, may be bounded or unbounded, etc. Moreover, time can be a linear concept or can have branches. Usually it should be structured, although there exist circular representations of time.

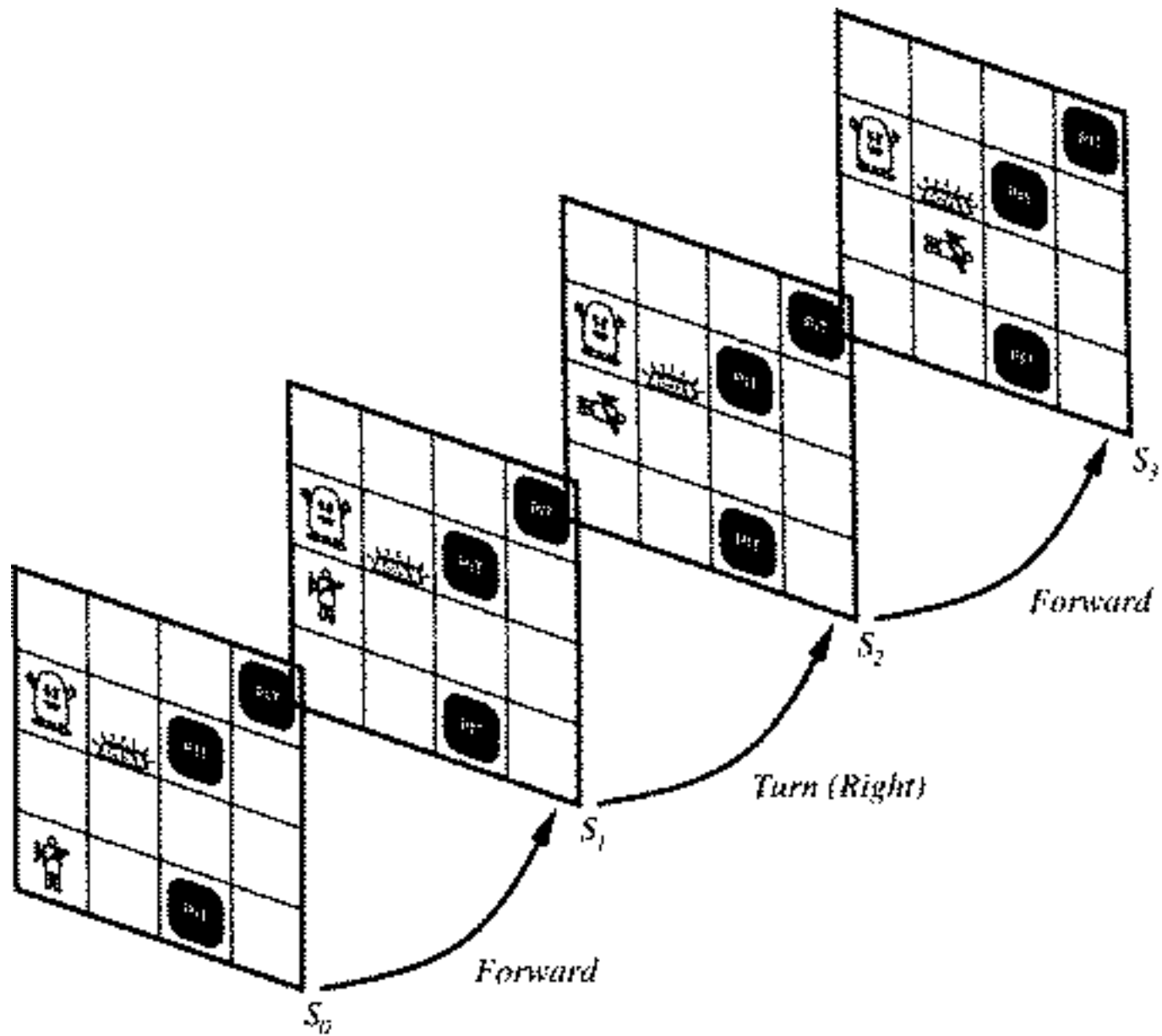
For each of these temporal logics, to be able to effectively reason about formulas created, which represent the phenomena that an intelligent agent faces, there must exist a proof procedure. The construction of such a procedure may be based on projections of the given theory to first-order predicate logic.

