

```

1  Tell(the initially empty  $KB$ , the background knowledge  $\beta$ )
2  repeat   $\phi \leftarrow$  a sentence describing the current percept  $\pi$ 
3            $Tell(KB, \phi)$ 
4            $\psi \leftarrow$  the question “which action should I do now?” as a sentence
5            $action \leftarrow Ask(KB, \psi)$ 
6           perform this suggested  $action$ 
7            $\chi \leftarrow$  a sentence saying “this  $action$  has been performed”
8            $Tell(KB, \chi)$ 
9  until  $action =$  “terminate yourself” (if ever)

```

Figure 52: A generic knowledge-based agent. (Russell and Norvig, 2003, Figure 7.1)

**Database** (e.g. relational) contains explicit *facts*, and answers queries by *fetching* and combining the required facts into results.

E.g. “list the names of all those *students* who have registered to both *courses* *TEK* and *LAI*”.

**Knowledge base** contains also *sentences* written in some *knowledge representation language*, and answers queries by *inferring* results from these facts and sentences.

- The idea is to give our agent in Figure 52 a (much) more flexible and powerful kind of memory than e.g. a database.
  - We choose the two logics discussed here as our two languages.
  - $Tell(KB, \phi)$  adds this sentence  $\phi$  into the knowledge base  $KB$ .
  - $Ask(KB, \psi)$  asks whether or not this sentence  $\psi$  *follows from* the sentences  $\phi$  in  $KB$ .
- In the generic knowledge-based agent shown as Figure 52, the
- background** knowledge  $\beta$  consists of whatever its designer thinks it must know beforehand. The designer can now express it
- declaratively** using the knowledge representation language instead of
- procedurally** using some programming language.
- question**  $\psi$  posed to  $KB$  depends on how the *goal* of the agent has been included into its design:
- Herein Figure 52 we assume that it has been written explicitly as part of the background knowledge  $\beta$ ; so we have a goal-based agent as in Section 2.5. Then inference can produce the answer.
  - Or it can be given as a separate utility function instead. Then the question  $\psi$  posed to  $KB$  is “What actions *could* I do now?” instead, and then the utility function selects some *action* from the answer.

## 4.1 The Wumpus World

- We shall use the *Wumpus world*, as in Figure 53, as our toy running example on using logic to build an intelligent agent.

- Its PEAS description is:

**Performance measure:**

- +1000 points for picking up the *gold* — this is the goal of the agent
- −1000 points for dying = entering a square containing a *pit* or a live *Wumpus* monster
- −1 point for each action taken, and
- −10 points for using the *arrow* trying to kill the Wumpus — so that the agent should avoid performing unnecessary actions.

**Environment:** A  $4 \times 4$  grid of squares with...

- the agent starting from square  $[1, 1]$  facing right
- the gold in one square
- the initially live Wumpus in one square, from which it never moves
- maybe pits in some squares.

The starting square  $[1, 1]$  has no Wumpus, no pit, and no gold — so the agent neither dies nor succeeds straight away.

**Actuators:** The agent can...

**turn**  $90^\circ$  left or right

**walk** one square forward in the current direction

**grab** an object in this square

**shoot** the single arrow in the current direction, which flies in a straight line until it hits a wall or the Wumpus.

**Sensors:** The agent has 5 TRUE/FALSE sensors which report a...

**stench** when the Wumpus is in an adjacent square — directly, not diagonally

**breeze** when an adjacent square has a pit

**glitter** when the agent perceives the glitter of the gold in the current square

**bump** when the agent walks into an enclosing wall (and then the action had no effect)

**scream** when the arrow hits the Wumpus, killing it.

- Let us consider an example gives as Figure 54 on how the agent might reason and act in this Wumpus world.

**Top left:** Agent **A** starts at  $[1, 1]$  facing right.

- Agent **A** is actually in the example world given as Figure 53, but does not know it.
- The background knowledge  $\beta$  assures agent **A** that he is at  $[1, 1]$  and that it is **OK** = *certainly* not deadly.
- Agent **A** gets the percept “**Stench** = **Breeze** = **Glitter** = **Bump** = **Scream** = FALSE”.
- Agent **A** *infers* from this percept and  $\beta$  that its both neighbouring squares  $[1, 2]$  and  $[2, 1]$  are also **OK**: “If *there* was a **Pit** (**Wumpus**), then *here* would be **Breeze** (**Smell**) — but isn’t, so...”.
- The *KB* enables agent **A** to discover certainties about parts of its environment — even without visiting those parts.

