

Knowledge and Communication: A First-Order Theory

Ernest Davis*
Courant Institute
New York University
davis@cs.nyu.edu

October 7, 2004

Abstract

This paper presents a theory of informative communications among agents that allows a speaker to communicate to a hearer truths about the state of the world; the occurrence of events, including other communicative acts; and the knowledge states of any agent — speaker, hearer, or third parties — any of these in the past, present, or future — and any logical combination of these, including formulas with quantifiers. We prove that this theory is consistent, and compatible with a wide range of physical theories. We examine how the theory avoids two potential paradoxes, and discuss how these paradoxes may pose a danger when this theory are extended.

Keywords: Communication, knowledge, paradox.

1 Introduction

In constructing a formal theory of communications between agents, the issue of expressivity enters at two different levels: the scope of what can be said *about* the communications, and the scope of what can be said *in* the communications. Other things being equal, it is obviously desirable to make both of these as extensive as possible. Ideally, a theory should allow a speaker to communicate to a hearer truths about the state of the world; the occurrence of events, including other communicative acts; the knowledge states of any agent — speaker, hearer, or third parties; any of these in the past, present, or future; and any logical combination of these. This paper presents a theory that achieves pretty much that.

A few examples of what can be expressed, together with their formal representation:

1. Alice tells Bob that all her children are asleep.

$$\begin{aligned} \exists_Q \text{ occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, Q)), s_0, s_1) \wedge \\ \forall_S \text{ holds}(S, Q) \Leftrightarrow \\ [\forall_C \text{ holds}(S, \text{child}(C, \text{alice})) \Rightarrow \text{holds}(S, \text{asleep}(C))]. \end{aligned}$$

2. Alice tells Bob that she doesn't know whether he locked the door.

*The research reported in this paper was supported in part by NSF grant IIS-0097537. The work described here comes out of and builds upon a project done in collaboration with Leora Morgenstern, stemming from a benchmark problem that she proposed.

$$\begin{aligned}
& \exists_Q \text{occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, Q)), s_0, s_1) \wedge \\
& \forall_S \text{holds}(S, Q) \Leftrightarrow \\
& \quad [\exists_{S_A} \text{k_acc}(\text{alice}, S, S_A) \wedge \\
& \quad \quad \exists_{S_{1A}, S_{2A}} S_{1A} < S_{2A} < S_A \wedge \\
& \quad \quad \text{occurs}(\text{do}(\text{bob}, \text{lock_door}), S_{1A}, S_{2A})] \wedge \\
& \quad [\exists_{S_A} \text{k_acc}(\text{alice}, S, S_A) \wedge \\
& \quad \quad \neg \exists_{S_{1A}, S_{2A}} S_{1A} < S_{2A} < S_A \wedge \\
& \quad \quad \text{occurs}(\text{do}(\text{bob}, \text{lock_door}), S_{1A}, S_{2A})].
\end{aligned}$$

3. Alice tells Bob that if he finds out who was in the kitchen at midnight, then he will know who killed Colonel Mustard. (Note: The interpretation below assumes that exactly one person was in the kitchen at midnight.)

$$\begin{aligned}
& \exists_Q \text{occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, Q)), s_0, s_1) \wedge \\
& \forall_S \text{holds}(S, Q) \Leftrightarrow \\
& \quad \forall_{S_2} [S_2 > S \wedge \\
& \quad \quad \exists_{PK} \forall_{S_{2A}} \text{k_acc}(\text{bob}, S_2, S_{2A}) \Rightarrow \\
& \quad \quad \quad \exists_{S_{3A}} S_{3A} < S_{2A} \wedge \text{midnight}(\text{time}(S_{3A})) \wedge \text{holds}(S_{3A}, \text{in}(PK, \text{kitchen}))] \Rightarrow \\
& \quad [\exists_{PM} \forall_{S_{2B}} \text{k_acc}(\text{bob}, S_2, S_{2B}) \Rightarrow \\
& \quad \quad \exists_{S_{3B}, S_{4B}} S_{3B} < S_{4B} < S_{2B} \wedge \text{occurs}(\text{do}(PM, \text{murder}(\text{mustard})), S_{3B}, S_{4B})].
\end{aligned}$$

4. Alice tells Bob that no one had ever told her she had a sister.

$$\begin{aligned}
& \exists_Q \text{occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, Q)), s_0, s_1) \wedge \\
& \forall_S \text{holds}(S, Q) \Leftrightarrow \\
& \quad \neg \exists_{S_2, S_3, Q_1, P_1} S_2 < S_3 < S \wedge \\
& \quad \quad \text{occurs}(\text{do}(P_1, \text{inform}(\text{alice}, Q_1)), S_2, S_3) \wedge \\
& \quad \quad \forall_{S_X} \text{holds}(S_X, Q_1) \Rightarrow \exists_{P_2} \text{holds}(S_X, \text{sister}(P_2, \text{alice})).
\end{aligned}$$

5. Alice tells Bob that he has never told her anything she didn't already know.

$$\begin{aligned}
& \exists_Q \text{occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, Q)), s_0, s_1) \wedge \\
& \forall_S \text{holds}(S, Q) \Leftrightarrow \\
& \quad \forall_{S_2, S_3, Q_1} \\
& \quad \quad [S_2 < S_3 \leq S \wedge \\
& \quad \quad \quad \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, Q_1)), S_2, S_3)] \Rightarrow \\
& \quad \quad \forall_{S_{2A}} \text{k_acc}(\text{alice}, S_2, S_{2A}) \Rightarrow \text{holds}(S_{2A}, Q_1).
\end{aligned}$$

These representations works as follows: The expression “do(AS ,inform(AH , Q))” denotes the action of speaker AS informing AH that fluent Q holds in the current situation. The content Q here is a *generalized fluent*, that is, a property of situations / possible worlds. Simple fluents are defined by ground terms, such as “in(mustard,kitchen).” In more complex cases, the fluent Q is characterized by a formula “ $\forall_S \text{holds}(S, Q) \Leftrightarrow \alpha(S)$ ” where α is some formula open in S . (Equivalently, Q could be defined using the lambda expression $Q = \lambda(S)\alpha(S)$.)

The above examples illustrate many of the expressive features of our representation:

- Example 1 shows that the content of a communication may be a quantified formula.
- Example 2 shows that the content of a communication may refer to knowledge and ignorance of past actions.

- Example 3 shows that the content of a communication may be a complex formula involving both past and present events and states of knowledge.
- Examples 4 and 5 show that the content of a communication may refer to other communications. They also show that the language supports quantification over agents and over the content of a communication, and thus allows the content to be partially characterized, rather than fully specified.

If we wish to reason about such informative actions — e.g. to be sure that they can be executed — then we must be sure, among other conditions, that the fluent denoting the content of the action exists. This requires a comprehension axiom that asserts that such a fluent exists for *any* such formula α . Comprehension axioms often run the risk of running into analogues of Russell’s paradox, but this one turns out to be safe. We will discuss two paradoxes that look dangerous for this theory, but the theory succeeds in side-stepping these. One of these is the well-known “unexpected hanging” paradox. To make sure that there are no further paradoxes in hiding that might be more destructive, we prove that our theory is consistent, and compatible with a wide range of physical theories.

