

Chapter 10

Society

The phrase aroused my interest because of its enigmatic quality: "He brought in a second actor." I stopped; I found that the subject of that mysterious action was Aeschylus and that, as we read in the fourth chapter of Aristotle's Poetics, he "raised the number of actors from one to two" ... With the second actor came the dialogue and the indefinite possibilities of the reaction of some characters on others. A prophetic spectator would have seen that multitudes of future appearances accompanied him: Hamlet and Faust and Segismundo and Macbeth and Peer Gynt and others our eyes cannot yet discern.

Jorge Luis Borges, "The Modesty of History" in *Other Inquisitions*

This final chapter deals with commonsense knowledge about the interrelations and interactions among agents. In this domain more than any other, commonsense understanding far outstrips any existing formal model in richness and sophistication. Such an understanding is vital for story understanding, in automated teaching and in intelligent interactive systems. In this chapter, we discuss some issues in human and social interactions that have been studied in the AI literature: common knowledge, plan interaction, communication, ethical values, and possession.

The last three of these issues — communication, ethics, and possession — share a number of characteristics. All three issues arise continually in life in a complex society. All are virtually universal among human societies¹ and, apparently, appear in rudimentary form among some nonhuman species. All three are determined, at least to

¹There are some nomadic hunting societies that carry virtually nothing with them on their travels. Property is an irrelevant concept to these. As far as I know, all societies have some language and some concept of the permitted and the prohibited.

some degree, by social convention,² but within the society they can often be extremely rigid and inescapable. It is easy to develop a representation for any of these issues that will cover the simplest cases; however, giving a deep theory of any raises profound and unresolved philosophical controversies.

10.1 Common Knowledge

Suppose that Tom and Huck are hiding together, and they witness Joe murder Dr. Robinson. Then Tom, Huck, and Joe each know the proposition “Joe murdered the doctor.” Moreover, Tom and Huck each know the proposition “Tom, Huck, and Joe each know that Joe murdered the doctor”; Joe, however, does not know this. Further, Tom and Huck each know the proposition “Tom and Huck each know that Tom and Huck each know that Joe murdered the doctor,” and so on *ad infinitum*. In short, Tom and Huck have a complete understanding between them about the fact.

The technical term for this is *common knowledge*: Tom and Huck have common knowledge that Joe murdered Dr. Robinson. To represent this, we introduce the syntactic operator “common_know(*SAA, P*)”, meaning that the set of agents denoted by string *SAA* have common knowledge of sentence *P*. (A third temporal argument can be added when necessary.) Thus, “Tom and Huck have common knowledge that Joe murdered Dr. Robinson” is represented

common_know(\prec
 $\{ \text{Tom, Huck} \} \succ, \prec \exists_I \text{ occur}(I, \text{murder(joe, robinson)}) \succ$)

It is possible for every pair of agents in a set to have common knowledge of a fact without there being common knowledge of the fact over the set as a whole. For example, if *A* and *B* share a secret ϕ , and each independently tells it to *C*, then each of the sets $\{ A, B \}$, $\{ B, C \}$, and $\{ A, C \}$ has common knowledge of ϕ , but the set $\{ A, B, C \}$ does not. For example, *A* does not know that *B* knows that *C* knows ϕ . (See Exercise 1 for an example where this distinction makes a big difference.)

Unlike simple knowledge, common knowledge is opaque in the agent argument³; the members of a set may share knowledge under one description of the set but not another. For example, the secret hand-

²In defense of the statement that ethical values or property are to some degree a matter of social convention, let me point out that I am including in these categories such rules as “Don’t drive through red lights” or “A person is considered to have paid a bill as of the moment the biller receives his check, not when the check clears.”

³Lesperance [1989] argues that the same should be true of single agents.

shake of the Masons is (I presume) common knowledge among the Masons: All Masons know it, know that all other Masons know it, know that other Masons know that all Masons know it, etc. Nonetheless, it does not follow that if Joseph and Dominic are both Masons, then Joseph knows that Dominic knows the secret handshake; Joseph may not know that Dominic is a Mason. Another example: In Renaissance Venice, the Council of Three were secretly chosen; they met masked and did not know each other's identity (according to Mark Twain). They presumably shared a large body of common knowledge. The same three individuals (being prominent Venetians) undoubtedly knew each other, and may, indeed, have shared common knowledge in mufti. The two bodies of common knowledge, however, were quite separate. In the first case, each of them knew that all the members of the Council of Three knew the facts; in the other case, each of them knew that each member of the set { Leonardo, Giuseppe, Pietro } knew the facts. Thus, in a modal theory, the agent argument is an opaque argument; in a syntactic theory, it is a string that denotes a set of agents; in a possible-worlds theory, it is a fluent ranging over sets of agents, whose value may vary from one possible world to another.

Table 10.1 shows some basic axioms of common knowledge in a syntactic language.

Axioms CK.1 and CK.2 establish the basic properties of common knowledge: If agents AA have common knowledge of ϕ , then they all know ϕ , and they all know that they all know ϕ , and so on. Axiom CK.1a, that common knowledge is veridical, is a direct consequence of CK.1 and KNOW.3, that knowledge is veridical. Axioms CK.3 and CK.4 are the usual consequential-closure principle, analogous to axioms KNOW.1 and KNOW.2 of knowledge. Axiom CK.5 states that, if agents AA have common knowledge of ϕ and agents BB have common knowledge that BB is a subset of AA , then agents BB have common knowledge of ϕ . We can use this rule to justify inferences like the following: If Masons have common knowledge of the Masonic handshake, and Joseph and Dominic have common knowledge that they are Masons, then Joseph and Dominic have common knowledge of the Masonic handshake.