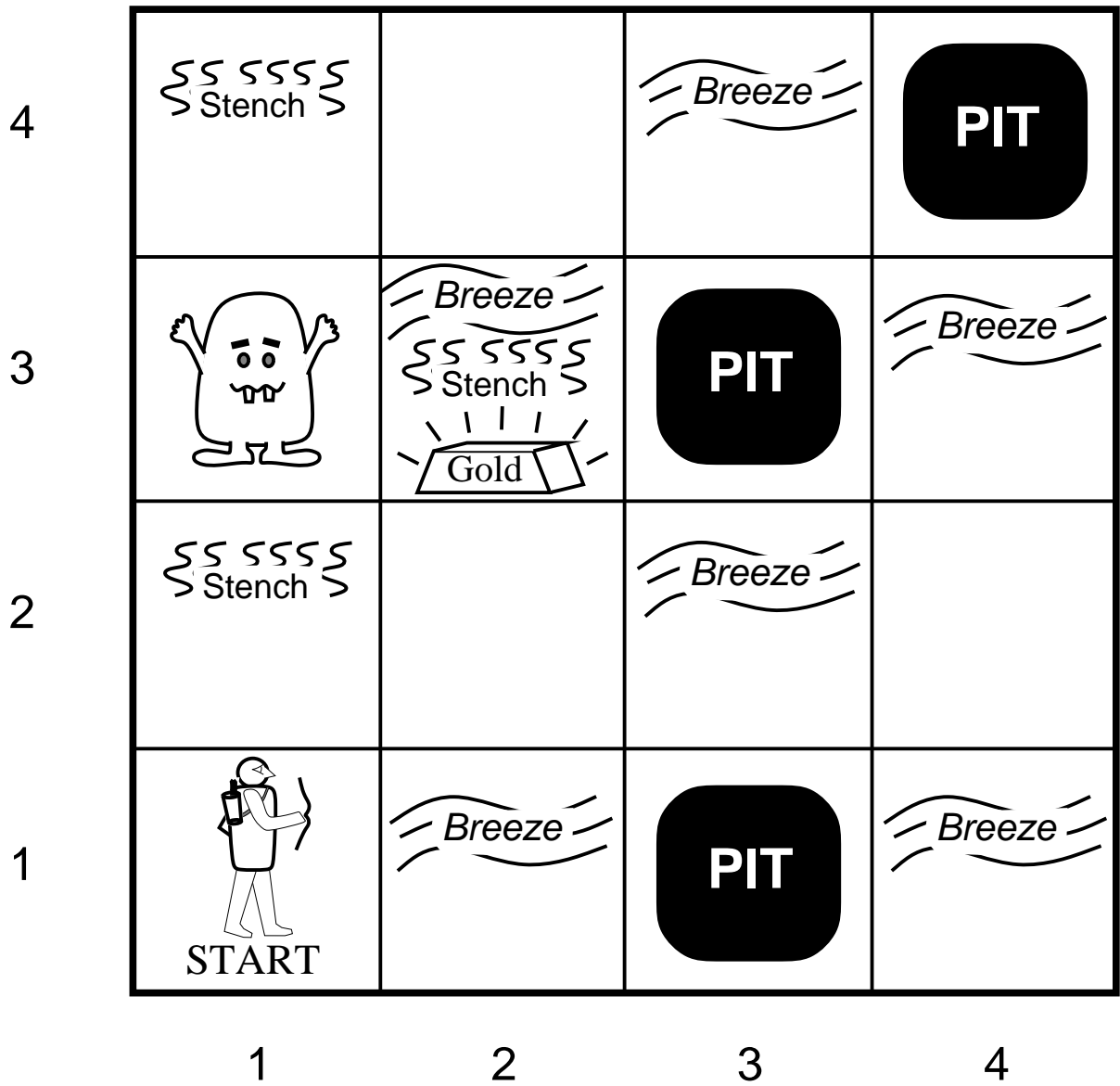


Figure 53: An example Wumpus world. (Russell and Norvig, 2003, Figure 7.2)

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

(a)

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

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

(b)

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

(a)

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

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

(b)

Figure 54: A Wumpus adventure. (Russell and Norvig, 2003, Figures 7.3 and 7.4)

**Top right:** Agent **A** is cautious, and will only move to **OK** squares.

- Agent **A** walks into [2, 1], because it is **OK**, and in the direction where agent **A** is facing, so it is cheaper than the other choice [1, 2]. Agent **A** also marks [1, 1] **Visited**.
- Agent **A** perceives a **Breeze** but nothing else.
- Agent **A** infers: “At least one of the adjacent squares [1, 1], [2, 2] and [3, 1] must contain a **Pit**. There is no **Pit** in [1, 1] by my background knowledge  $\beta$ . Hence [2, 2] or [3, 1] or both must contain a **Pit**.”
- Hence agent **A** cannot be certain of either [2, 2] or [3, 1], so [2, 1] is a dead end for a cautious agent like **A**.

**Bottom left:** Agent **A** has turned back from the dead end [2, 1] and walked to examine the other **OK** choice [1, 2] instead.

- Agent **A** preceives a **Stench** but nothing else.
- Agent **A** infers using also *earlier percepts*: “The **Wumpus** is in an adjacent square. It is not in [1, 1]. It is not in [2, 2] either, because then I would have sensed a **Stench** in [2, 1]. Hence it is in [1, 3].”
- Agent **A** infers using also *earlier inferences*: “There is no **Breeze** here, so there is no **Pit** in any adjacent square. In particular, there is no **Pit** in [2, 2] after all. Hence there is a **Pit** in [3, 1].”
- Agent **A** finally infers: “[2, 2] is **OK** after all — now it is certain that it has neither a **Pit** nor the **Wumpus**.”

This reasoning is too complicated for many animals — but not for the logical agent **A**.

**Bottom right:**

1. Agent **A** walks to the only unvisited **OK** choice [2, 2]. There is no **Breeze** here, and since the the square of the **Wumpus** is now known too, [2, 3] and [3, 2] are **OK** too.
2. Agent **A** walks into [2, 3] and senses the **Glitter** there, so he grabs the gold and succeeds.

## 4.2 Principles behind Logic(s)

- Let us now review the fundamental principles behind logic(s) from our intelligent agent viewpoint.
- From this viewpoint, the actual *syntax* in which the logical *statements* are written as *formulas* is just a “print-out” of the underlying data structure.
  - This underlying data structure is the parse tree of the *formula*.
  - From a programming viewpoint, implementations of logical inference create from existing such structures new ones in ways which respect their *meaning*.
- This meaning or *semantics* of a *statement* defines whether it is TRUE or FALSE in the given *possible world*.
  - In our agent viewpoint, possible world = possible state of the environment.

- Formal semantics for logic(s) uses the term *interpretation* or *model* = a possible world  $w$  + a “dictionary” telling how the different parts of the *statement* are to be understood in  $w$ .