We should note at the outset the limitations of our theory. The theory deals only with informative acts (and not, for example, with requests) and assumes that the following conditions are true and universally known: If *AS* communicates *Q* to *AH*, then

1. *AS* knows that *Q* is true at the time that he initiates the communication.
2. From the time that he initiates the communication, *AS* knows that he is carrying out a communication; he knows that the content is *Q*; and he knows that the recipient is *AH*.
3. Similarly, when the communication is complete, *AH* knows that he has received a communication; he knows that the content was *Q*; and he knows that the sender was *AS*.
4. When the communication is complete, *AS* knows that the communication is complete and *AH* knows the time at which the communication was initiated.

The paradigmatic example of a form of communication satisfying conditions 2, 3, and 4 is direct speech. Another example could be mail, assuming that

- All messages are time-stamped with the time of sending, and signed by the sender.
- There is a universally known maximal delay *D* between the time of sending and the time of receiving a message. (“Receiving” here means the time when the hearer reads the message, not the time that it arrives in his mailbox.)

In this case, if we define a communication to be “complete” at the time of sending plus *D*, then the above conditions are met.

Many aspects of the theory can be applied to communications that do not meet condition (4), but I have not been able to find a plausible axiomatization of this more general case that I can prove to be consistent. Also, I cannot prove that the theory is consistent unless time is taken to be discrete. These are discussed further in section 8.

The paper proceeds as follows: Section 2 reviews the theories of time and of knowledge, which are not new here. Section 3 presents our language and axioms of communication. Section 4 illustrates the power of the theory by showing how it supports three example inferences. Section 5 describes two apparent paradoxes — a paradox analogous to Russell’s paradox and the “unexpected hanging” paradox — and explains why these do not cause inconsistencies in the theory. Section 6 gives the

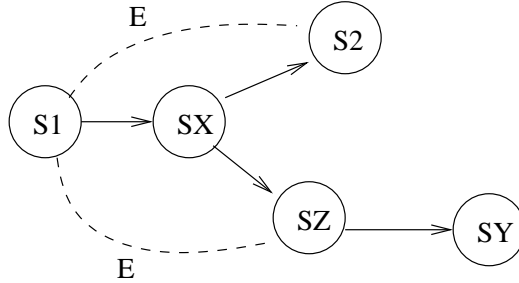


Figure 1: Axiom T.9

statement of Theorems 1 and 2, which assert that the theory is internally consistent and compatible with a wide range of physical theories. Section 7 discusses related work. Section 8 discusses open problems and summarizes our conclusions. Appendix A gives the proofs of theorems 1 and 2.

2 Framework

We use a situation-based, branching theory of time; an interval-based theory of multi-agent actions; and a possible-worlds theory of knowledge. This is all well known, so the description below is brief.

2.1 Time and Action

We use a situation-based theory of time. Time can be either continuous¹ or discrete, but it must be *branching*, like the situation calculus. The branching structure is described by the partial ordering “ $S1 < S2$ ”, meaning that there is a timeline containing $S1$ and $S2$ and $S1$ precedes $S2$. It is convenient to use the abbreviations “ $S1 \leq S2$ ” and “ $\text{ordered}(S1, S2)$.” The predicate “ $\text{holds}(S, Q)$ ” means that fluent Q holds in situation S .

Each agent has, in various situations, a choice about what action to perform next, and the time structure includes a separate branch for each such choice. Thus, the statement that action E is feasible in situation S is expressed by asserting that E occurs from S to $S1$ for some $S1 > S$.

Following (McDermott 1982), actions are represented as occurring over an interval; the predicate $\text{occurs}(E, S1, S2)$ states that action E occurs starting in $S1$ and ending in $S2$. However, the whole theory could be recast without substantial change into the situation calculus extended to permit multiple agents, after the style of (Reiter, 2001).

Table ?? shows the axioms of our temporal theory. Throughout this paper, we use a sorted first-order logic with equality, where the sorts of variables are indicated by their first letter. The sorts are clock-times (T), situations (S), Boolean fluents (Q), actions (E), agents (A), and actionals (Z). (The examples at the beginning of this paper use some terms of other sorts *ad hoc*; these are self-explanatory.) An *actional* is a characterization of an action without specifying the agent. For example, the term “ $\text{puton}(\text{blocka}, \text{table})$ ” denotes the actional of someone putting block A on the table. The term “ $\text{do}(\text{john}, \text{puton}(\text{blocka}, \text{table}))$ ” denotes the action of John putting block A on the table. Free variables in a formula are assumed to be universally quantified.

Our theory does not include a representation of what *will* happen from a given situation as

¹As will be discussed below, I cannot prove the theory consistent for continuous theories of time except in special cases; however, nothing in the form of the representation or in the axioms is inherently unusable in or inconsistent with a continuous model of time.

Primitives:

$T1 < T2$ — Time $T1$ is earlier than $T2$.

$S1 < S2$ — Situation $S1$ precedes $S2$, on the same time line. (We overload the $<$ symbol.)

$\text{time}(S)$ — Function from a situation to its clock time.

$\text{holds}(S, Q)$ — Fluent Q holds in situation S .

$\text{occurs}(E, S1, S2)$ — Action E occurs from situation $S1$ to situation $S2$.

$\text{do}(A, Z)$ — Function. The action of agent A doing actional Z .

Definitions:

TD.1 $S1 \leq S2 \equiv S1 < S2 \vee S1 = S2$.

TD.2 $\text{ordered}(S1, S2) \equiv$
 $S1 < S2 \vee S1 = S2 \vee S2 < S1$.

TD.3 $\text{feasible}(E, S) \Leftrightarrow \exists S2 \text{ occurs}(E, S, S2)$.

Axioms:

T.1 $T1 < T2 \vee T2 < T1 \vee T1 = T2$.

T.2 $\neg[T1 < T2 \wedge T2 < T1]$.

T.3 $T1 < T2 \wedge T2 < T3 \Rightarrow T1 < T3$.
 (Clock times are linearly ordered)

T.4 $S1 < S2 \wedge S2 < S3 \Rightarrow S1 < S3$. (Transitivity)

T.5 $(S1 < S \wedge S2 < S) \Rightarrow \text{ordered}(S1, S2)$.
 (Forward branching)

T.6 $S1 < S2 \Rightarrow \text{time}(S1) < \text{time}(S2)$.
 (The ordering on situations is consistent with the orderings of their clock times.)

T.7 $\forall S, T1 \exists S1 \text{ ordered}(S, S1) \wedge \text{time}(S1)=T1$.
 (Every time line contains a situation for every clock time.)

T.8 $\text{occurs}(E, S1, S2) \Rightarrow S1 < S2$.
 (Events occur forward in time.)

T.9 $[\text{occurs}(E, S1, S2) \wedge S1 < SX < S2 \wedge SX < SY] \Rightarrow$
 $\exists SZ \text{ ordered}(SY, SZ) \wedge \text{occurs}(E, S1, SZ)$.
 (If action E starts to occur on the time line that includes SY , then it completes on that time line. (Figure ??))

Table 1: Temporal Axioms

opposed to what *can* happen. This will be important in our discussion of the paradox of the unexpected hanging.