Finding an appropriate set description for a set of agents with common knowledge can raise substantial difficulties. For example, consider three strangers who are standing together and together witness a rabbit eating a carrot. Then it is clear that the three people have common knowledge that the rabbit ate the carrot, but under what description of one another? One solution would be to say that the common knowledge is shared among "the people here now," but that involves positing a common language of indexicals among agents. It

Table 10.1 Axioms of Common Knowledge

CK.1. If P is common knowledge among SAA , then every agent in SAA knows P .

$$[\text{common_know}(SAA, P) \wedge A \in \text{denotation}(SAA)] \Rightarrow \text{know}(A, P).$$

CK.1a (Veridicality) If P is common knowledge among SAA , then P is true.

$$\text{common_know}(SAA, P) \Rightarrow \text{true}(P).$$

CK.2. (Positive introspection) If P is common knowledge among SAA , then it is common knowledge among SAA that P is common knowledge among SAA .

$$\text{common_know}(SAA, P) \Rightarrow \text{common_know}(SAA, \prec \text{common_know}(!SAA!, !P!) \succ).$$

CK.3. (Consequential closure) Common knowledge is closed under implication.

$$[\text{common_know}(AA, P) \wedge \text{common_know}(AA, \prec \downarrow P \Rightarrow \downarrow Q \downarrow \succ)] \Rightarrow \text{common_know}(AA, Q)$$

CK.4. (Common knowledge of the axioms) Basic axioms are common knowledge.

If P spells out an axiom (of logic, knowledge, time, or whatever), then

$$\forall_{SAA} \text{common_know}(SAA, P)$$

is an axiom.

CK.5. (Common knowledge among a subset)

$$[\text{common_know}(SAA, P) \wedge \text{common_know}(SBB, \prec \downarrow SBB \downarrow \subseteq \downarrow SAA \downarrow \succ)] \Rightarrow \text{common_know}(SBB, P).$$

is worthwhile distinguishing situations where this question does not arise, because all the agents involved know each other. We will represent the state in which all the agents in set AA know one another by common names and have common knowledge of this fact using the state type "acquainted(AA)." To represent common knowledge in a possible-worlds semantics, we use an accessibility relation "ck_acc($FAA, W1, W2$)," where $W1$ and

W_2 are possible worlds and FAA is a fluent that ranges over sets of agents in different worlds. FAA thus corresponds to a particular description of some set of agents. The relation $ck_acc(FAA, W_1, W_2)$ holds if W_2 is consistent with everything that is common knowledge in W_1 to the set of agents described by FAA . The statement that agents AA have common knowledge of ϕ in W_1 is expressed by asserting that ϕ is true in every accessible world.

10.2 Multiagent Plans

An agent in a society of many agents must consider his neighbors' actions and reactions in forming and evaluating his plans. Interactions among plans can be divided into a number of basic categories:

- *Cooperation*: A set of agents with a common goal cooperate on a common plan. The structure and analysis of such a plan may be considerably more complicated than the analysis of a single-agent plan. In many cases, there will not be any single agent that knows all steps of the plan. Step-by-step validation of the plan, such as we have seen in single-agent plans, therefore becomes largely irrelevant; there is no agent who knows enough to carry out such a validation. Instead, the analysis centers around the ways in which subtasks are divided among the various agents. Another new issue that becomes important in cooperative plans is the communication of relevant information between agents. This type of planning has much in common with the programming of distributed computer systems [Smith and Davis 1981].
- *Influence*: In a plan of influence, the planner tries to get others to perform actions that further his goals. Here, the other agents are analogous to external physical events that the agent is trying to control. The difference is that the agents are governed by psychological and social rules, rather than physical rules. A major aspect of this planning is predicting how another agent will respond to a request, discussed in Section 10.3.2. Section 10.3.3 shows how a simple plan of influence, involving making a request, can be validated.
- *Contingent interaction*: Other agents intend to perform actions, purely for their own purposes, that may aid or hinder the plan of the planner. Again, as in plans of influence, the other agents may be considered as external events that can be affected in a number of ways.

- *Direct conflict*: Another agent specifically wishes to prevent the planner from accomplishing his goal (or vice versa). Here, the planner must assume that the opponent will always perform the action that he (the opponent) considers most obstructive. Such planning often involves an exchange of hostile, destructive, or deceptive actions.

A single plan may involve different kinds of interactions. For example, a group of agents may cooperate on a plan that involves influencing some other agents and conflict with others. There are also many cases that fall on the borderline. For example, the workings of a company are, in some respects, a cooperative activity of the employees; in some respects, an influencing of the workers by the boss; and in some respects, a conflict between the workers and the boss.

Research in representing and reasoning about such plans is currently in a primitive stage, though some intriguing work has been done. Citations are given in the chapter reference list.

10.3 Communication

The most central type of interactions among intelligent agents are communications, also known as *speech acts*. (We use the terms associated with speech — “speech,” “speaker,” “hearer,” and so on — to apply to any mode of communication: speech, writing, telegraphy, etc.) Speech acts may be divided into five categories: *declarative* acts, which convey information, such as stating “Paris is the capital of France”; *interrogative* acts, which ask for information, such as asking “What is the capital of France?”; *imperative* acts, which make a request or issue a command, such as asking “Please go away”; *exclamatory* acts, which express an emotion, such as crying “Alas!”; and *performative* acts, which, by social convention, bring about a condition, such as “I now proclaim you man and wife” or “I hereby cede my rights to Sherwood Forest.” These categories are (nearly) mutually exclusive and exhaustive; almost every speech act belongs to exactly one of the categories.

Another well-known categorization associated with speech acts is the division into *locutionary*, *illocutionary* and *perlocutionary* acts [Austin 1961]. Unlike the previous categorization, these are categories of speech-act *types*, rather than of *tokens*; any given speech-act token may be described in any of these ways. A locutionary speech act describes the act in terms of the physical characteristics of the acts and the symbols used. An illocutionary speech act describes the act in

Table 10.2 Speech Act Categories

Locutionary:	Shouting “I’m mad as hell!” Whispering “Te amo.” Writing “John Hancock.”
Illocutionary:	Declaring allegiance to the king. Asking for a raise. Proposing a merger between General Soap and Urumchi University Press.
Perlocutionary:	Boring the audience. Converting a student to Marxism. Charming a guest.

terms of its content. A perlocutionary speech act describes the act in terms of the effect on the hearer. Table 10.2 gives some examples of each:

In this section, we will study locutionary and illocutionary descriptions of declarative, interrogative, and imperative speech acts.