We rarely need to make this distinction.

(The meaning of a modal logic involves not just one but several possible worlds — the different viewpoints.)

- Then a given *statement*  $\phi$  stands for the set of all those possible worlds  $w$  such that  $\phi$  is TRUE in  $w$ .

Formal semantics calls these  $w$  the *models* of  $\phi$  and denotes this set as  $\text{Mod}(\phi)$ .

- E.g.  $\text{Mod}(KB)$  after agent **A** has walked into  $[2, 1]$  and perceived only the **Breeze** there is shown in Figure 55 and consists of all the worlds such that
  - there is **Breeze** in  $[2, 1]$ , and
  - all the possible **Pit** placements around it, according to his background knowledge  $\beta$ , but
  - the other squares of the worlds are not drawn, because agent **A** does not know anything more than  $\beta$  about them (yet).
- Here  $KB$  is considered to be the combined *statement*

$$\beta \text{ and } \alpha_1 \text{ and } \alpha_2 \text{ and } \alpha_3 \text{ and } \dots \text{ and } \alpha_t$$

where these  $\alpha_i$  are the individual *statements* that have been added into  $KB$  using tell *Tell* until now.

Formal semantics often considers the *theory*  $\{\beta, \alpha_1, \alpha_2, \alpha_3, \dots, \alpha_t\}$  instead, to permit also infinite theories.

- The central semantic relation between two *statements* is

$$\gamma \text{ entails } \delta$$

which is written as

$$\gamma \models \delta$$

and means that “**if**  $\gamma$  is TRUE, **then**  $\delta$  must also be TRUE”. For worlds, this is the same thing as

$$\text{Mod}(\gamma) \subseteq \text{Mod}(\delta).$$

E.g.

$$KB \models \alpha_1$$

for the *statement*

$$\alpha_1 = “[1, 2] \text{ has no } \mathbf{Pit}”$$

in Figure 55. Figure 56 shows another example where  $KB$  does *not* entail  $\alpha_2$ .

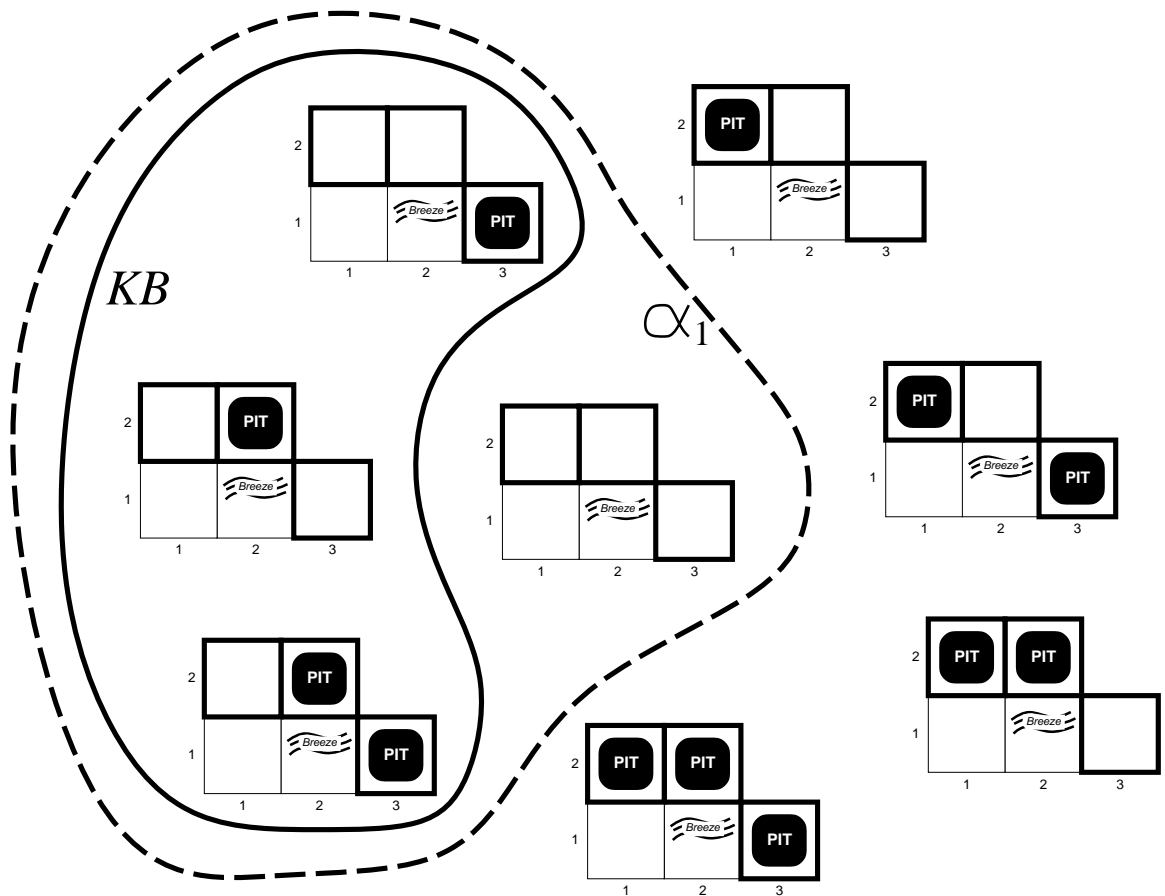


Figure 55: The possible worlds according to  $KB$  and an entailed  $\alpha_1$ . (Russell and Norvig, 2003, Figure 7.5(a))