2.2 Knowledge

As first proposed by Moore (1980,1985a) and widely used since, knowledge is represented by identifying temporal situations with epistemic possible worlds and positing a relation of knowledge accessibility between situations. The relation $k_acc(A, S, SA)$ means that situation SA is accessible from S relative to agent A 's knowledge in S ; that is, as far as A knows in S , the actual situation could be SA . The statement that A knows ϕ in S is represented by asserting that ϕ holds in every situation that is knowledge accessible from S for A . As is well known, this theory enables the expression of complex interactions of knowledge and time; one can represent both knowledge about change over time and change of knowledge over time.

Again following Moore (1985a), the state of agent A knowing *what something is* is expressed by using a quantifier of larger scope than the universal quantification over accessible possible worlds. For example, the statement, "In situation $s1$, John knows who the President is" is expressed by asserting that there exists a unique individual who is the President in all possible worlds accessible for John from $s1$.

$$\exists X \forall_{S1A} k_acc(john,s1,S1A) \Rightarrow \text{holds}(S1A,president(X)).$$

For convenience, we posit an S5 logic of knowledge; that is, the knowledge accessibility relation, restricted to a single agent, is in fact an equivalence relation on situations. This is expressed in axioms K.1, K.2, and K.3 in table ???. Three important further axioms govern the relation of time and knowledge.

- K.4. Axiom of memory: If A knows ϕ in S , then in any later situation, he remembers that he knew ϕ in S .
- K.5. A knows all the actions that he has begun, both those that he has completed and those that are ongoing. That is, he knows a *standard identifier* for these actions; if Bob is dialing (212) 998-3123 on the phone, he knows that he is dialing (212) 998-3123 but he may not know that he is calling Ernie Davis. At any time, A knows what actions are feasible for him now.
- K.6 Knowledge accessibility relations do not cross in the time structure. I have not found any natural expression of this axiom, but certainly a structure that violated it would be a very odd one. (Figure ??.) In a discrete theory of time, axiom K.6 is a consequence of the axiom of memory K.4, as we shall show in theorem 3 below. (Knowledge accessibility relations that violate this condition have sometimes been used in the literature for agents who do not satisfy the axiom of memory.)

The theory includes a forms of common knowledge, restricted to two agents. Agents $A1$ and $A2$ have *shared knowledge* of ϕ if they both know ϕ , they both know that they both know ϕ and so on. We represent this by defining a further accessibility relation, " $sk_acc(A1, A2, S, SA)$ " (SA is accessible from S relative to the shared knowledge of $A1$ and $A2$). This is defined as the transitive closure of links of the form $k_acc(A1, \cdot, \cdot)$ together with links of the form $k_acc(A2, \cdot, \cdot)$. (Of course, transitive closure cannot be exactly defined in a first-order theory; axioms K.7 and K.8 define an approximation that is adequate for our purposes.)

Primitives:

$k_acc(A, SA, SB)$ — SB is accessible from SA relative to A 's knowledge in SA .

$sk_acc(A1, A2, SA, SB)$ — SB is accessible from SA relative to the shared knowledge of $A1$ and $A2$ in SA .

Axioms

K.1 $\forall_{A,SA} k_acc(A, SA, SA)$.

K.2 $k_acc(A, SA, SB) \Rightarrow k_acc(A, SB, SA)$

K.3 $k_acc(A, SA, SB) \wedge k_acc(A, SB, SC) \Rightarrow k_acc(A, SA, SC)$.

(K.1 through K.3 suffice to ensure that the knowledge of each agent obeys an S5 logic: what he knows is true, if he knows ϕ he knows that he knows it; if he doesn't know ϕ , he knows that he doesn't know it.)

K.4 $[k_acc(A, S2A, S2B) \wedge S1A < S2A] \Rightarrow \exists_{S1B} S1B < S2B \wedge k_acc(A, S1A, S1B)$.

(Axiom of memory: If agent A knows ϕ at any time, then at any later time he knows that ϕ was true.)

K.5 $[\text{occurs}(\text{do}(A, Z), S1A, S2A) \wedge S1A \leq SA \wedge \text{ordered}(SA, S2A) \wedge k_acc(A, SA, SB)] \Rightarrow \exists_{S1B, S2B} \text{occurs}(\text{do}(A, Z), S1B, S2B) \wedge S1B \leq SB \wedge$

$[S2A < SA \Rightarrow S2B < SB] \wedge$

$[S2A = SA \Rightarrow S2B = SB] \wedge$

$[SA < S2A \Rightarrow SB < S2B] \wedge$

$[S1A = SA \Rightarrow S1B = SB]$

(An agent knows which actions he has completed, which actions he has begun, and which actions are now feasible.)

K.6 $\neg \exists_{A, S1A, S1B, S2A, S2B}$

$S1A < S2A \wedge S1B < S2B \wedge k_acc(A, S1A, S2B) \wedge k_acc(A, S2A, S1B)$.

(Knowledge accessibility links do not cross in the time structure (Figure ??).)

K.7 $sk_acc(A1, A2, SA, SB) \Leftrightarrow$

$[k_acc(A1, SA, SB) \vee k_acc(A2, SA, SB) \vee$

$sk_acc(A1, A2, SB, SA) \vee$

$sk_acc(A2, A1, SA, AB) \vee$

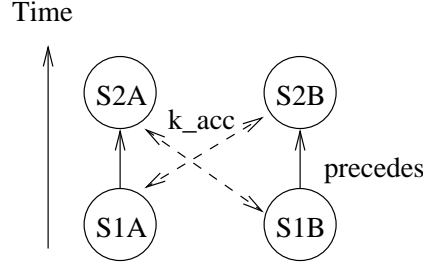
$\exists_{SC} sk_acc(A1, A2, SA, SC) \wedge sk_acc(A1, A2, SC, SB)]$.

Definition of sk_acc as a equivalence relation, symmetric in $A1, A2$, that includes the k_acc links for the two agents $A1, A2$.

K.8 (Induction from k_acc links to sk_acc links.) Let $\Phi(S)$ be a formula with a free situational variable S . Then the closure of the formula

$$[\forall_{AS, AH} [[\forall_{SA, SB} \phi(SA) \wedge k_acc(AS, SA, SB) \Rightarrow \phi(SB)] \wedge [\forall_{SA, SB} \phi(SA) \wedge k_acc(AH, SA, SB) \Rightarrow \phi(SB)]] \Rightarrow [\forall_{SA, SB} \phi(SA) \wedge sk_acc(AS, AH, SA, SB) \Rightarrow \phi(SB)].$$

Table 2: Axioms of Knowledge



Axiom K.6 prohibits this structure.

Figure 2: Axiom K.6

3 Communication

We now introduce the function “inform”, taking two arguments, a agent AH and a fluent Q . The term “ $\text{inform}(AH, Q)$ ” denotes the actional of informing AH that Q ; the term “ $\text{do}(AS, \text{inform}(AH, Q))$ ” thus denotes the action of speaker AS informing AH that Q . Our theory here treats “ $\text{do}(AS, \text{inform}(AH, Q))$ ” as a primitive actions; in a richer theory, it would be viewed as an illocutionary description of an underlying locutionary act (not here represented) — the utterance or writing or broadcasting of a physical signal.

We also add a second actional “ $\text{communicate}(AH)$ ”. This alternative characterization of a communicative act, which specifies the hearer but not the content of the communication, enables us to separate out *physical* constraints on a communicative act from *contentive* constraints. Thus, we allow a purely physical theory to put constraints on the occurrence of a communication, or even to posit physical effects of a communication, but these must be independent of the information content of the communication.