10.3.1 Locutionary Descriptions

Locutionary acts are physical actions that produce sequences of signs in speech, writing, sign language, or other physical medium. We introduce the function “*speak(P)*,” mapping a string of phonemes *P* to the action of speaking *P*, and the predicate “*pronunciation(S, P, L)*,” meaning that phonemes *P* are an acceptable pronunciation of string *S* in language *L*. Thus, we can express the statements “Humphrey Bogart said ‘Here’s looking at you’” and “Lincoln spoke the Gettysburg Address” in the forms

$$\begin{aligned} \exists_P \text{pronunciation}(P, \langle \text{Here's_looking_at_you} \rangle, \text{english}) \wedge \\ \text{occur}(\text{i202}, \text{do}(\text{bogart}, \text{speak}(P))). \\ \exists_P \text{pronunciation}(P, \text{gettysburg_address}, \text{english}) \wedge \\ \text{occur}(\text{i624}, \text{do}(\text{lincoln}, \text{speak}(P))). \end{aligned}$$

where “*gettysburg_address*” is a constant denoting the string $\langle \text{Four score and seven years ago ...} \rangle$.

Speech can accomplish communication only to someone who is present to hear it. Since hearing is a perception, we may apply the theory of perception developed in Section 8.7. For example, we can express a rule such as “If hearer AH is within distance d_hear of speaker AS , then he can hear whether AS is speaking and what he is saying” in the following form:

$$\begin{aligned} & [\forall_{S \in I} \text{distance}(\text{value_in}(S, \text{place}(AH)), \\ & \quad \text{value_in}(S, \text{place}(AS))) \leq d_hear \wedge \\ & \quad \text{bpc}(AH, B, \text{behavior}(I))] \Rightarrow \\ & [\text{occurs}(I, \text{do}(AS, \text{speak}(P))) \Leftrightarrow \exists_{AS1} \text{occurs}(B, \text{do}(AS1, \text{speak}(P)))]. \end{aligned}$$

10.3.2 Illocutionary Speech Acts

An illocutionary description of a speech act characterizes its content. We introduce the event type, “ $\text{illoc}(AS, AH, M, P)$ ”, the event of AS communicating P in mode M to AH . (For simplicity, we assume a single speaker and hearer.) In a declarative act, P will be a sentence in a formal language, and M will be the constant “declarative.” (Of course, declarative speech acts may be carried out without speaking an entire natural-language sentence. For instance, a question can often be answered in a single word. We assume, however, that the content of any declarative speech act can be expressed as a sentence in a formal language.) In an imperative act, P will be a term in a formal language denoting the action type that the speaker wants the hearer to perform, and M will be “imperative.” We will treat interrogatives as special types of imperatives in which the speaker is requesting the hearer to communicate an answer to his question.

For example, the statement “Becky told Amelia that all crows are black” may be represented

$$\begin{aligned} & \text{occur}(i202, \text{illoc}(becky, amelia, \text{declarative}, \\ & \quad \neg \forall_X \text{crow}(X) \Rightarrow \text{black}(X))). \end{aligned}$$

The statement “Becky told Amelia where Joseph was” is interpreted as “Becky told Amelia a true statement of the form ‘Joseph is at X ’ for some constant X ,” and is thus represented

$$\begin{aligned} & \exists_{X,P} P = \neg \text{value_in}(\text{start}(i202), \text{place}(joseph)) = \downarrow X \downarrow \wedge \\ & \quad \text{constant}(X) \wedge \text{true}(P) \wedge \\ & \quad \text{occur}(i202, \text{illoc}(becky, amelia, \text{declarative}, P)). \end{aligned}$$

The statement “Becky asked Amelia to pass the salt” is represented

occur(i202, illoc(becky, amelia, imperative,
 do(amelia.pass(salt42))))

The statement, "Amelia asked Becky where Joseph was" is interpreted as "Amelia requested that Becky communicate a true sentence of the form 'Joseph is at X ,' where X is a constant," which is represented

occur(i202, illoc(amelia, becky, imperative,
 $\nwarrow \{ I \mid \exists_{X,P} \text{occur}(I, \text{iloc}(becky, amelia, declarative, } P)) \wedge$
 $P = \neg \text{value_in}(\text{start}(i202), \text{place}(joseph)) = \downarrow X \uparrow \neg \wedge$
 $\text{true}(P) \wedge \text{constant}(X) \} \nwarrow \}).$

(Note: the down arrows surrounding X in the formula above are scoped to the inner string delimiters.) Table 10.3 shows some basic properties of illocutionary acts.

Axioms IL.1 and IL.3 characterize the relation between locutionary and illocutionary acts. They can serve as the basis for executing an illocutionary act as a locutionary act in a task-reduction planner. The functions “meaning_of(S, K)” and “mode_of(S, K)” in these axioms specify the relation between the string spoken and the content communicated; the content of these functions would be specified by a theory of natural-language semantics. Axioms IL.6 and IL.8 state some necessary conditions that declarative and imperative speech acts be sincere. Additional conditions could be added. For example, we could define a declaration of P to be sincere if the speaker both believes P and believes that the hearer will believe that the speaker is sincere. We could define a request of P to be sincere if the speaker desires P and believes that his requesting P will make the hearer more likely to perform P [Searle 1969].

By combining these axioms with the axioms of belief in Chapter 8, we can justify⁴ the plausible inference that if A declares ϕ to B , then B will believe ϕ . Assume that A declares ϕ to B . By IL.2, A and B have common knowledge that A has declared ϕ to B . By axiom CK.1, B knows that A has declared ϕ to B . Using rule BEL.14, we can presume that B will apply IL.4 and infer that A believes ϕ , and that he will further apply rule BEL.13, and infer that ϕ is true.

Stronger initial information will support stronger conclusions. For example, assume (i) that A declares ϕ to B ; (ii) that A and B have common knowledge that A is speaking sincerely; and (iii) that A and B have common knowledge that A knows whether ϕ is true. Then it

⁴Modulo the limitation that our notation “plausible” doesn’t actually refer to any particular theory of plausible reasoning.