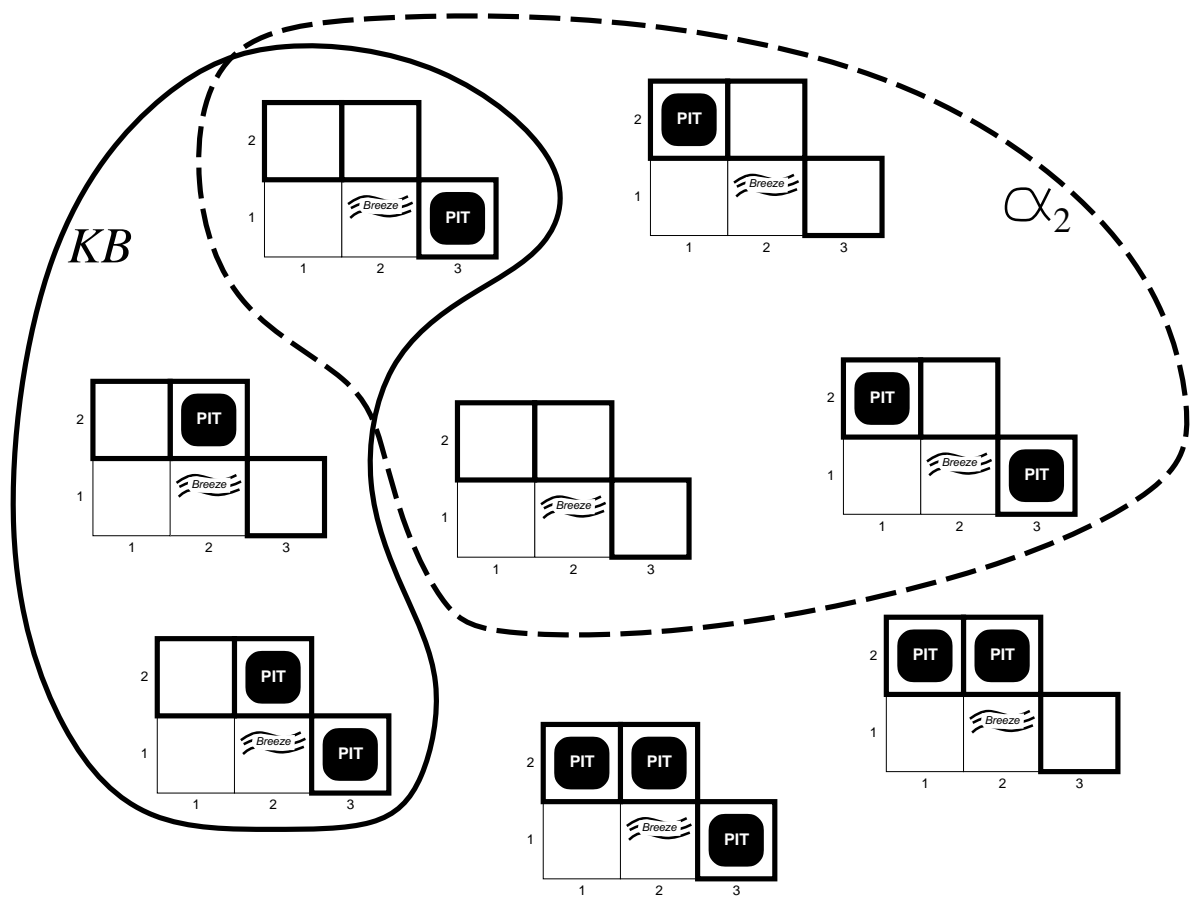


Figure 56:  $KB \not\models \alpha_2$ . (Russell and Norvig, 2003, Figure 7.5(b))



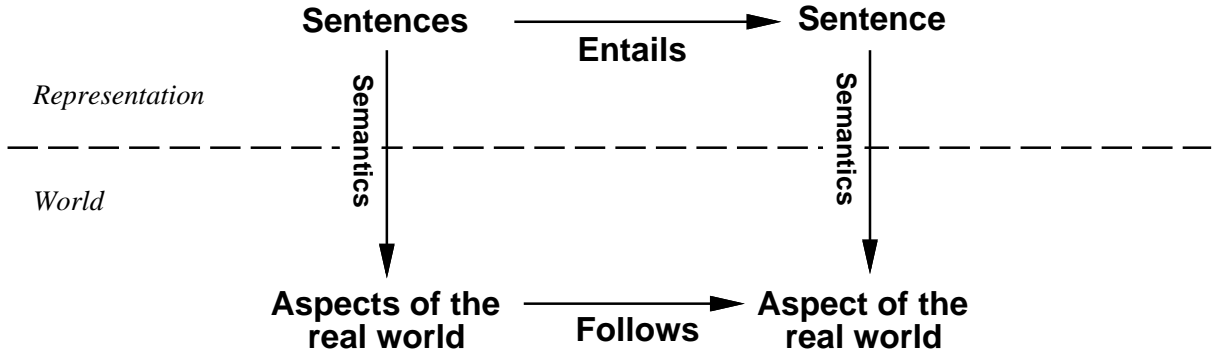


Figure 57: The world and its representation. (Russell and Norvig, 2003, Figure 7.6)

- These sets  $\text{Mod}(\phi)$  are the same *belief* sets  $B_\phi$  as in the agent's belief state space  $\mathcal{B}$ .
  - Conversely,  $\phi$  is one logical representation for this belief set  $B_\phi \in \mathcal{B}$ .
  - The agent believes that the actual world  $w_{\text{actual}}$  is one of the worlds  $w \in \text{Mod}(KB)$ , according to the information  $KB$  it now has about  $w_{\text{actual}}$ .
  - The other worlds  $w' \in \text{Mod}(KB)$  than  $w_{\text{actual}}$  represent the incompleteness of this information, since some things are different in these other  $w'$  than in  $w_{\text{actual}}$ .
  - The philosophical problem of *grounding* is “Why would it be the case that  $w_{\text{actual}} \in \text{Mod}(KB)$ ?”