We posit five axioms of communication, summarized in table ?? . Some of these are straightforward; others much less so. We discuss them below in increasing order of complexity. We also put forward a sixth axiom, a frame axiom for ignorance, but its status is much more dubious, for reasons that we will discuss.

3.1 Relation between informing and communication

Axiom I.1: Any inform act is a communication.

$$\begin{aligned} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) &\Rightarrow \\ \text{occurs}(\text{do}(AS, \text{communicate}(AH)), S1, S2). \end{aligned}$$

Axiom I.2: If a speaker AS can communicate with a hearer AH , then AS can inform AH of some specific Q if and only if A knows that Q holds at the time he begins speaking.

$$\begin{aligned} \text{feasible}(\text{do}(AS, \text{communicate}(AH)), S1) &\Rightarrow \\ [\forall Q \text{ feasible}(\text{do}(AS, \text{inform}(AH, Q)), S1) &\Leftrightarrow \\ [\forall_{S1A} \text{ k_acc}(AS, S1, S1A) \Rightarrow \text{holds}(S1A, Q)]] \end{aligned}$$

By virtue of these two axioms, the preconditions for an AS informing AH that Q are just that it is feasible for AS to communicate to AH and that AS knows that Q is true. The content Q may not affect the feasibility in any other way. Axiom I.1 further guarantees that any other physical constraints over communications, such as the duration of a communication or its physical effects must apply also to inform acts; that is, that the physical characteristics of any inform act must

I.1 Any inform act is a communication.

$$\text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \Rightarrow \\ \text{occurs}(\text{do}(AS, \text{communicate}(AH)), S1, S2).$$

I.2. If a speaker AS can communicate with a hearer AH , then AS can inform AH of some specific Q if and only if A knows that Q holds at the time he begins speaking.

$$\text{feasible}(\text{do}(AS, \text{communicate}(AH)), S1] \Rightarrow \\ [\forall_Q \text{feasible}(\text{do}(AS, \text{inform}(AH, Q)), S1) \Leftrightarrow \\ [\forall_{S1A} \text{k_acc}(AS, S1, S1A) \Rightarrow \text{holds}(S1A, Q)]]$$

I.3. If AS informs AH of Q from $S1$ to $S2$, then in $S2$, AH knows that AS has informed him of Q .

$$\forall_{S1, S2, S2A} [\text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \wedge \text{k_acc}(AH, S2, S2A)] \Rightarrow \\ \exists_{S1A} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1A, S2A) \wedge \text{k_acc}(AH, S1, S1A).$$

I.4. If AS informs AH of $Q1$ over $[S1, S2]$ and the shared knowledge of AS and AH in $S1$ implies that $\text{holds}(S1, Q1) \Leftrightarrow \text{holds}(S1, Q2)$, then AS has also informed AH of $Q2$ over $[S1, S2]$. Conversely, the two actions “ $\text{do}(AS, \text{inform}(AH, Q1))$ ” and “ $\text{do}(AS, \text{inform}(AH, Q2))$ ” co-occur only if $Q1$ and $Q2$ are related in this way.

$$\text{occurs}(\text{do}(AS, \text{inform}(AH, Q1)), S1, S2) \Rightarrow \\ [\text{occurs}(\text{do}(AS, \text{inform}(AH, Q2)), S1, S2) \Leftrightarrow \\ [\forall_{S1A} \text{sk_acc}(AS, AH, S1, S1A) \Rightarrow \\ [\text{holds}(S1A, Q1) \Leftrightarrow \text{holds}(S1A, Q2)]]]]$$

I.5. Axiom of comprehension: any property of situations that can be stated in the language is a fluent.

Let $\alpha(S)$ be a first-order formula with exactly one free variable S of sort “situation”. (α may have other free variables of other sorts.) Then the closure of the following formula is an axiom:

$$\exists_Q \forall_S \text{holds}(S, Q) \Leftrightarrow \alpha(S).$$

(The closure of a formula β is β scoped by universal quantifications of all its free variables.)

I.6 Frame axiom for ignorance. See discussion in text below.

Table 3: Axioms of Communication

be consistent with the physical constraints on communications. These axioms do not rule out the possibility that the content could affect other physical aspects of the inform act — for example, that a complex content takes longer to communicate than a simple content — but I have not shown that any such constraints lead to a consistent theory.

Note that axiom I.2 requires, conversely, that any fluent Q that is known to be true can be communicated; that is, there is a branch in the time structure corresponding to the communication of Q .

3.2 Epistemic effect of communication

Since we require the strong conditions mentioned in section 1, we can posit the following axiom:²

Axioms I.3: If AS informs AH of Q from $S1$ to $S2$, then in $S2$, AH knows that AS has informed him of Q .

$$\begin{aligned} \forall_{S1, S2, S2A} [\text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \wedge \\ \text{k_acc}(AH, S2, S2A)] \Rightarrow \\ \exists_{S1A} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1A, S2A) \wedge \text{k_acc}(AH, S1, S1A) \end{aligned}$$

Lemmas 3.1 and 3.2 are important consequences of I.3 together with the preceding axioms:

Lemma 3.1: If AS informs AH of Q then, when the communication is in complete, the AS and AH have shared knowlege that the communication has taken place.

$$\begin{aligned} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \wedge \text{sk_acc}(AS, AH, S2, S2A) \Rightarrow \\ \exists_{S1A} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1A, S2A) \end{aligned}$$

Proof: By K.5, AS knows when he has completed a communication.

$$\begin{aligned} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \wedge \text{k_acc}(AS, S2, S2A) \Rightarrow \\ \exists_{S1A} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1A, S2A) \end{aligned}$$

By I.3, AH knows when he has received a communication.

$$\begin{aligned} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \wedge \text{k_acc}(AH, S2, S2A) \Rightarrow \\ \exists_{S1A} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1A, S2A) \end{aligned}$$

Choosing $\Phi(S)$ to be the formula “ $\text{occurs}(\text{do}(AS, \text{inform}(AH, Q)))$ ”, the formula in Lemma 3.1 then follows from axiom K.8.

Lemma 3.2: If AS informs AH of Q then, when the communication is in complete, the AS and AH have shared knowlege that Q was true when the communication began.

$$\begin{aligned} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1, S2) \wedge \text{sk_acc}(AS, AH, S2, S2A) \Rightarrow \\ \exists_{S1A} \text{occurs}(\text{do}(AS, \text{inform}(AH, Q)), S1A, S2A) \wedge \text{holds}(S1A, Q) \end{aligned}$$

Proof:

Let $as, ah, q, s0, s1, s2a$ satisfy the left side of the above implication.

By I.3 there exists $s1a$ such that $\text{occurs}(\text{do}(as, \text{inform}(ah, q)), s1a, s2a)$.

By K.1, $\text{k_acc}(as, s1a, s1a)$.

By I.2, $\text{holds}(s1a, q)$.

²The statement of this axiom in the KR-2004 paper (Davis, 2004) was not correct.

3.3 Axiom of comprehension

The axiom of comprehension states that there is a fluent corresponding to any property of situations definable in the language. The content of this axiom therefore depends on the overall language \mathcal{L} . We state this as an axiom schema i.e. an infinite set of axioms.