# The situation calculus

An alternative to temporal logics is a direct recording of time moments in the representation language. An example of such an approach is the **situation calculus**:

$$At(Agent, [1, 1], S_0) \wedge At(Agent, [1, 2], S_1) \wedge S_1 = Result(Forward, S_0)$$





# The situation calculus (contd.)

The situation calculus uses the concepts of: **situations**, **actions**, and **fluents**:

**situations:** a situation is the initial state  $s_0$ , and for any situation  $s$  and action  $a$  a situation is also  $Result(a, s)$ ; situations correspond to sequences of actions, and are thus different from states, i.e. an agent may be in some state through different situations,

**fluents:** functions and relations which can vary from one situation to the next are called fluents; by convention their last argument is the situation argument,

**possibility axioms:** describe preconditions of actions, e.g. for action *Shoot*:  
 $Have(Agent, Arrow, s) \Rightarrow Poss(Shoot, s)$

**successor-state axioms:** describe for each fluent what happens depending on the action taken, e.g. for action *Grab* the axiom should assert, that after properly executing the action the agent will end up holding whatever she grabbed; but we must also remember about situations when the fluent was unaffected by some action:

$$\begin{aligned} Poss(a, s) \Rightarrow \\ (Holding(Agent, g, Result(a, s)) \Leftrightarrow \\ a = Grab(g) \vee (Holding(Agent, g, s) \wedge a \neq Release(g))). \end{aligned}$$

**unique action axioms:** because of the presence of the action inequality clauses on the successor-state axioms, we must enable to agent to effectively derive such facts, by adding the unique action axioms; for each pair of action symbols  $A_i$  and  $A_j$  we must state the (seemingly obvious) axiom  $A_i \neq A_j$ ; also, for actions with parameters we must also state:

$$A_i(x_1, \dots, x_n) = A_j(y_1, \dots, y_n) \Leftrightarrow x_1 = y_1 \wedge \dots \wedge x_n = y_n$$



# Example: monkey and bananas — the axioms

- general knowledge of the world and operators (partial and simplified):

$$\begin{aligned} A1: & \forall p \forall p_1 \forall s [\text{At}(\text{BOX}, p, s) \Rightarrow \text{At}(\text{BOX}, p, \text{goto}(p_1, s))] \\ A2: & \forall p \forall p_1 \forall s [\text{At}(\text{BANANAS}, p, s) \Rightarrow \text{At}(\text{BANANAS}, p, \text{goto}(p_1, s))] \\ A3: & \forall p \forall s [\text{At}(\text{MONKEY}, p, \text{goto}(p, s))] \\ A4: & \forall p \forall p_1 \forall s [\text{At}(\text{BOX}, p, s) \wedge \text{At}(\text{MONKEY}, p, s) \Rightarrow \text{At}(\text{BOX}, p_1, \text{move}(\text{BOX}, p, p_1, s))] \\ A5: & \forall p \forall p_1 \forall p_2 \forall s [\text{At}(\text{BANANAS}, p, s) \Rightarrow \text{At}(\text{BANANAS}, p, \text{move}(\text{BOX}, p_1, p_2, s))] \\ A6: & \forall p \forall p_1 \forall s [\text{At}(\text{MONKEY}, p, s) \Rightarrow \text{At}(\text{MONKEY}, p_1, \text{move}(\text{BOX}, p, p_1, s))] \\ A7: & \forall s [\text{Under}(\text{BOX}, \text{BANANAS}, s) \Rightarrow \text{Under}(\text{BOX}, \text{BANANAS}, \text{climb}(\text{BOX}, s))] \\ A8: & \forall p \forall s [\text{At}(\text{BOX}, p, s) \wedge \text{At}(\text{MONKEY}, p, s) \Rightarrow \text{On}(\text{MONKEY}, \text{BOX}, \text{climb}(\text{BOX}, s))] \\ A9: & \forall s [\text{Under}(\text{BOX}, \text{BANANAS}, s) \wedge \text{On}(\text{MONKEY}, \text{BOX}, s) \Rightarrow \text{Havebananas}(\text{grab}(\text{BANANAS}, s))] \\ A10: & \forall p \forall s [\text{At}(\text{BOX}, p, s) \wedge \text{At}(\text{BANANAS}, p, s) \Rightarrow \text{Under}(\text{BOX}, \text{BANANAS}, s)] \end{aligned}$$

- specific case data:

$$A11: [\text{At}(\text{MONKEY}, P_1, S_0) \wedge \text{At}(\text{BOX}, P_2, S_0) \wedge \text{At}(\text{BANANAS}, P_3, S_0)]$$

- theorem to prove:

$$\exists s (\text{Havebananas}(s))$$

---

<sup>2</sup>The solution to the monkey and bananas problem presented here is based on an example in book „Artificial Intelligence“ by Philip C. Jackson Jr.