Table 10.3 Axioms for Illocutionary Acts

IL.1. $\text{token_of}(K, \text{illoc}(AS, AH, M, P)) \Rightarrow \exists_S \text{token_of}(K, \text{do}(AS, \text{speak}(S))) \wedge P = \text{meaning_of}(S, K) \wedge M = \text{mode_of}(S, K)$.
 An illocutionary act communicating P in mode M involves speaking a string whose meaning is P and mode is M .

IL.2. $[\text{occur}(I, \text{illoc}(AS, AH, M, P)) \wedge \text{true_in}(\text{start}(I), \text{acquainted}(\{AS, AH\}))] \Rightarrow \text{common_know}(\neg\{AS, AH\}, \neg\text{occur}(I, \text{illoc}(AS, AH, M, P)), \text{end}(I))$.
 The occurrence of an illocutionary act is a matter of common knowledge to speaker and hearer.

IL.3. $\text{true_in}(\text{start}(I), \text{acquainted}(\{AS, AH\})) \Rightarrow [\text{occur}(I, \text{illoc}(AS, AH, M, P)) \Leftrightarrow \text{common_know}(\neg\{AS, AH\}, \neg\exists_{S,K} \text{token_of}(K, \text{do}(AS, \text{speak}(S))) \wedge \text{@}I = \text{time_of}(K) \wedge P = \text{meaning_of}(S, K) \wedge M = \text{mode_of}(S, K), \text{end}(I))]$.
 A communication of P from AS to AH occurs just if they have common knowledge that AS spoke some string S that means P .

IL.4. $\text{plausible}(\text{token_of}(K, \text{illoc}(AS, AH, M, P)), \text{sincere}(K))$.
 Illocutionary acts are typically sincere.

IL.5. $\text{token_of}(K, \text{illoc}(AS, AH, \text{declarative}, P)) \Rightarrow \text{sentence}(P)$.
 The content of a declarative act is a sentence.

IL.6. $[\text{token_of}(K, \text{illoc}(AS, AH, \text{declarative}, P)) \wedge \text{sincere}(K)] \Rightarrow \text{believe}(AS, P, \text{start}(\text{time_of}(K)))$.
 The speaker believes a sincere declarative speech act.

IL.7. $\text{token_of}(K, \text{illoc}(AS, AH, \text{imperative}, P)) \Rightarrow \text{sort_of}(\text{denotation}(P)) = \text{event} \wedge \text{actor_of}(\text{denotation}(P)) = AH$.
 The content of an imperative act is a description of an action by the hearer.

IL.8. $[\text{token_of}(K, \text{illoc}(AS, AH, \text{imperative}, P)) \wedge \text{sincere}(K)] \Rightarrow \text{goal}(AS, P, \text{start}(\text{time_of}(K)))$.
 If an imperative act is sincere, then the speaker wants the specified action to be carried out.

is possible to infer that A and B have common knowledge of ϕ when the communication is complete (Exercise 3).

In an analogous way, one would like to support the inference that if A requests something from B then, under suitable conditions, B will perform it. Such an inference is necessary to construct plans that involve the cooperation of agents who do not necessarily share the planner's goals. The problem is to spell out the requisite conditions. So far, this has only been done for very idealized microworlds. Section 10.3.3 shows the validation of a simple plan of influence involving a request, using an *ad hoc* assumption that this particular hearer always carry out the requests of the speaker immediately, if possible. (Schank and Abelson [1977] discuss a sequence of request forms, in increasing order of power and decreasing order of ease: asking; invoking a theme relating the two actors (e.g., "We've always been friends"); informing the hearer of the reason for the request; bargaining with the hearer; threatening the hearer; and overpowering the hearer.)

10.3.3 Sample Verification of a Plan of Influence

In this section, we illustrate how a simple plan of influence in a toy world can be verified. This verification combines primitives and rules from Chapters 5, 8, 9, and 10.

Assume that, in the blocks world of Section 5.3, block A is clear and on top of block B and the hand is above them. Fred manipulates the hand. Jack wishes block B to be clear. Show that Jack can satisfy his goal by asking Fred to pick up block A.

We must here assume that branches in time correspond only to Jack's actions, not to Fred's. Therefore, feasibility for blocks-world actions should be defined in domain-specific axioms, not in terms of branching time. For example, the feasibility of a pickup would be defined in the following rule:

$$\begin{aligned} \text{true_in}(S, \text{feasible}(\text{pickup})) \Leftrightarrow \\ \text{true_in}(S, \text{clear}(\text{hand})) \wedge \exists_X \text{ true_in}(S, \text{under_hand}(X)) \end{aligned}$$

An alternative approach would be to distinguish between what Fred could do in given situation (any physically feasible action) and what he would do (what Jack asks). See Chapter 5, exercise 7.

Assumptions:

i. Starting situation.

$$\begin{aligned} \text{true_in}(s0, \text{clear}(\text{hand})) \wedge \text{true_in}(s0, \text{under_hand}(\text{blocka})) \wedge \\ [\forall_X \text{ true_in}(s0, \text{beneath}(\text{blockb}, X)) \Leftrightarrow X = \text{blocka}] \end{aligned}$$