Axiom I.5: The comprehension axiom for fluents in a language \mathcal{L} is this: Let $\alpha(S)$ be a first-order formula in \mathcal{L} with exactly one free variable S of sort “situation”. (α may have other free variables of other sorts.) Then the closure of the following formula is an axiom:

$$\exists Q \forall S \text{ holds}(S, Q) \Leftrightarrow \alpha(S).$$

(The closure of a formula β is β scoped by universal quantifications of all its free variables.)

Let us first discuss the significance of free variables in the formula α . The reason to allow free variables that are not situations is to deal with examples such as the following: We want to be able to posit that a speaker can say, for example, that some specific block is either red or blue without requiring that the language \mathcal{L} have a constant symbol for each block, or even a formula that uniquely identifies each block.³

This axiom achieves this. We choose $\alpha(S)$ to be the formula “ $\text{holds}(S, \text{red}(X)) \vee \text{holds}(S, \text{blue}(X))$ ”. The axiom schema then state

$$\forall X \exists Q \forall S \text{ holds}(S, Q) \Leftrightarrow \text{holds}(S, \text{red}(X)) \vee \text{holds}(S, \text{blue}(X))$$

That is, for every object X there is a fluent Q that corresponds to the situations in which X is either red or blue.

The reason to exclude formulas that have other situational free variables in addition to S is that it doesn’t seem to mean anything to have this kind of *de re* reference to situations. A situation is meaningful only in relation to the current situation; there is no other way to meaningfully refer to a situation. It may be noted that the consistency proof for the theory (theorem 1 below) does not depend on this restriction.

The content of the comprehension axiom depends on the overall language \mathcal{L} . In general, one supposes that the language \mathcal{L} will contain many domain and problem specific symbols beyond those that are used in the axioms enumerated here. Theorem 1 shows that these axioms are consistent when \mathcal{L} is a physical language augmented with the symbols from the theory of knowledge and communication described here. In (Davis and Morgenstern, 2004) we consider a language that includes also agent commitments and requests. In that setting, the above formulation of the axiom turns out to be too strong; we have to limit the comprehension axiom to apply only to formulas that do not include symbols describing commitment and requests.

In view of this comprehension axiom, axiom K.8 could be restated as a single axiom (rather than an axiom schema) as follows;

$$\text{K.8.A } \forall_{Q, AS, AH} [[\forall_{S, SA} \text{ holds}(S, Q) \wedge \text{k_acc}(AS, S, SA) \Rightarrow \text{holds}(SA, Q)] \wedge \\ [\forall_{S, SA} \text{ holds}(S, Q) \wedge \text{k_acc}(AH, S, SA) \Rightarrow \text{holds}(SA, Q)]] \Rightarrow \\ \forall_{S, SA} \text{ holds}(S, Q) \wedge \text{sk_acc}(AS, AH, S, SA) \Rightarrow \text{holds}(SA, Q).$$

However we did not use this formulation originally because we did not want K.8 to be dependent on I.5.

³You might well ask, “If you can’t refer to the block in \mathcal{L} , how is the speaker talking about it?” Perhaps he is pointing. Perhaps he is using a slightly richer language with more constant symbols. It is not a very important point, but it does make the theory more elegant and easier to use if one assumes that a speaker can refer *de re* to any entity other than a situation.

3.4 Independent actions

In a temporal representation, like ours, that permits the concurrent execution of actions, it does not suffice just to describe what actions can be executed; one must also, to greater or lesser extent, describe what combinations of actions can be executed concurrently. At the minimum, if two actions are independent, it should be possible to execute the one without the other. In the case of “inform” acts, the natural axiom would be that, if AS knows ϕ , then he can choose to carry out the single act of informing AH of ϕ and not doing anything else. One might suppose that this could be expressed in the following two axioms:

WRONG.1 $\text{feasible}(\text{do}(AS, \text{inform}(AH, Q)), S_0) \Rightarrow$
 $\exists_{S_2} \text{occurs}(\text{do}(AS, Z), S_1, S_2) \Leftrightarrow Z = \text{do}(AS, \text{inform}(AH, Q)).$

WRONG.2 $\text{do}(AS_1, \text{inform}(AH_1, Q_1)) = \text{do}(AS_2, \text{inform}(AH_2, Q_2)) \Rightarrow$
 $AS_1 = AS_2 \wedge AH_1 = AH_2 \wedge Q_1 = Q_2.$

However, as my labels subtly suggest,⁴ this is not an acceptable formulation. In fact, as we shall show in section ??, these are inconsistent with the axiom of comprehension I.5.

The problem, intuitively, is this: The comprehension axiom asserts that there exists a fluent for every set of situations; axiom WRONG.1 asserts that there exists a separate branch in time for every fluent. Therefore, if you try to construct a model of these axioms combined, you first have to construct all sets of situations; then add branches for each of these, which gives a whole bunch more resultant situations; these in turn generate vast numbers of new sets of situations . . . There is no way to make this construction converge. (I’m being a little loose here, but one can make this tight. The decisive proof that this can’t be made to work is the “mised” paradox of section ??.)

Therefore, we have to weaken axiom WRONG.1.⁵ The approach we take is as follows: In general, it is only necessary to distinguish an occurrence of action A1 from an occurrence of action A2 if they have different causal consequences. For instance, in the blocks world, if all you are interested in is the position of blocks, then all that matters in discriminating actions is the ending position of the block being moved; the trajectory through which it moves is immaterial.

Now, in the case of informative acts, the causal consequence of concern is the effect on knowledge states. Assuming axiom I.3, the main effect of AS informing AH of Q is that, when the communication is complete, AS and AH have shared knowledge that Q held at the beginning of the communication. Therefore, if Q_1 and Q_2 are two informative contents such that the effects on the shared knowledge of AS and AH following a communication of Q_1 from AS to AH are the same as those effects following a communication of Q_2 , then we can treat the communication of Q_1 and the communication of Q_2 as *the same action*; they, so to speak, attain the same end state via different trajectories. And a sufficient condition to ensure this is that AS and AH have shared knowledge at the *start* of the communication that Q_1 and Q_2 are equivalent.

For example, if Jack and Jane share the knowledge that George Bush is the President and that 1600 Pennsylvania Avenue is the White House, then the action of Jack informing Jane that Bush is at the White House is identical to the act of Jack informing Jane that the President is at 1600 Pennsylvania Avenue. If they do not share this knowledge, then these two acts are different.

⁴One thing I have learned in twenty years of teaching is that, if you write down a wrong formula on the blackboard for purposes of discussion, you have to label it WRONG in large letters. Otherwise, students copy it into their notebooks . . . Similarly, if someone is skimming through this paper looking for formal axioms, I do not want him to use these.

⁵Weakening axiom WRONG.2 does not work. In fact, WRONG.2 ends up being true in the model we will construct, but its truth won’t actually matter much once we have correctly reformulated WRONG.1.

This, then, is our axiom: If it is feasible in $S1$ for AS to inform AH of $Q1$, then there exists a branch of the time structure in which the only informative action that AS starts in $S1$ are those that are “equivalent” to $Q1$ in the above sense.⁶

$$\begin{aligned} \text{I.4: } & \text{occurs}(\text{do}(AS, \text{inform}(AH, Q1)), S1, S2) \Rightarrow \\ & [\text{occurs}(\text{do}(AS, \text{inform}(AH, Q2)), S1, S2) \Leftrightarrow \\ & [\forall_{S1A} \text{sk_acc}(AS, AH, S1, S1A) \Rightarrow \\ & [\text{holds}(S1A, Q1) \Leftrightarrow \text{holds}(S1A, Q2)]]]] \end{aligned}$$

As we shall see in section ??, in a model of discrete time this is sufficient to avoid the contradiction.