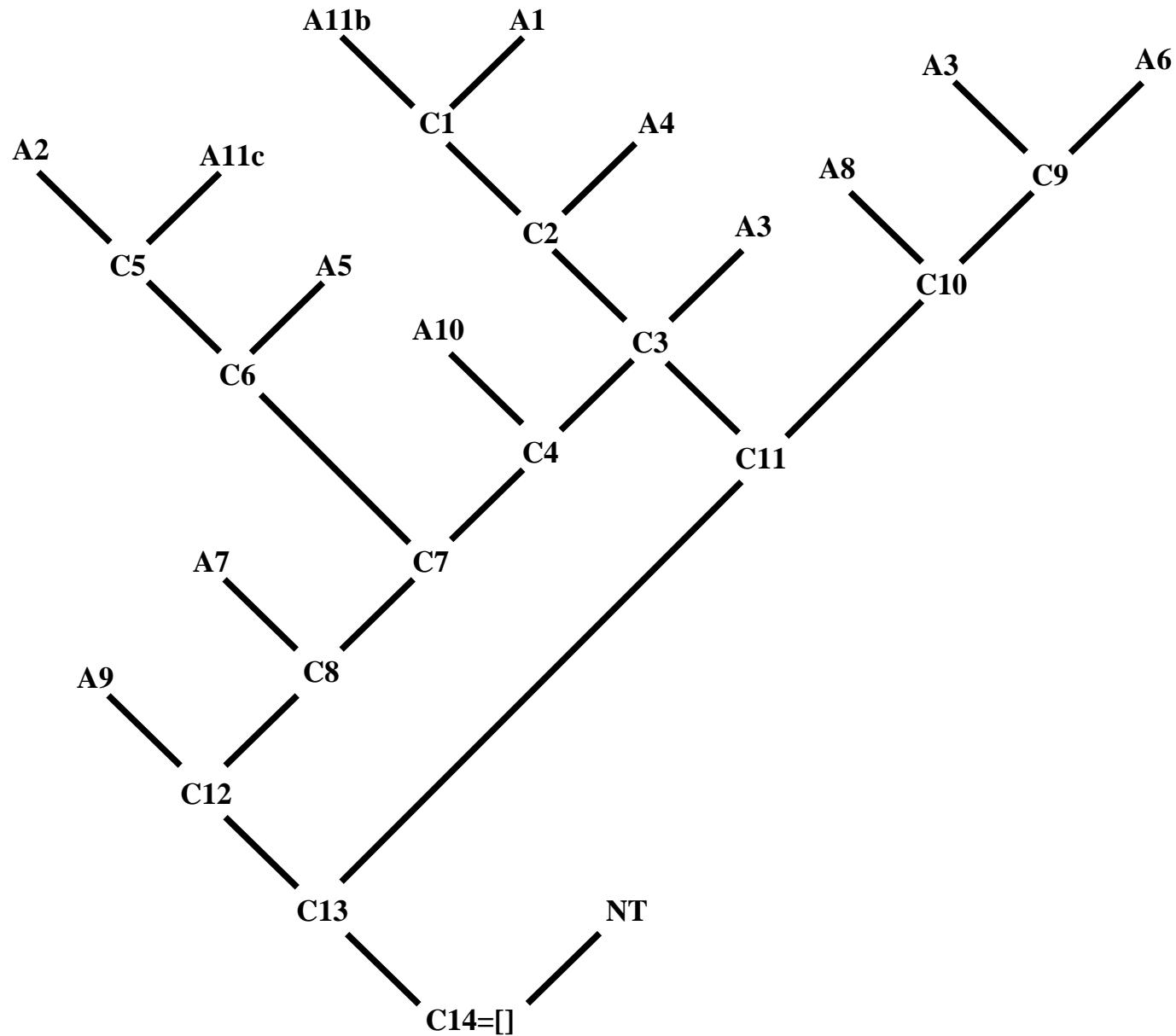
# Example: monkey and bananas — the clauses