The agent approach assumes that its sensors create the connection.

- If you think it is strange that the agent's *knowledge* is a *belief* set, then recall the philosopher Plato's (“Platon” in Finnish) (Greece, 429-4347 BCE) dictum

“Knowledge is true justified belief”.

- In summary, as Figure 57 shows:
  - The *sentences* in the  $KB$  represent various (assumed to be) TRUE *aspects* of the real world; The *semantics* explains how.
  - Then  $KB \models \psi$  ensures that the *aspect* represented by  $\psi$  is also TRUE.
- $\text{Ask}(KB, \psi)$  on Step 5 of the knowledge-based agent in Figure 52 is to first **verify** that  $KB \models \psi$  **provide** some appropriate *answer* explaining why.

This agent view is constructive, even when it uses a classical logic: a blunt “because it just is” would not be a sufficient *answer* for action.

- Hence we need a way to *compute*

$$KB \models \psi$$

and an *answer* by manipulating the data structures for  $KB$  and  $\psi$  instead of considering the possible worlds  $w$ .

There are just too many  $w$  to consider — in general, even infinitely many!

- We say “statement  $\psi$  can be *derived* (or *inferred*) from  $KB$  by the *proof system*  $\mathcal{P}$ ”, in symbols

$$KB \vdash_{\mathcal{P}} \psi$$

when this computation can be done using a certain specialized kind of search problem  $\mathcal{P}$  not discussed yet.

- For this search problem  $\mathcal{P}$ ...

**soundness** is the *crucial* logical property that  $\mathcal{P}$  must not lie:

$$\text{if } KB \vdash_{\mathcal{P}} \psi \text{ then } KB \models \psi.$$

**completeness** is the opposite *desirable* logical property that everything which is TRUE should also be derivable by  $\mathcal{P}$ :

$$\text{if } KB \models \psi \text{ then } KB \vdash_{\mathcal{P}} \psi.$$

**effectiveness** is the desirable *computational* property that the search “is  $KB \vdash_{\mathcal{P}} \psi$ ?” always terminates with a YES/NO answer.

- An early (1931) celebrated deep logical result is Kurt Gödel’s *1st incompleteness theorem*:

There cannot be any such proof system  $\mathcal{N}$  with all these three properties for the natural numbers  $\mathbb{N}$ .

## 4.3 Propositional Logic

- The smallest “addressable unit” of propositional logic (“propositio- eli lauselogiikka” in Finnish) is a *whole sentence* which can be either TRUE or FALSE.
- That is, any expression  $\eta$  such that

“Is it the case that  $\eta$ ?”

is a meaningful question with a YES/NO answer.

- E.g. asking this about  $\eta =$

$$\text{Romeo loves Juliet.} \tag{23}$$

is meaningful, so it is a sentence in this sense.

- On the other hand, e.g. neither “Is it the case that **loves**?” nor “Is it the case that **Juliet**?” is a meaningful question.

We shall later present another logic (the *predicate* logic) where we can address also such **parts** of sentences.

- Note also that “Is it the case that **Romeo loves**?” is another meaningful question — but this is a different sentence than Eq. (23).

- The syntax used in this course and its book Russell and Norvig (2003, Figure 7.7) for propositional logic is:

$$\begin{array}{lcl}
\textit{Sentence} \rightarrow \textit{Atomic} & | & \textit{Complex} \\
\textit{Atomic} \rightarrow \text{TRUE} & | & \text{FALSE} & | & \textit{Symbol} \\
\textit{Complex} \rightarrow \neg \textit{Sentence} & & & & \\
& | & (\textit{Sentence} \wedge \textit{Sentence}) \\
& | & (\textit{Sentence} \vee \textit{Sentence}) \\
& | & (\textit{Sentence} \Rightarrow \textit{Sentence}) \\
& | & (\textit{Sentence} \Leftrightarrow \textit{Sentence})
\end{array}$$

- The predicate *Symbols* are variable names which stand for different *sentences*. E.g. we might choose the name  $X$  to stand for Eq. (23) when we are writing the background knowledge  $\beta$ , and so on.
- We suppress writing nested parentheses by stipulating that the *connectives*  $\neg, \wedge, \vee, \Rightarrow, \Leftrightarrow$  are in descending binding power.
- We also stipulate that  $\wedge, \vee$  are associative, but  $\Rightarrow, \Leftrightarrow$  are not.
- In the semantics, a possible world  $w$  is just a function

$$w: \textit{Symbol} \mapsto \underbrace{\{\text{TRUE}, \text{FALSE}\}}_{\text{call this } \mathbb{B}}$$

which gives for each *Symbol*  $Y$  its truth value  $w(Y)$ .

- This  $w$  extends from single *Symbols* to whole *Sentences*:

$$w(c) = c \text{ if } c \in \mathbb{B}$$

$$w(\neg\phi) = \begin{cases} \text{FALSE} & \text{if } w(\phi) = \text{TRUE} \\ \text{TRUE} & \text{if } w(\phi) = \text{FALSE} \end{cases}$$

$$w(\phi \oplus \psi) = \text{look it up from the right line and column from}$$

line		column $\oplus$			
$w(\phi)$	$w(\psi)$	$\wedge$	$\vee$	$\Rightarrow$	$\Leftrightarrow$
FALSE	FALSE	FALSE	FALSE	TRUE	TRUE
FALSE	TRUE	FALSE	TRUE	TRUE	FALSE
TRUE	FALSE	FALSE	TRUE	FALSE	FALSE
TRUE	TRUE	TRUE	TRUE	TRUE	TRUE

These are the *truth tables* for our binary connectives (Russell and Norvig, 2003, Figure 7.8).