Note: The above axiom is not sufficient to rule out models in which the informative actions of one agent are linked to the concurrent actions of another agent. The easiest way to insure independence between agents is to posit an axiom of “anti-synchrony” that no two agents begin two actions at the same time (Reiter, 2001).

$$\text{T.10 } \text{occurs}(\text{do}(A1, Z1), S1, S2) \wedge \text{occurs}(\text{do}(A2, Z2), S1, S3) \Rightarrow A1 = A2.$$

However, since this axiom is part of the physical theory, and not all physical theories may wish to use it, we have not made it part of our standard set of temporal axioms.

Two alternative formulations of axiom I.4 should be mentioned. We can modify I.4 to read that communicating $Q1$ and $Q2$ co-occur just if they coincide over all situations of the same time as the beginning of the situation.

$$\begin{aligned} \text{I.4.A: } & \text{occurs}(\text{do}(AS, \text{inform}(AH, Q1)), S1, S2) \Rightarrow \\ & [\text{occurs}(\text{do}(AS, \text{inform}(AH, Q2)), S1, S2) \Leftrightarrow \\ & [\forall_{S1A} \text{time}(S1A) = \text{time}(S1) \Rightarrow \\ & [\text{holds}(S1A, Q1) \Leftrightarrow \text{holds}(S1A, Q2)]]]] \end{aligned}$$

The consistency proof in Appendix A requires only a small modification to deal with this new version. However, this version seems to me harder to justify than the previous version.

A second alternative, which is in effect equivalent to axiom I.4.A, is to use the axioms WRONG.1 and WRONG.2 and modify the comprehension axiom to state that there is a fluent corresponding to every property of situations at some particular time T :

I.5.B: Let $\alpha(S)$ be a first-order formula in \mathcal{L} with exactly one free variable S of sort “situation”. (α may have other free variables of other sorts.) Then the closure of the following formula is an axiom:

$$\forall_T \exists_Q \forall_S \text{holds}(S, Q) \Leftrightarrow \alpha(S) \wedge \text{time}(S) = T.$$

3.5 The frame inference

Finally, it would be desirable to carry out the frame inference over knowledge and ignorance.

The frame axiom over knowledge is just the axiom of memory, axiom K.4; if A knows in S that ϕ is true, then he remembers in all later situations that ϕ was true. Since we have no actions or events

⁶This is slightly more general than the formulation given in the KR-2004 paper (Davis, 2004).

that cause forgetting, this simple formulation suffices. Note that “knowing ϕ ” is represented as “all worlds in which ϕ is false are inaccessible.” Hence preserving knowledge amounts to saying that if two worlds are inaccessible one from the other, any of their descendants are likewise inaccessible one from the other.

The frame axiom over ignorance is the reverse: Given that A does not know ϕ in $S0$, and given that nothing occurs between $S0$ and $S1$ that would cause him to learn ϕ , we wish to infer that he still does not know ϕ in $S1$. Since “not knowing ϕ in S ” is represented as “there are possible worlds accessible from S in which ϕ is false,” this frame inference should have the following general form: If $S0A$ is accessible from $S0$, $S1 > S0$, $S1A > S0A$, and as far as A ’s sources of knowledge are concerned, the interval between $S0$ and $S1$ is indistinguishable from the interval between $S0A$ and $S1A$, then $S1A$ is accessible from $S1$.

Stating this formally is mostly a matter of collecting all the necessary sources of knowledge. Our theory requires that agent A gains knowledge in S under the following circumstances

1. If A begins action E in $S1$, and $S2$ is on a branch in which E is executed, then in $S2$, A knows that E is executed. If E is completed at or before $S2$, then in $S2$ A knows when it was completed.
2. If action E is feasible for A in situation S , then A knows that E is feasible in S .
3. If A receives a communication from AS in S then A knows in S that he has received a communication. If we assume axiom A.3.B, then A and AS have shared knowledge in S that A received a communication.

We also assume that there are domain-specific axioms of knowledge production. In an S5 logic, it is reasonable to assume that these are all of the following form: In all situations S , A knows whether $\Phi(A, S)$, where Φ is a formula that can refer only to present or past *physical* states or to past (but not present) knowledge states.⁷ Formally, we impose the following conditions on $\Phi(A, S)$:

- The only free variables in $\Phi(A, S)$ are A and S .
- If $S1$ is a quantified variable other than S appearing in Φ , and $S1$ is used as either the second-to-last or last argument for either k_acc or sk_acc , then the quantification of $S1$ imposes the restriction $S1 < S$.
- If $S1$ is a quantified variable other than S appearing in Φ , and $S1$ is not used as an argument for either k_acc or sk_acc , then the quantification of $S1$ imposes the restriction $S1 \leq S$.

Thus we assume the existence of a finite collection of axioms of the form

$$\forall_{A,S} [[\forall_{SA} k_acc(A, S, SA) \Rightarrow \Phi_i(A, S)] \vee [\forall_{SA} k_acc(A, S, SA) \Rightarrow \neg\Phi_i(A, S)]]$$

For example, (Scherl and Levesque, 2003) propose the use of an action “SENSE $_Q$ ” which informs the actor whether fluent Q is true. We can achieve that in the above framework by choosing $\Phi(A, S)$ to be the condition that A has executed SENSE $_Q$ and Q holds:

$$\Phi(A, S) \Leftrightarrow \exists_{S1} \text{occurs}(\text{SENSE}_Q, S1, S) \wedge \text{holds}(S, Q).$$

⁷Actually, I conjecture that these restrictions are not necessary, and that it is consistent to allow Φ to be any formula, but I have not proven it.

We now posit that every agent always knows whether $\Phi(A, S)$. Since, by axiom K.5, an agent always knows whether he has executed SENSE_Q , it follows that, if an agent has executed SENSE_Q , then he knows whether Q is true.

So now we can state the frame axiom asserting that if a knowledge accessibility relation disappears then one of the above conditions must have been met.

$$\begin{aligned}
\text{I.6: } & [\text{k_acc}(A, S0A, S0B) \wedge S0A < S1A \wedge S0B < S1B \wedge \\
& \text{time}(S1B) = \text{time}(S0B) \wedge \neg \text{k_acc}(A, S1A, S1B)] \Rightarrow \\
& [[\exists Z \neg [\exists_{SZA} \text{occurs}(\text{do}(A, Z), S1A, SZA) \wedge \text{ordered}(SZA, S2A)] \Leftrightarrow \\
& \quad [\exists_{SZB} \text{occurs}(\text{do}(A, Z), S1B, SZB) \wedge \text{ordered}(SZB, S2B)]]] \wedge \\
& \quad [\exists_{Z, TZ} \neg [\exists_{SZA} \text{occurs}(\text{do}(A, Z), S1A, SZA) \wedge SZA \leq S2A \wedge \text{time}(SZA) = TZ] \Leftrightarrow \\
& \quad \quad [\exists_{SZB} \text{occurs}(\text{do}(A, Z), S1B, SZB) \wedge SZB \leq S2B \wedge \text{time}(SZB) = TZ]]] \vee \\
& \quad [\exists_{AS, Q} \neg [\exists_{SQA} \text{occurs}(\text{do}(AS, \text{inform}(A, Q)), SQA, S2A) \Leftrightarrow \\
& \quad \quad [\exists_{SQB} \text{occurs}(\text{do}(AS, \text{inform}(A, Q)), SQB, S2B)]]] \vee \\
& \quad \exists_{S3A, S3B} S1A < S3A \leq S2A \wedge S1B < S3B \leq S2B \wedge \text{time}(S3A) = \text{time}(S3B) \wedge \\
& \quad \quad \bigvee_i \neg [\Phi_i(A, S3A) \Leftrightarrow \Phi_i(A, S3B)] \\
&]
\end{aligned}$$