- ii. (Definition) A hand command is “pickup” or “putdown” or “move(L).”
 $\text{hand_command}(E) \Leftrightarrow$
 $[E = \langle \text{pickup} \rangle \vee E = \langle \text{putdown} \rangle \vee \exists L E = \langle \text{move}(\downarrow L \downarrow) \rangle].$
- iii. All the hand’s activities are deliberate activities of Fred’s.
 $\text{hand_command}(E) \Rightarrow$
 $[\text{token_of}(K, \text{denotation}(E)) \Leftrightarrow \text{token_of}(K, \text{deliberate}(\text{fred}, E))].$
- iv. Fred executes a hand command if Jack asks him to and it is physically possible to do so.
 $[\text{hand_command}(P) \wedge E = \text{denotation}(P)] \Rightarrow$
 $[\text{occurs}(I, E) \Leftrightarrow$
 $\exists_{I1} \text{meet}(I1, I) \wedge \text{occur}(I1, \text{iloc}(\text{jack}, \text{fred}, \text{imperative}, P)) \wedge$
 $\text{can_do}(\text{fred}, P, \text{start}(I))].$
- v. Fred and Jack are permanently acquainted.
 $\text{true_in}(S, \text{acquainted}(\text{fred}, \text{jack})).$
- vi. The sentence “Please execute a pickup” is under all circumstances a request for a pickup.
 $\text{pronounciation}(P, \langle \text{Please_execute_a_pickup} \rangle, \text{english}) \Rightarrow$
 $\text{mode_of}(P, K) = \text{imperative} \wedge \text{meaning_of}(P, K) = \text{pickup}.$
- vii. Jack can always speak any string.
 $\forall_{S, P} \text{true_in}(S, \text{feasible}(\text{do}(\text{jack}, \text{speak}(P)))).$
- viii. If Jack speaks a string, then it will be common knowledge between him and Fred that he has spoken the string. (Note: This assumption is a quick and dirty device to get around specifying the conditions under which a speech of Jack’s is known by Fred, and the frame axioms associated with those conditions.)
 $\text{occur}(I, \text{do}(\text{jack}, \text{speak}(P))) \Rightarrow$
 $\text{common_know}(\langle \{\text{jack}, \text{fred} \} \rangle,$
 $\langle \text{occur}(@I@, \text{do}(\text{jack}, \text{speak}(!P!))) \rangle, \text{end}(I)).$
- ix. If Jack now says “Please execute a pickup,” then the only illocutionary act that he will complete at any time during his speech is to request a pickup.
 $[\text{pronunciation}(P, \langle \text{Please_execute_a_pickup} \rangle, \text{english}) \wedge$
 $\text{token_of}(K, \text{do}(\text{jack}(\text{speak}(P)))) \wedge \text{start}(\text{time_of}(K)) = s0 \wedge$
 $\text{token_of}(K1, \text{iloc}(\text{jack}, A2, M, R)) \wedge$
 $\text{end}(\text{time_of}(K1)) \in \text{time_of}(K)] \Rightarrow$
 $K1 = K.$

- x. In s_0 , Fred has completed any blocks-world action that he started previously.
 $[E=\text{pickup} \wedge E=\text{putdown} \wedge E=\text{move}(L) \wedge \text{occur}(I, E) \wedge \text{start}(I) < s_0] \Rightarrow \text{end}(I) \leq s_0$.
- xi. Jack knows about all blocks-world events.
 $[\text{hand_command}(E) \wedge \text{occurs}(I, \text{denotation}(E))] \Rightarrow \text{know}(\text{jack}, \neg \text{occur}(@I@, \downarrow E \downarrow), \text{end}(I))$.
- xii. Fred and Jack have common knowledge of all the blocks-world axioms plus assumptions i through xi, above.

To prove: Jack can accomplish block B being clear.
 $\text{can_achieve}(\text{jack}, \text{clear}(\text{blockb}), s_0)$.

Proof: Let plan0 be the plan,

$\neg \text{sequence}(\text{do}(\text{jack}, \text{speak}(\neg \text{Please_execute_a_pickup})), \text{wait_while}(\text{pickup}))$.

By the definition of “can_achieve” (axioms KPS.8, KPS.7), we can prove the desired result “ $\text{can_achieve}(\text{jack}, \text{clear}(\text{blockb}), s_0)$ ” by showing that Jack knows that plan0 is feasible, that it leads to $\text{clear}(\text{blockb})$, and that its knowledge preconditions are satisfied. Since Jack knows all the axioms (assumption xii and CK.1) and we assume consequential closure on knowledge (axiom KNOW.1), it suffices to show that it follows from the axioms that plan0 is a valid plan to accomplish $\text{clear}(\text{blockb})$ and that all the knowledge preconditions are met.

From PL1.a it is easily shown that a plan consisting of two actions in sequence is feasible just if the first step is feasible in the starting situation, and the second step is feasible at the end of the first step. In our case, the first step is “ $\text{do}(\text{jack}(\text{speak}(\neg \text{Please_execute_a_pickup})))$,” which, by assumption vii, is always feasible. Given the frame axioms on the blocks world, the state of the blocks changes only if a pickup, putdown, or move is executed. Assumption iii states that such events are actions of Fred’s. Assumption iv states that Fred will perform such an action only if Jack asks him to. Assumption ix states that Jack does not complete any such request between s_0 and the completion of speaking “Please execute a pickup.” Assumption x states that in s_0 , Fred is not in the middle of an action. Therefore, it follows that Fred does not perform any hand actions until the completion of Jack’s speaking. Therefore, using the frame axioms of the blocks world, the blocks remain as they were at the start throughout Jack’s speech.

Assumptions vi, viii, and xii guarantee that the conditions of IL.3 are satisfied, and that Jack’s speech is a request to Fred for a pickup.

By assumption iv, Fred will execute the requested pickup if he can do so; that is, (by definition of "can_do"), if he knows that it is feasible and that its knowledge preconditions are satisfied. Since "pickup" is a constant, the knowledge preconditions are always satisfied. By the axioms of the blocks world, the pickup is currently feasible. Moreover, since Fred knew the state of the blocks at the start (assumption xii), and he knows that he has not committed any deliberate acts (KPG.2), applying assumptions iii and xi, KNOW.1, and the blocks-world frame axioms he knows that the blocks are in the starting position and that the pickup is feasible. Since he knows that the pickup is feasible and that its knowledge preconditions are satisfied, he can execute the pickup; hence, by assumption iv, he will. By the blocks-world axioms, this will result in block B being clear.