- If we think that  $\text{FALSE} < \text{TRUE}$ , then

‘ $\wedge$ ’ is ‘min’.

‘ $\vee$ ’ is ‘max’.

‘ $\Rightarrow$ ’ is ‘ $\leq$ ’.

This  $\phi \Rightarrow \psi$  is *material* implication: the *antecedent*  $\phi$  does not have to be in any way “relevant” or “a cause” for the *consequent*  $\psi$ .

(Philosophers have developed so-called relevance logics which try to capture the idea that  $\phi$  should be somehow relevant to  $\psi$ .)

‘ $\Leftrightarrow$ ’ is ‘=’.

(This connection between logical proofs and calculating with truth values is behind the *calculational proofs* approach, which is handy in Computer Science, e.g. in formal program construction (Backhouse, 2003).)

- The part of the background knowledge  $\beta$  for our Wumpus world which deals with pits can be written in propositional logic as follows:
  - Let the *Symbol*  $P_{i,j}$  (or  $B_{i,j}$ ) stand for “There is a **P**it (or **B**reeze) in square  $[i, j]$ ”.
  - “There is no pit in  $[1, 1]$ ” is then simply

$$\neg P_{1,1}. \quad (24)$$

- We have the general rule that

$$\text{“this square has a Breeze exactly when at least one of its adjacent squares has a Pit”}. \quad (25)$$

- But now we encounter a shortcoming of propositional logic:
  - \* Rule (25) combines together two whole *sentences* with “*exactly when*”. This we *can* do with ‘ $\Leftrightarrow$ ’.
  - \* But rule (25) is trying to combine them so that “this” and “it” are the same square.  
Alas, propositional logic does not allow us to access such *parts* of whole *statements*.
- Since our Wumpus worlds are finite ( $4 \times 4$ -square grids), we can get around this problem with a sledgehammer solution: Write a separate copy of rule (25) for each different “this” square. E.g.

$$\text{corner square: } B_{1,1} \Leftrightarrow P_{1,2} \vee P_{2,1} \quad (26)$$

$$\text{edge square: } B_{2,1} \Leftrightarrow P_{1,1} \vee P_{2,2} \vee P_{3,1} \quad (27)$$

$$\text{other square: } B_{2,2} \Leftrightarrow P_{1,2} \vee P_{2,1} \vee P_{2,3} \vee P_{3,2}. \quad (28)$$

All these *sentences* go into the background knowledge  $\beta$ , and they are TRUE in all Wumpus worlds.

- Then we add into  $KB$  also the Breeze percepts at the top right Wumpus world of Figure 54:

$$\neg B_{1,1} \quad (29)$$

$$B_{2,1}. \quad (30)$$

#### 4.3.1 Propositional Inference

- The aim is to verify that  $KB \models \alpha$  by computing  $KB \vdash_{\mathcal{P}} \alpha$  using some sound proof system  $\mathcal{P}$ .
- For our classical propositional logic, we can even select this  $\mathcal{P}$  to be also complete and effective.
- Here we take this  $\mathcal{P}$  to be some program which the agent is running, and consider some approaches to writing such a program.

- We say that a formula  $\alpha$  is...

**satisfiable** if  $\alpha$  is TRUE in *at least one* possible world  $w$ .

- The computational problem “Is this  $\alpha$  satisfiable?” is called the *SATisfiability* problem.
- For our classical propositional logic, SAT is “the mother of all **NP**-complete problems”.
- As a search problem, it is “try to assign TRUE or FALSE to the *Symbols* in  $\alpha$  in such a way that  $\alpha$  becomes TRUE”.

**valid** if  $\alpha$  is TRUE in *every* possible world  $w$ . Then this  $\alpha$  is...

- TRUE solely because of its own logical form, regardless of how the world(s) happen to be. E.g. the formalization

$$\phi \vee \neg\phi$$

of Eq. (22) is classically valid, regardless of what  $\phi$  “really says” and whether what it says is TRUE or FALSE about the agent’s current environment.

- also called a *tautology*.

- The *deduction theorem* connects entailment and validity:

$$\gamma \models \delta \text{ if and only if } \gamma \Rightarrow \delta \text{ is valid.} \quad (31)$$

- That is, the implication ‘ $\Rightarrow$ ’ is at the *language* level the same concept as entailment ‘ $\models$ ’ is at its *semantic* level.
- It is a starting point for developing new logics: the concept “entails” can be understood in several ways (such as “ $\gamma$  must be somehow relevant to  $\delta$ ”) — what kind of “implication” arises if we understand it like this?
- It also means that we get our program  $\mathcal{P}$  for answering the question “ $\gamma \models \delta$ ?” if we develop an algorithm for the question “Is  $\gamma \Rightarrow \delta$  valid?”.
- This question can in turn be answered by looking only at the logical form of this formula  $\gamma \Rightarrow \delta$ .

- We also get the principle of proof by *refutation* or *contradiction* (or *reductio ad absurdum*):

$$\gamma \models \delta \text{ if and only if } \gamma \wedge \neg\delta \text{ is not satisfiable.} \quad (32)$$

That is: “If  $\delta$  was FALSE, then this would contradict  $\gamma$ ” (= our *KB*, which we maintain satisfiable).

- For propositional logic, this question “Is  $\gamma \Rightarrow \delta$  valid?” is in turn the *complement* of a search problem:

- “Try to assign values TRUE or FALSE to the *Symbols* in  $\gamma \Rightarrow \delta$  so that  $\gamma \Rightarrow \delta$  becomes FALSE. If there are *no* such values, then  $\gamma \Rightarrow \delta$  is valid.”
- Classical propositional validity is **co-NP**-complete.