- A1:  $\{\neg \text{At}(\text{BOX}, p, s_1), \text{At}(\text{BOX}, p, \text{goto}(p_1, s_1))\}$   
A2:  $\{\neg \text{At}(\text{BANANAS}, q, s_2), \text{At}(\text{BANANAS}, q, \text{goto}(q_1, s_2))\}$   
A3:  $\{\text{At}(\text{MONKEY}, r, \text{goto}(r, s_3))\}$   
A4:  $\{\neg \text{At}(\text{BOX}, u, s_4), \neg \text{At}(\text{MONKEY}, u, s_4), \text{At}(\text{BOX}, u_1, \text{move}(\text{BOX}, u, u_1, s_4))\}$   
A5:  $\{\neg \text{At}(\text{BANANAS}, t, s_5), \text{At}(\text{BANANAS}, t, \text{move}(\text{BOX}, t_2, t_3, s_5))\}$   
A6:  $\{\neg \text{At}(\text{MONKEY}, v_1, s_6), \text{At}(\text{MONKEY}, v_2, \text{move}(\text{BOX}, v_1, v_2, s_6))\}$   
A7:  $\{\neg \text{Under}(\text{BOX}, \text{BANANAS}, s_7), \text{Under}(\text{BOX}, \text{BANANAS}, \text{climb}(\text{BOX}, s_7))\}$   
A8:  $\{\neg \text{At}(\text{MONKEY}, w, s_8), \neg \text{At}(\text{BOX}, w, s_8), \text{On}(\text{MONKEY}, \text{BOX}, \text{climb}(\text{BOX}, s_8))\}$   
A9:  $\{\neg \text{Under}(\text{BOX}, \text{BANANAS}, s_9), \neg \text{On}(\text{MONKEY}, \text{BANANAS}, s_9),$   
Havebananas( $\text{grab}(\text{BANANAS}, s_9)$ )\}  
A10:  $\{\neg \text{At}(\text{BOX}, p, s_{10}), \neg \text{At}(\text{BANANAS}, p, s_{10}), \text{Under}(\text{BOX}, \text{BANANAS}, s_{10})\}$   
A11a:  $\{\text{At}(\text{MONKEY}, P_1, S_0)\}$   
A11b:  $\{\text{At}(\text{BOX}, P_2, S_0)\}$   
A11c:  $\{\text{At}(\text{BANANAS}, P_3, S_0)\}$   
NT:  $\{\neg \text{Havebananas}(z)\}$

# Example: monkey and bananas — the proof

C1(A1,A11b)	$\{ \text{At}(\text{BOX}, P_2, \text{goto}(p_1, S_0)) \}$
C2(C1,A4)	$\{ \neg \text{At}(\text{BANANAS}, P_2, \text{goto}(p_1, S_0)),$ $\text{At}(\text{BOX}, u_1, \text{move}(\text{BOX}, P_2, u_1, \text{goto}(p_1, S_0))) \}$
C3(C2,A3)	$\{ \text{At}(\text{BOX}, u_1, \text{move}(\text{BOX}, P_2, u_1, \text{goto}(P_2, S_0))) \}$
C4(C3,A10)	$\{ \neg \text{At}(\text{BANANAS}, u_1, \text{move}(\text{BOX}, P_2, u_1, \text{goto}(P_2, S_0))),$ $\text{Under}(\text{BOX}, \text{BANANAS}, \text{move}(\text{BOX}, P_2, u_1, \text{goto}(P_2, S_0))) \}$
C5(A2,A11c)	$\{ \text{At}(\text{BANANAS}, P_3, \text{goto}(q_1, S_0)) \}$
C6(C5,A5)	$\{ \text{At}(\text{BANANAS}, P_3, \text{move}(\text{BOX}, t_2, t_3, \text{goto}(q_1, S_0))) \}$
C7(C6,C4)	$\{ \text{Under}(\text{BOX}, \text{BANANAS}, \text{move}(\text{BOX}, P_2, P_3, \text{goto}(P_2, S_0))) \}$
C8(C7,A7)	$\{ \text{Under}(\text{BOX}, \text{BANANAS}, \text{climb}(\text{BOX}, \text{move}(\text{BOX}, P_2, P_3, \text{goto}(P_2, S_0)))) \}$
C9(A3,A6)	$\{ \text{At}(\text{MONKEY}, v_2, \text{move}(\text{BOX}, r, v_2, \text{goto}(r, r_1))) \}$
C10(C9,A8)	$\{ \text{At}(\text{BOX}, v_2, \text{move}(\text{BOX}, r, v_2, \text{goto}(r, r_1))),$ $\text{On}(\text{MONKEY}, \text{BOX}, \text{climb}(\text{BOX}, \text{move}(\text{BOX}, r, r_2, \text{goto}(r, r_1)))) \}$
C11(C10,C3)	$\{ \text{On}(\text{MONKEY}, \text{BOX}, \text{climb}(\text{BOX}, \text{move}(\text{BOX}, P_2, u_1, \text{goto}(P_2, S_0)))) \}$
C12(C8,A9)	$\{ \neg \text{On}(\text{MONKEY}, \text{BOX}, \text{climb}(\text{BOX}, \text{move}(\text{BOX}, P_2, P_3, \text{goto}(P_2, S_0))))),$ $\text{Havebananas}(\text{grab}(\text{BANANAS},$ $\text{climb}(\text{BOX}, \text{move}(\text{BOX}, P_2, P_3, \text{goto}(P_2, S_0)))) \}$
C13(C11,C12)	$\{ \text{Havebananas}(\text{grab}(\text{BANANAS},$ $\text{climb}(\text{BOX}, \text{move}(\text{BOX}, P_2, P_3, \text{goto}(P_2, S_0)))) \}$
C14(C13,NT)	$\{ \}$