Thus, the plan plan0 is feasible and results in block B being clear. By axiom KPS.5, the knowledge preconditions of plan0 are that the knowledge preconditions of the speech are satisfied in the initial situation, and that Jack know that the knowledge preconditions of the "wait" are satisfied in the middle situation. But the speech involves only a constant argument, and so has no knowledge preconditions. By axiom KPS.4, the knowledge preconditions of the action "wait while the pickup takes place" are that Jack is sure that the pickup will take place, and that he will know when it has taken place. Jack can determine that the pickup will take place from the argument above, since Jack knows all the premises of that argument. Assumption xi says that Jack will know when the pickup does take place. Thus, the knowledge preconditions are satisfied.

Finally, again, Jack knows all the premises of the above argument (assumption xi) so, by axiom KNOW.1, he knows that plan0 is feasible, that it accomplishes the goal, and that its knowledge precondition are satisfied. Thus, Jack can achieve block B being clear. Q.E.D.

10.4 Ethics

The responsible builder of an autonomous robot will design its planning module so that it perfers doing right to doing wrong. Further, though the prudent robot will not assume that all other agents are as moral as itself, in order to avoid paranoia it must assume that most other agents rarely break important prohibitions.

A crude representation of ethical valuation can be developed along the following lines: We will consider ethical qualities as characteristics of an action type E in a situation S . We introduce the predicates

Table 10.4 Axioms for Ethical Predicates

ETH.1.	If E is obligatory in S , then it is permitted. Equivalently, it is never obligatory both to do E and not to do it. $\text{obligatory}(E, S) \Rightarrow \text{permitted}(E, S)$.
ETH.2.	It is always possible to avoid all prohibited actions. $\forall_S \exists_C \forall_{I,E} [I \subset C \wedge \text{start}(I)=S \wedge \text{occurs}(I, E)] \Rightarrow \text{permitted}(E, S)$.
ETH.3	Obligatory actions are better than prohibited ones. $[\text{obligatory}(EO, S) \wedge \neg\text{permitted}(EP, S)] \Rightarrow \text{goodness}(EO, S) > \text{goodness}(EP, S)$.
ETH.4	Actions can be assumed to be permitted, unless there is reason to believe them prohibited. $\text{plausible}(\text{true}, \text{permitted}(E, S))$.
ETH.5	Prohibited actions are generally avoided. $\text{plausible}(\neg\text{permitted}(E, S), \forall_I [I \subset \text{real_chronicle} \wedge \text{start}(I)=S] \Rightarrow \neg\text{occurs}(I, E))$.

“permitted(E, S)” and “obligatory(E, S)”. (Our use of first-order predicates implies that we assume that the ethical qualities of an action are transparent under substitution of equals.) Obligation and permission are dual relations; it is obligatory to do E in S just if it is not permitted not to do E . Formally, we can express this dual relation as follows:

$$\text{obligatory}(E, S) \Leftrightarrow \neg\text{permitted}(\text{nonoccurrence}(E), S).$$

Another representation of ethical values uses a measure space of “goodness” and a function “goodness(E, S)” mapping an event type E and a situation S to a measure of goodness.

Table 10.4 shows some plausible, though perhaps optimistic, axioms on obligation, permission, and goodness.

The main body of a theory of ethical values would be a specification of what acts are permitted, obligatory, prohibited, good, or bad. This would involve both specific rules (e.g., “Killing people is prohibited except under very rare circumstances”) and general rules (e.g., “If doing

A will leave everyone happier than doing *B*, then, in general, *A* is a better action than *B*.”) Formulating such a collection of rules is far beyond the scope of this book.

10.5 Possession

The following crude theory of possession will handle simple cases. We assume that an object *O* is owned by at most one agent at a time. The owner of *O* can change only by virtue of a transfer by its current owner. Use of *O* by agents other than the current owner are permitted only if the owner permits it. Table 10.5 shows the formal statement of this theory.

This theory is completed by a specification of what actions constitute the use of an object and how the owner of an object permits some particular use of it.

10.6 Appendix A: Conceptual Dependency

Conceptual dependency (CD) is one of the oldest and most influential representations of commonsense knowledge. It focuses on the representation of primitive human actions for the purpose of representing narrative text. Since it deals with issues spanning a number of our chapters, we have postponed discussing it to this final appendix.

CD was developed by Roger Schank and his associates and students over a number of years [Schank 1969, 1975; Schank and Abelson 1977] during which time it underwent continual development. It was not defined within a formal setting. Therefore, our account below, which integrates CD with the theories we have developed in this book, is to some extent a *post hoc* reconstruction. (Our account closely follows the formalization of CD in [Charniak and McDermott 1985].)

CD is primarily concerned with the characterization of primitive human actions. The basic function is the “*action(ACTOR, ACTION, OBJECT, SOURCE, DESTINATION)*,” which maps its arguments to an event type. *ACTOR* is the actor of the event. *ACTION* is one of eleven⁵ constants denoting different types of primitive actions. The meanings of *OBJECT*, *SOURCE*, and *DESTINATION* vary with particular action types. The primitive acts and the meanings of these arguments are shown in Table 10.6.

⁵The exact number and list varied over time.

Table 10.5 Axioms of Possession

Nonlogical symbols:

`owner_of(O)`: Function. The fluent of the owner of O over time.

`transfer(O, A)`: Function. The action of transferring possession of O to agent A .

`use_of(E, A, O)`: Event type E constitutes a use of object O by agent A .

`permits(A, E)`: State type. Agent A permits event type E .