- The brute-force method shown as Figure 58 for solving SAT is to start

**computing** its truth table

= **trying out** all different assignments of TRUE or FALSE into the *Symbols* appearing in the given *formula* — the other *Symbols* not in it don't need any value and stopping with the answer YES if we get the *formula* to be TRUE.

- Validity can be solved in the same way, but stopping with NO if we get the *formula* to be FALSE.
- This method collects into an initially empty data structure  $v$  the values currently assigned to *Symbols*.
  - If a symbol  $X$  has not been assigned any value yet, then  $v[X] = \perp$ .
  - Then  $RowValue(\alpha, v)$  = the which the given *formula*  $\alpha$  gets with these currently assigned values  $v[Y]$  — but not all *Symbols*  $Y$  have been assigned a value yet.
  - Hence we must permit this ' $\perp$ ' also in our truth tables:

		$\wedge$	$\vee$	$\Rightarrow$	$\Leftrightarrow$
FALSE	<i>other</i>	FALSE	<i>other</i>	TRUE	$\overline{other}$
<i>other</i>	FALSE	FALSE	<i>other</i>	<i>other</i>	$\overline{other}$
TRUE	<i>other</i>	<i>other</i>	TRUE	<i>other</i>	<i>other</i>
<i>other</i>	TRUE	<i>other</i>	TRUE	TRUE	<i>other</i>

(33)

Here  $\overline{other}$  means “flip the *other* value around if it is defined”:

$$\begin{aligned}\overline{TRUE} &= FALSE \\ \overline{FALSE} &= TRUE \\ \overline{\perp} &= \perp.\end{aligned}$$

- This brute-force method takes up to

$\mathcal{O}(m \cdot 2^n)$  steps, where

$m$  = the size of the (parse tree for) formula  $\alpha$

$n$  = the number of distinct *Symbols* in  $\alpha$

simply because it might try all the  $2^n$  distinct assignments  $v$  to the *Symbols* in  $\alpha$ .

- Fortunately we do have algorithms, which take much less in practice.
- However, there are no SAT algorithms which could always guarantee a YES/NO answer in less than exponentially many steps in  $n$  (unless  $\mathbf{P} = \mathbf{NP}$ ).
- One way to present a proof system  $\mathcal{P}$  is as a collection of *inference rules*, like the following *Modus Ponens*:

$$\frac{\alpha \quad \alpha \Rightarrow \beta}{\beta} \quad (34)$$

**Premisses** above the line list all that must hold before this rule can be applied.

**Conclusion** below the line gives what can then be inferred.

- An inference rule like Eq. (34) can be read in two directions.

```

TruthTable( $\alpha, v$ )
1   $r \leftarrow \text{RowValue}(\alpha, v)$ 
2  if  $r = \perp$ 
3      then  $X \leftarrow$  some Symbol in  $\alpha$  such that  $v[X] = \perp$ 
4          return TruthTable( $\alpha, v$  extended with  $v[X] \leftarrow \text{TRUE}$ ) or
               TruthTable( $\alpha, v$  extended with  $v[X] \leftarrow \text{FALSE}$ )
5  else return  $r$ 

RowValue( $\beta, v$ )
1  if  $\beta$  is a Symbol
2      then return  $v[\beta]$ 
3  if  $\beta$  is of the form  $\neg\phi$ 
4      then return  $\overline{\text{RowValue}(\phi, v)}$ 
5  if  $\beta$  is of the form  $\phi \otimes \psi$ 
6      then return the entry for RowValue( $\phi, v$ ), RowValue( $\psi, v$ ) and ‘ $\otimes$ ’
               in Eq. (33)
7  return  $\beta$  (since  $\beta \in \mathbb{B}$  must now be the case)

```

Figure 58: SAT via truth-table enumeration. Modified from Russell and Norvig (2003, Figure 7.10).

**Forward** direction reads “if I have already inferred these premisses, then I can infer this conclusion too”.

**Backward** direction reads “I want to infer this conclusion, so I will try to infer these premisses next”.

This gives two opposite approaches to computing  $KB \vdash_{\mathcal{P}} \beta$ .

- Soundness of  $\mathcal{P}$  requires that its inference rules like Eq. (34) satisfy the condition  
the conjunction of all its premisses  $\models$  its conclusion. (35)

This can in turn be checked (e.g.) by computing the corresponding truth table.

- E.g. our agent could infer as follows in the initial top left situation of the Wumpus world shown as Figure 54:

$$\begin{array}{c}
 \text{background Eq. (26)} \\
 \hline
 (B_{1,1} \Rightarrow P_{1,2} \vee P_{2,1}) \wedge (P_{1,2} \vee P_{2,1} \Rightarrow B_{1,1}) \quad \Leftrightarrow\text{-elimination} \\
 \hline
 P_{1,2} \vee P_{2,1} \Rightarrow B_{1,1} \quad \wedge\text{-elimination} \\
 \hline
 \text{percept Eq. (29)} \quad \neg B_{1,1} \Rightarrow \neg(P_{1,2} \vee P_{2,1}) \quad \text{contraposition} \\
 \hline
 \neg(P_{1,2} \vee P_{2,1}) \quad \text{Eq. (34)} \\
 \hline
 \neg P_{1,2} \wedge \neg P_{2,1} \quad \text{de Morgan}
 \end{array} \tag{36}$$

Here the agent used the following inference rules:

**$\Leftrightarrow$ -elimination** or the definition of ‘ $\Leftrightarrow$ ’ as ‘ $\Rightarrow$ ’ in both directions

$$\frac{\alpha \Leftrightarrow \beta}{(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)} \tag{37}$$

**$\wedge$ -elimination** on the right

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

**contraposition** or flipping ' $\Rightarrow$ ' around with ' $\neg$ '

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

**de Morgan** for moving ' $\neg$ ' in and out of ' $\wedge$ ' or ' $\vee$ '

$$\frac{\neg \left( \begin{array}{c} \vee \\ \alpha \quad \beta \\ \wedge \end{array} \right)}{(\neg \alpha) \wedge (\neg \beta)} \quad (38)$$

(and vice versa).