# Example: monkey and bananas — the resolution tree



# The frame problem

As we could see in the wumpus and the monkey and bananas examples, a correct logical description of a problem requires explicitly stating the axioms for the effects of actions on the environment, as well as other effect (like rain). It is also necessary to write the axioms to conclude the lack of change:

$$\begin{aligned}\forall a, x, s \text{ Holding}(x, s) \wedge (a \neq \text{Release}) &\Rightarrow \text{Holding}(x, \text{Result}(a, s)) \\ \forall a, x, s \neg \text{Holding}(x, s) \wedge (a \neq \text{Grab}) &\Rightarrow \neg \text{Holding}(x, \text{Result}(a, s))\end{aligned}$$

Unfortunately, in a world more complex than the wumpus world, there will be many fluents, and the description must represent their changes, as well as invariants, both as direct and indirect consequences of the actions.

These axioms, called the **frame axioms**, are hard to state in a general way, and they significantly complicate the representation.

Of course, during the course of work, the agent must state and answer many questions, and prove theorems. The multiplicity of axioms causes a rapid expansion of her database, which slows down further reasoning, and can result in a total paralysis.

# Short review

1. Write a situation calculus based representation for the wumpus world, as described at the beginning of this document.

# Problems with the lack of information

The logic-based methods presented so far assumed that all information necessary to carry out logical reasoning is available to the agent. Unfortunately, this is not a realistic assumption.

**One problem is that of incomplete information.** Agent may not have full information about the problem, allowing him to draw categorical conclusions. She may, however, have partial information, such as:

- “typical” facts,
- “possible” facts,
- “probable” facts,
- exceptions to the generally valid facts.

Having such information is often crucial for making the right decisions. Unfortunately, the classical predicate logic cannot make any use of them.

**Another problem is the uncertainty of information.** An agent may have data from a variety of not fully reliable sources. In the absence of certain information, those unreliable data should be used. She should reason using the best available data, and estimate the reliability of any conclusions obtained this way. Again, classical logic does not provide such tools.

# Common sense reasoning

Consider what information a human knows for sure, making decisions in everyday life. Getting up in the morning, her intention is to go to work. But what if there is a large-scale failure of public transportation? She should, in fact, get up much earlier, and first check whether the buses are running. The day before, she bought products to make breakfast. But can she be sure that her breakfast salad is still in the refrigerator, or if it did not spoil, or perhaps if someone has not sneaked to steal it, etc?

Conclusion: a logically reasoning agent needs 100% certain information to conduct her actions, and sooner or later she will be paralyzed by the perfect correctness of her inference system. In the real world she will never be able to undertake any action, until she has complete and correct information about the surrounding world.

However, in the world full of incomplete and uncertain information, guesses, defaults and exceptions, people are doing quite well. How do they do it?

We have to accept the fact that **people use a way of reasoning different from the fully rigorous mathematical logic**. We call this mechanism the **common sense reasoning**.

**Common sense reasoning allows humans to reach more conclusions than would be possible with the strict mathematical logic. However, some of these conclusions may later turn out to be incorrect when (some of) the previously made assumptions turn out to be false.**



# Nonmonotonic logics

Part of the blame for the problems of inference using classical logic bears its property known as **monotonicity**. In classical logic, the more we know, the more we can deduce using inferencing.

Humans use a different model of inference, much more flexible, utilizing typical facts, default facts, possible facts, and even lack of information. This kind of reasoning seems not to have the monotonicity property.