POS.1. Causal axiom of transference: If the owner of O transfers O to A_2 , then A_2 becomes the owner.

$$\text{occurs}(I, \text{do}(\text{value_in}(\text{start}(I), \text{owner_of}(O)), \text{transfer}(O, A_2))) \Rightarrow \\ A_2 = \text{value_in}(\text{end}(I), \text{owner_of}(O)).$$

POS.2. Only the owner can transfer possession.

$$\text{occurs}(I, \text{do}(A_1, \text{transfer}(O, A_2))) \Rightarrow \\ A_1 = \text{value_in}(\text{start}(I), \text{owner_of}(O)).$$

POS.3. Frame axiom of transference: Ownership can change only by an act of transference.

$$[\text{precedes}(S_1, S_2) \wedge \\ A_1 = \text{value_in}(S_1, \text{owner_of}(O)) \neq \\ \text{value_in}(S_2, \text{owner_of}(O))] \Rightarrow \\ \exists_{IT, A_2} \text{intersect}(IT, [S_1, S_2]) \wedge \\ \text{occurs}(IT, \text{do}(A_1, \text{transfer}(O, A_2))).$$

POS.4. Use by others is permissible only if the owner allows it.

$$[A_2 \neq \text{value_in}(S, \text{owner_of}(O)) \wedge \text{use_of}(E, A_2, O)] \Rightarrow \\ [\text{permitted}(E, S) \Leftrightarrow \\ \text{true_in}(S, \text{permits}(\text{value_in}(S, \text{owner_of}(O)), E))].$$

The object of an `mbuild` or `mtrans` action is a mental construct. The same problems of referential opacity that we have seen in connection with knowledge and belief therefore arise here as well. We will assume a syntactic theory, in which these objects are strings of characters. CD theory distinguishes between the "mental places" of long-term memory and short-term memory. Thus, the action `mtrans` may be used either for remembering a fact (transferring information from long-term to short-term memory) or for communication (transferring it from one agent to another).

Table 10.6 Primitive Action Types in CD

ptrans	— To move <i>OBJECT</i> from <i>SOURCE</i> to <i>DESTINATION</i> .
propel	— To exert a force on <i>OBJECT</i> in direction <i>DESTINATION</i> .
move	— To move body part <i>OBJECT</i> to <i>DESTINATION</i> .
grasp	— To grasp <i>OBJECT</i> with body part <i>DESTINATION</i> .
ingest	— To consume <i>OBJECT</i> with body orifice <i>DESTINATION</i> .
expel	— To emit <i>OBJECT</i> from body orifice <i>SOURCE</i> .
speak	— To make sounds <i>OBJECT</i> .
mtrans	— To transfer information <i>OBJECT</i> from <i>SOURCE</i> to <i>DESTINATION</i> .
mbuild	— To mentally construct (imagine, consider, come to believe) <i>OBJECT</i> .
attend	— To focus sensory organ <i>OBJECT</i> on object <i>DESTINATION</i> .
atrans	— To transfer possession or control of <i>OBJECT</i> from <i>SOURCE</i> to <i>DESTINATION</i> .

Second, CD defines a number of fluents. These are less central to the theory than the primitive acts, and are often generated *ad hoc*. The following are commonly used:

- $\text{place}(O)$ — The place of physical object O .
- $\text{mloc}(M, P)$ — Mental location P contains mental object M .
- $\text{owner_of}(O)$ — Agent owning object O .
- $\text{health-val}(A)$ — Physical state of agent A . Characterized in terms of an integral measure space, running from “dead” through “sick” to “perfect-health”.
- $\text{happiness}(A)$ — Emotional state of agent A . Characterized in terms of an integral measure space, running from “miserable” to “indifferent” to “content” to “ecstatic”.

We define the function “ $\text{change}(F, D)$ ” to be the event of fluent F changing in direction D , where D is either “pos” or “neg”. For example, “ $\text{change}(\text{happiness}(\text{abel}), \text{pos})$ ” is the event of Abel becoming happier.

Finally, CD defines a number of causal connectives. The arguments of these connectives fall into three categories:

- Tokens of physical events or actions.

Table 10.7 Causal Connectives in CD

- $\text{result}(E, S)$ — Action E brings about state S .
- $\text{enable}(S, E)$ — State S makes action E possible.
- $\text{disable}(S, E)$ — State S makes action E impossible.
- $\text{initiate}(S, M)$ — Action or state S initiates mental state M .
- $\text{reason}(M, E)$ — Mental state E is the reason for action E .
- $\text{instrumental}(E1, E2)$ — Action $E1$ is instrumental to action $E2$.

- Tokens of physical states or state changes.
- Tokens of mental actions, mental states, or changes to mental state.

Table 10.7 shows the causal connectives used in CD. We abbreviate the above category as “actions,” the second as “states,” and the third as “mental states.”

Using these primitives, augmented with primitives defined earlier or defined *ad hoc*, we can now represent some simple narratives. Table 10.8 shows the CD representation of the following story:

Marie saw Jessica pushing a toy wagon. Marie wanted to have the wagon. Marie took the wagon from Jessica. Jessica was upset and started to cry.

10.7 References

Common knowledge: Common knowledge has been studied in a variety of fields, including psychology, economics, game theory, philosophy, and computer science, particularly the theory of distributed systems. [Halpern and Moses 1984] is a particularly important and readable paper. Among other results, it shows that a distributed system cannot attain common knowledge of a fact in a series of communications unless it can attain common knowledge of the fact in a single communication. It also discusses several weakenings of the concept

Table 10.8 CD Representation of a Story

Events and states:

EV1: Jessica pushes a wagon.
 $\text{token_of}(\text{ev1}, \text{action}(\text{jessica}, \text{propel}, \text{wagon1}, \perp, \perp)).$

EV2: Marie looks at Jessica push the wagon.
 $\text{token_of}(\text{ev2}, \text{action}(\text{marie}, \text{attend}, \text{eyes}(\text{marie}), \perp, \text{ev1})).$

MS3: Marie knows that Jessica pushes the wagon.
 $\text{token_of}(\text{ms3}, \text{mloc}(\text{mind_of}(\text{marie}), \text{ev1})).$

MS4: Marie wants to have the wagon.
 $\text{token_of}(\text{ms4}, \text{goal}(\text{marie}, \text{eql}(\text{owner_of}(\text{wagon1}), \text{marie}))).$

EV5: Marie takes the wagon from Jessica.
 $\text{token_of}(\text{ev5}, \text{action}(\text{marie}, \text{atrans}, \text{wagon1}, \text{jessica}, \text{marie})).$

ST6: Marie has the wagon.
 $\text{token_of}(\text{st6}, \text{eql}(\text{owner_of}(\text{wagon1}), \text{marie})).$