- Searching for a proof can be much faster than computing the truth table, because the search can ignore much of *KB* as irrelevant for this current task.

E.g. here all knowledge about the **Wumpus** was ignored, because this current task was to reason about **Pits** only.

#### 4.3.2 Propositional Resolution

- Using a proof system  $\mathcal{P}$  with many rules is fine for humans.
- For a computer, the less rules  $\mathcal{P}$  has the simpler the corresponding search algorithm gets.
- Just *one* rule called *resolution* turns out to be enough — if we
  1. first rewrite our *KB* into a certain fixed form
  2. then seek refutation proofs, as in Eq. (32).
- The general idea of resolution is

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

or “if the same disjunct  $\beta$  occurs both

**positively** (= outside ' $\neg$ ') in one disjunction and

**negatively** (= inside ' $\neg$ ') in another disjunction

then you can combine these two disjunctions into one *resolvent* and drop their common  $\beta$ ”.

- E.g. our agent could infer as follows in the bottom left situation of the Wumpus world shown as Figure 54:

$$\frac{\frac{\frac{B_{2,1} \quad B_{2,1} \Rightarrow P_{1,1} \vee P_{2,2} \vee P_{3,1} \text{ like in Eq. (36)}}{P_{1,1} \vee P_{2,2} \vee P_{3,1}} \quad \neg P_{1,1}}{P_{2,2} \vee P_{3,1}} \text{ resolution} \quad \frac{\neg P_{2,2} \text{ like in Eq. (36)}}{P_{3,1}} \text{ resolution}$$

The associativity and commutativity of ' $\vee$ ' lets us pick any disjunct in these resolution steps.



- In both of these two resolution steps, the common disjunct  $\beta$  is just a single *Symbol*.
  - We are going to rewrite our *KB* so that this will always be the case.
  - The reason is that an unmodified *KB* would not contain very many possible choices for  $\beta$ .

- The form into which we are going to rewrite our *KB* can be defined as follows:

**Literal** is *Symbol* or  $\neg \text{Symbol}$ .

**Clause** is a (possibly empty) disjunction of literals.

Then our *KB* is in *Conjunctive Normal Form (CNF)* if it is a conjunction of clauses.

- The full resolution rule takes the form “you can form the resolvent of two clauses if the same *Symbol*  $X$  occurs positively in one of them and negatively in the other”:

$$\frac{\overbrace{p_1 \vee p_2 \vee p_3 \vee \cdots \vee p_m \vee X}^{\text{one clause}} \quad \overbrace{\neg X \vee q_1 \vee q_2 \vee q_3 \vee \cdots \vee q_n}^{\text{the other clause}}}{p_1 \vee p_2 \vee p_3 \vee \cdots \vee p_m \vee q_1 \vee q_2 \vee q_3 \vee \cdots \vee q_n}$$

By the associativity and commutativity of ‘ $\vee$ ’ this common symbol  $X$  can appear anywhere inside these two clauses.

- We must also *factor* the resolvent: if some literal would appear more than once, only one copy is retained.
- Or alternatively we can think that a clause is a *set* of literals:

$$\frac{\{p_1, p_2, p_3, \dots, p_m, X\} \quad \{\neg X, q_1, q_2, q_3, \dots, q_n\}}{\{p_1, p_2, p_3, \dots, p_m, q_1, q_2, q_3, \dots, q_n\}}$$

This form includes the factoring rule implicitly.

- A given *formula* can be converted into CNF with the following 4 steps:
  1. Replace each occurrence of ‘ $\Leftrightarrow$ ’ with the corresponding two occurrences of ‘ $\Rightarrow$ ’ as in Eq. (37).

E.g. Eq. (26) of our Wumpus world background knowledge becomes

$$(B_{1,1} \Rightarrow P_{1,2} \vee P_{2,1}) \wedge (P_{1,2} \vee P_{2,1} \Rightarrow B_{1,1}). \quad (39)$$

2. Replace each occurrence of  $\alpha \Rightarrow \beta$  with the (classically) equivalent  $\neg \alpha \vee \beta$ .

E.g. Eq. (39) becomes

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg(P_{1,2} \vee P_{2,1}) \vee B_{1,1}). \quad (40)$$

3. Move each ‘ $\neg$ ’ towards the *Symbols* using Eq. (38). If you get a double negation like  $\neg \neg \alpha$ , then erase them, leaving only  $\alpha$ .

E.g. Eq. (40) becomes

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge ((\neg P_{1,2} \wedge \neg P_{2,1}) \vee B_{1,1}). \quad (41)$$

4. Finally move each ‘ $\wedge$ ’ from under any ‘ $\vee$ ’ by using their distributivity, which permits replacing  $\alpha \vee (\beta \wedge \gamma)$  with  $(\alpha \vee \beta) \wedge (\alpha \vee \gamma)$  and so on.

E.g. Eq. (41) becomes

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg P_{1,2} \vee B_{1,1}) \wedge (\neg P_{2,1} \vee B_{1,1}). \quad (42)$$

If we now undo step 2, then we see that Eq. 42 does indeed say the same thing as the original Eq. (26), but in a different way:

$$(B_{1,1} \Rightarrow P_{1,2} \vee P_{2,1}) \wedge (P_{1,2} \Rightarrow B_{1,1}) \wedge (P_{2,1} \Rightarrow B_{1,1}).$$