For example, lacking good information about a situation a human is likely to make assumptions and derive some conclusions. After having acquired more complete information, she might not be able to conduct the same reasoning and work out the same solutions.<sup>3</sup>

Hence, different models of inference, designed to overcome these problems, and following a more flexible reasoning model similar to that of humans, are collectively called **nonmonotonic logics**.

---

<sup>3</sup>The solution produced earlier, in the absence of information, turns out to be wrong now, but perhaps it was better than the lack of any action. But not necessarily.

# Nonmonotonic logics — example

Minsky's challenge: to design a system, which would permit to correctly describe a well-known fact, that the birds can fly.

$$\forall x[\text{BIRD}(x) \rightarrow \text{CANFLY}(x)]$$

In order to accommodate the exceptions, e.g. ostriches, the preceding formula must be modified per case.

$$\forall x[\text{BIRD}(x) \wedge \neg \text{OSTRICH}(x) \rightarrow \text{CANFLY}(x)]$$

But there are more exceptions: birds bathed in spilled crude oil, wingless birds, sick birds, dead birds, painted birds, abstract birds, ...

An idea: we introduce a modal operator **M**:

$$\forall x[\text{BIRD}(x) \wedge \mathbf{M} \text{ CANFLY}(x) \rightarrow \text{CANFLY}(x)]$$

Now the exceptions can be introduced modularly:

$$\forall x[\text{OSTRICH}(x) \rightarrow \neg \text{CANFLY}(x)]$$

For the following set of facts:

$$\Delta = \{\text{BIRD}(\textit{Tweety}), \text{BIRD}(\textit{Sam}), \text{OSTRICH}(\textit{Sam})\}$$

we can deduce:  $\neg \text{CANFLY}(\textit{Sam})$

so it should not be possible to derive:

$$\mathbf{M} \text{ CANFLY}(\textit{Sam}) \quad \text{nor} \quad \text{CANFLY}(\textit{Sam})$$

However, using the normal proof procedure we cannot prove the *Tweety's* ability to fly:

$$\mathbf{M} \text{ CANFLY}(\textit{Tweety}), \text{CANFLY}(\textit{Tweety})$$

In order to do this, a proof procedure is needed, capable of effective (and automatic) proving theorems in predicate logic extended with the modal operator **M**, consistent with the following inference rule:

$$\frac{\text{Not}(\vdash \neg p)}{\mathbf{M} p}$$

# Nonmonotonic logics — what proof procedure?

Leaving aside the restrictions resulting from the reference to the proof procedure in the above definition, such a procedure may be neither effective computationally, nor decidable nor even semidecidable, as are the proof procedures for the predicate logic.

The premise of the above inference rule contains the statement, that some formula is impossible to prove. To start, this may not be possible to determine at all. And to find a positive confirmation of this fact, it will certainly be necessary to carry out global inferencing over the entire database. For how else we could say that something can not be proved.

In contrast, proofs in first-order predicate calculus are local in nature. If, for example, we are lucky to choose the appropriate premises, we can obtain the proof in several steps, even if the data base contains thousands of facts.

# Problems with logic-based methods

The approach to knowledge representation based on first-order logic at one time created much excitement and hope for building powerful and universal systems of artificial intelligence. There are, however, important considerations which significantly limit the practical applications of this approach:

- **combinatorial explosion** of the proof procedure; while there exist speedup strategies, they do not help much; at the same time it is hard to incorporate in the proof procedure any heuristic information available
- **undecidability** and the Gödel's **incompleteness** of the predicate calculus
- reasoning about **changes** — situation calculus, temporal logics
  - reasoning about change exhibits the **frame problem** — besides determining what has changed, it is essential to keep track of what has not
- reasoning with **incomplete** and **uncertain information**, truly challenging for the formal approaches, but seems crucial for the human reasoning
  - taking into account uncertain information leads to **nonmonotonic** reasoning, a property of the human reasoning, while the traditional (mathematical) logic is strictly monotonic

# Applications of the logic-based methods

The above problems with logic-based methods significantly burden their application as a platform for implementing intelligent agents. Nevertheless, the first order predicate language itself is commonly used in artificial intelligence for representing facts.

Still, in some specific applications it is possible to use this methodology, and the above problems are not critical. Some of these applications are:

- computer program synthesis and verification, software engineering
- design and verification of computing hardware, including the VLSI design
- theorem proving in mathematics; which help seek proofs for any postulated theorems, for which efforts failed to find a proof in the traditional way