MS7: Jessica is upset.
 $\text{token_of}(\text{ms7}, \text{change}(\text{happiness}(\text{jessica}), \text{neg})).$

EV8: Jessica cries.
 $\text{token_of}(\text{ev1}, \text{action}(\text{jessica}, \text{expel}, \text{tears}, \text{eyes}(\text{jessica}), \perp)).$

Causal connections:

$\text{initiate}(\text{ev2}, \text{ms3}).$

$\text{initiate}(\text{ms3}, \text{ms4}).$

$\text{reason}(\text{ms4}, \text{ev5}).$

$\text{result}(\text{ev5}, \text{st6}).$

$\text{initiate}(\text{st6}, \text{ms7}).$

$\text{reason}(\text{ms7}, \text{ev8}).$

Actual occurrence of these events in this order:

$[\text{time_of}(\text{ev1}) \cup \text{time_of}(\text{ev2}) \cup \text{time_of}(\text{ms3}) \cup \text{time_of}(\text{ms4}) \cup \text{time_of}(\text{ev5}) \cup \text{time_of}(\text{st6}) \cup \text{time_of}(\text{ms7}) \cup \text{time_of}(\text{ev8})] \subset \text{real_chronicle}.$

$\text{contains}(\text{time_of}(\text{ev1}), \text{time_of}(\text{ev2})).$

$\text{start}(\text{time_of}(\text{ev2})) \leq \text{start}(\text{time_of}(\text{ms3})) \leq \text{start}(\text{time_of}(\text{ms4})).$

$\text{end}(\text{time_of}(\text{ev1})) \leq \text{start}(\text{time_of}(\text{ev5})).$

$\text{start}(\text{time_of}(\text{ms4})) \leq \text{start}(\text{time_of}(\text{ev5})).$

$\text{meet}(\text{time_of}(\text{ev5}), \text{time_of}(\text{st6})).$

$\text{start}(\text{time_of}(\text{st6})) \leq \text{start}(\text{time_of}(\text{ms7})).$

$\text{contains}(\text{time_of}(\text{ms7}), \text{time_of}(\text{ev8})).$

that are more easily attainable. [Vardi 1988] contains a number of other papers on the subject. Common knowledge has not been extensively studied in the context of commonsense reasoning. As far as I know, the observation that the agent argument of common knowledge is opaque is original here.

Multiagent planning: [Bond and Gasser 1988] contains a number of important papers on multiagent planning. See particularly [Smith and Davis 1981; Georgeff 1983, 1984, 1986; Rosenschein 1982 and Stuart 1985] on cooperative planning, and [Morgenstern 1987] on plans of influence. The use of a theory of plan interaction in narrative understanding is studied in [Wilensky 1978, 1980] and [Dyer 1985].

Communication: High-level AI theories of communication have been extensively studied from many points of view and different degrees of formality. The representation discussed here derives largely from [Morgenstern 1988]. Other important works include [Appelt 1982; Perrault and Allen 1980; Cohen and Levesque 1987; Cohen and Pollack 1987], and [Litman 1985]. The classic philosophical studies of speech acts are [Austin 1961] and [Searle 1969]. See also [Wittgenstein 1958] to get an idea of what a complete commonsense theory of communication and language would have to include.

Ethics: [Sanders 1989] discusses a representation for ethical valuations, possession, and emotions. The formal treatment of ethical values is called *deontic* logic. Deontic logic has been applied to legal reasoning by McCarty [1983]. Philosophical studies of deontic logic include [Von Wright 1968; Prior 1967].

Possession: Many representational systems for narratives (e.g., [Schank 1975]) have used a theory of possession essentially equivalent to the simple one here. Theories that deal with certain aspects of property in greater detail have been developed for use in legal reasoning: see, for example, [McCarty and Sridharan 1980, 1981].

Conceptual dependency: Conceptual dependency was first introduced in [Schank 1969]; this version of the theory had the basic form of actions, with actor, object, source, and destination, but had only one action primitive: "trans." The standard version of CD is presented in [Schank 1975] and [Schank and Abelson 1977].

10.8 Exercises

(Starred problems are more difficult.)

1. * The following is known as the "cheating husbands" problem.

There was once a community with a distinctive system of conventions for dealing with marital infidelity.

- (a) Every woman knew of every man in the community, except her own husband, whether or not he was faithful to his wife.
- (b) A woman who knew that her husband was unfaithful was obliged to shoot him at midnight as soon as possible when she found out. If a woman shot her husband, the fact was immediately common knowledge.
- (c) All women had common knowledge of rules a and b. All women were perfect logical reasoners.

At the time in question, there were k adulterous husbands. One day, at a meeting of the entire community, one of the husbands made the following public statement, witnessed in common by the entire community: "It must be admitted that at least one of us husbands has committed adultery."

k nights later, every adulterous husband was shot by his wife.

Provide a formal justification for the actions of the women.

2. Express the axioms of common knowledge CK.1–CK.5, using the accessibility relation ck_acc .
3. Justify the following inference: Assume (i) that A declares ϕ to B ; (ii) that A and B have common knowledge that A is speaking sincerely; and (iii) that A and B have common knowledge that A knows whether ϕ is true. Infer that A and B have common knowledge of ϕ when the communication is complete.
4. Express the following statements using the primitives of Section 10.3.
 - (a) Mr. Martin said to Mrs. Barrows "I'm sitting in the catbird seat."
 - (b) Mr. Martin told Mr. Barrows that he had taken cocaine.
 - (c) Mrs. Barrows told Mr. Martin to leave her apartment.
 - (d) Mrs. Barrows asked Mr. Martin whether he was drunk.
 - (e) Mr. Fitweiler told Mr. Martin that Mrs. Barrows had told him (Fitweiler) that Martin had taken cocaine.
5. * Validate the following plan of influence [Morgenstern 1988].

John wishes to call Mary, but doesn't know her phone number. However, he does know that Harry knows her number, and that Harry will tell him the number if asked. He therefore constructs the following plan: Ask Harry for Mary's number; wait for Harry to answer the question; dial the number that Harry states.