That is: If $S0B$ is knowledge accessible from $S0A$ but later $S1B$ is not knowledge accessible from $S1A$ — that is, something was learned in between the two times to distinguish these — then one of the following conditions was met:

- A started some action Z on one branch but not the other. That is, it is not true that he started Z on one branch if and only if he started it on the other.
- A was informed of something on one branch but not the other.
- One of the domain specific conditions held at some time on one branch but not the other.

Well, there it is. It is not a candidate for any “Top 10 most elegant axioms” lists.

A more serious problem is that it doesn’t give us what we want. What we want is: Given that in s_0 , Sam doesn’t know whether Herbert Hoover invented the vacuum cleaner (P), and given that the only thing that happens between s_0 and s_1 is that Jack tells Sam that tea is selling for \$2 a pound in Shanghai (Q), we should be able to infer that Jack still doesn’t know whether Herbert Hoover invented the vacuum cleaner. But that inference is not valid. The problem is that it is consistent with the givens that Sam knows $\neg P \Leftrightarrow Q$, and so, when Jack tells him Q he finds out $\neg P$.

The problem here is not with the frame axiom; the frame axiom is fine. The problem is with the specification of the initial state. What you want to say is something like “All agents are as ignorant as possible, consistent with the givens,” but it is not easy to characterize what kind of possible worlds structure that would entail, let alone to formulate that characterization in a set of first-order axioms. Of course, in any particular case, one can get around this by adding more givens. You can specify that, in s_0 , Jack does not know that $Q \Rightarrow \neg P$; this approach is taken in the “Persistence of ignorance” theorem of (Scherl and Levesque, 2003). If the class of physical fluents is finite, you can assert that there is a possible world for each possible valuation consistent with the axioms, and that any two such possible worlds are knowledge accessible, unless the axioms rule this out. Since this is a finite structure, this, at least in principle, can be stated. But (a) if there are infinitely many possible states of the world, then it is not at all clear that this can be stated in first-order logic; and (b) it does not achieve *utter* ignorance, because it is common knowledge among all agents that none of them know anything about the physical fluents. In a system of this kind, Sam knows that Jack does not know P , whereas in the desired state of utter ignorance, presumably Sam doesn’t know whether or not Jack knows P . How to characterize a state of minimal knowledge of this kind is, as far as I know, an open problem.

4 Sample Inferences

We illustrate the power of the above theory with three toy problems.

4.1 Sample Inference 1:

Given:

- X.1 Sam knows in s_0 that it will be sunny on July 4.
 $[k_acc(sam, s_0, S_0A) \wedge S_0A < S_1A \wedge time(S_1A)=july4] \Rightarrow$
 $holds(S_1A, sunny).$
- X.2 In any situation, if it is sunny, then Bob can play tennis.
 $\forall_S holds(S, sunny) \Rightarrow feasible(occurs(do(bob, tennis), S))$
- X.3 Sam can always communicate with Bob.
 $\forall_{S_1} feasible(do(sam, communicate(bob)), S_1).$

Infer:

- X.P Sam knows that there is an action he can do (e.g. tell Bob that it will be sunny) that will cause Bob to know that he will be able to play tennis on July 4.

$$k_acc(sam, s_0, S_0A) \Rightarrow$$

$$\exists_{Z, S_1A} occurs(do(sam, Z), S_0A, S_1A) \wedge$$

$$\forall_{S_2A, S_2B, S_3B} [occurs(do(sam, Z), S_0A, S_2A) \wedge k_acc(bob, S_2A, S_2B) \wedge$$

$$S_2B < S_3B \wedge time(S_3B)=july4] \Rightarrow$$

$$feasible(do(bob, tennis), S_3B).$$

Proof:

By the comprehension axiom I.5 there is a fluent q_1 that holds in any situation S just if it will be sunny on July 4 following S .

$$P.1: holds(S, q_1) \Leftrightarrow [\forall_{S_1} [S < S_1 \wedge time(S_1)=july4] \Rightarrow holds(S_1, sunny)].$$

Let $z_1 = inform(bob, q_1)$. By axioms I.2, X.1, and X.3, $do(sam, z_1)$ is feasible in s_0 ;

$$P.2: feasible(do(sam, z_1), s_0).$$

By axiom K.5, Sam knows in s_0 that $do(sam, z_1)$ is feasible.

$$P.3: \forall_{S_0A} k_acc(s_0, S_0A) \Rightarrow feasible(do(sam, z_1), S_0A).$$

Let s_0a be any situation such that $k_acc(sam, s_0, s_0a)$.

By P.3, there exists a situation s_1a such $occurs(do(sam, z_1), s_0a, s_1a)$.

Let s_2a be any situationn such that $occurs(do(sam, z_1), s_0a, s_2a)$.

Let s_2b be any situation such that $k_acc(bob, s_2a, s_2b)$.

By Lemma 3.2, there exists s_1b such that $occurs(do(sam, z_1), s_1b, s_2b)$ and $holds(s_1b, q_1)$.

Let s_3b be any situation such that $s_2b < s_3b$ and $time(s_3b)=july4$.

By T.8 and T.4, $s_1b < s_3b$.

By P.2, $holds(s_3b, sunny)$.

By X.2, $feasible(do(bob, tennis), s_3b)$.

Applying universal abstraction over s_0a , s_2a , s_2b , and s_3b and existential abstraction over z_1 and s_1a gives us formula X.P.

4.2 Sample Inference 2

Given: Bob tells Alice that he has cheated on her. Alice responds by telling Bob that he has never told her anything she did not already know.

Y.1 Bob confesses to Alice that he has cheated on her.

$$\begin{aligned} & \exists_Q \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, Q)), s_0, s_1) \wedge \\ & \forall_S \text{holds}(S, Q) \Leftrightarrow \exists_{S_2, S_3} S_3 < S \wedge \text{occurs}(\text{do}(\text{bob}, \text{cheat}), S_2, S_3). \end{aligned}$$

Y.2 Alice responds that Bob has never told her anything she didn't already know. (Equivalently, whenever he has told her anything, she already knew it.)

$$\begin{aligned} & \exists_Q \text{occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, Q)), s_1, s_2) \wedge \\ & \forall_S \text{holds}(S, Q) \Leftrightarrow \\ & \quad \forall_{S_3, S_4, Q_1} [S_3 < S_4 \leq S \wedge \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, Q_1)), S_3, S_4)] \Rightarrow \\ & \quad \forall_{S_3A} \text{k_acc}(\text{alice}, S_3, S_3A) \Rightarrow \text{holds}(S_3A, Q_1). \end{aligned}$$

Infer: Bob now knows that Alice knew before he spoke that he had cheated on her.

Y.P Bob now knows that Alice had already known, before he spoke, that he had cheated on her.

$$\begin{aligned} & \forall_{S_2A} \text{k_acc}(\text{bob}, s_2, S_2A) \Rightarrow \\ & \quad \exists_{S_0A, S_1A, Q_1} S_1A < S_2A \wedge \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, Q_1)), S_0A, S_1A) \wedge \\ & \quad [\forall_{S_0B} \text{k_acc}(\text{alice}, S_0A, S_0B) \Rightarrow \\ & \quad \exists_{S_3B, S_4B} S_4B < S_0B \wedge \text{occurs}(\text{do}(\text{bob}, \text{cheat}), S_3B, S_4B)]. \end{aligned}$$

Proof:

Let q1 be the content of Bob's statement in Y.1, and let q2 be the content of Alice's statement in Y.2.

By K.4, Bob knows in s2 that he has informed Alice of q1.

$$\begin{aligned} \text{Q.1: } & \forall_{S_2A} \text{k_acc}(\text{bob}, s_2, S_2A) \Rightarrow \\ & \quad \exists_{S_0A, S_1A} S_1A < S_2A \wedge \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, q_1)), S_0A, S_1A). \end{aligned}$$

By Lemma 3.2, Bob knows in s2 that q2 held when Alice started to speak.

$$\begin{aligned} \text{Q.2: } & \text{k_acc}(\text{bob}, s_2, S_2A) \Rightarrow \\ & \quad \exists_{S_1A} \text{occurs}(\text{do}(\text{alice}, \text{inform}(\text{bob}, q_2)), S_1A, S_2A) \wedge \text{holds}(S_1A, q_2). \end{aligned}$$

Let s2a be any situation such that $\text{k_acc}(\text{bob}, s_2, s_2a)$, and let s1a be a corresponding value of S1A satisfying Q.2. Then $\text{holds}(s_1a, q_2)$.

By definition of q2, we have that in s1a, whenever Bob had previously told Alice anything (Q3), she had already known it.

$$\begin{aligned} \text{Q.3: } & \forall_{S_3, S_4, Q_3} [S_3 < S_4 \leq s_1a \wedge \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, Q_3)), S_3, S_4)] \Rightarrow \\ & \quad \forall_{S_3A} \text{k_acc}(\text{alice}, S_3, S_3A) \Rightarrow \text{holds}(S_3A, Q_1). \end{aligned}$$

By K.4 and Y.3, Bob knows in s1 that he has informed Alice of q1.

$$\begin{aligned} \text{Q.4: } & \forall_{S_1A} \text{k_acc}(\text{bob}, s_1, S_1A) \Rightarrow \\ & \quad \exists_{S_0A} \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, q_1)), S_0A, S_1A). \end{aligned}$$

In particular, therefore, Q.4 is true of $S_1A=s_1a$.

$$\text{Q.5: } \exists_{S_0A} \text{occurs}(\text{do}(\text{bob}, \text{inform}(\text{alice}, q_1)), S_0A, s_1a).$$

Let $s0a$ be a situation satisfying Q.5. Applying Q.3, with $S3 \rightarrow s0z$, $S4 \rightarrow s1a$, and $Q3 \rightarrow q1$, gives Q.6. $\forall_{S0B} k_acc(alice,s0a,S0B) \Rightarrow holds(S0B,q1)$.

Applying the definition of $q1$, we get the desired result.

4.3 Sample Inference 3

Given:

Z.1: Anne does not know that she has a brother.

$\neg[\forall_{S0A} k_acc(anne,s0,S0A) \Rightarrow \exists_Y holds(S0A,brother(Y,anne))]$.

Z.2: Anne knows that, if she had a brother, someone would have told her about him.

$\forall_{S0A} k_acc(anne,s0,S0A) \Rightarrow$
 $\forall_Y holds(S0A,brother(Y,anne)) \Rightarrow$
 $[\exists_{S1A,S2A,AS} S2A \leq S0A \wedge occurs(do(AS,inform(anne,brother(Y,anne)),S1A,S2A))]$

Z.3: Brotherhood is forever.

$S0 < S1 \wedge holds(S0,brother(X,Y)) \Rightarrow holds(S1,brother(X,Y))$

Infer: Anne knows that she has no brother.

Z.4: $\forall_{S0A} k_acc(anne,s0,S0A) \Rightarrow \neg \exists_Y holds(S0A,brother(Y,anne))$.

Note: This is a monotonic variant of the “auto-epistemic” inference (Moore, 1985b).

Proof by contradiction: Suppose that Z.4 is false and Anne does not know that she does not has a brother — in other words, as far as she knows she might have a brother.

R.1: $\exists_{S0A,Y} k_acc(anne,s0,S0A) \wedge holds(S0A,brother(Y,anne))$

Let sb and yb be values satisfying R.1. Thus $k_acc(anne,s0,sb)$ and $holds(sb,brother(yb,anne))$. By Z.2, in sb someone would have already told her that she had a brother.

R.2: $\exists_{S1A,S2A,AS} S2A \leq sb \wedge$
 $occurs(do(AS,inform(anne,brother(yb,anne)),S1A,S2A))$

Let $s1b$, $s2b$, as satisfy R.2. By Lemma 3.2, Anne would know in sb that she had previously had a brother. R.3: $\forall_{SC} k_acc(anne,sb,SC) \Rightarrow$

$\exists_{S1C} S1C < SC \wedge holds(S1C,brother(yb,anne))$.

Let $s0x$ be any situation such that $k_acc(anne,s0,s0x)$. By K.2 and K.3 $k_acc(anne,sb,s0x)$. By R.3 and X.3, $holds(s0x,brother(yb,anne))$. Applying universal abstraction to $s0x$ we have

R.4: $\forall_{S0X} k_acc(anne,S0,S0X) \Rightarrow holds(S0X,brother(yb,anne))$.

But this contradicts X.1.

5 Paradox

The following Russell-like paradox seems to threaten our theory:

Paradox: Let Q be a fluent. Suppose that over interval $[S0, S1]$, agent $a1$ carries out the action of informing $a2$ that Q holds. Necessarily, Q must hold in $S0$, since agents are not allowed to lie (axiom I.2). Let us say that this communication is *immediately obsolete* if Q no longer holds in $S1$. For example, if it is raining in $s0$, the event of $a1$ telling $a2$ that it is raining occurs over $[s0,s1]$, and it has stopped raining in $s1$, then this communication is immediately obsolete. Now let us say

that situation S is “misled” if it is the end of an immediately obsolete communication. Since “being misled” is a property of a situation, by the comprehension axiom it should be definable as a fluent. Symbolically,

$$\begin{aligned} \text{holds}(S, \text{misled}) &\equiv \\ \exists_{Q, A1, A2, S0} \text{occurs}(\text{do}(A1, \text{inform}(A2, Q)), S0, S) \wedge \neg \text{holds}(S, Q) \end{aligned}$$

Now, suppose that, as above, in s_0 it is raining; from s_0 to s_1 , a_1 tells a_2 that it is raining; and in s_1 it is no longer raining and a_1 knows that it is no longer raining. Then a_1 knows that “misled” holds in s_1 . Therefore, (axiom I.2) it is feasible for a_1 to tell a_2 that “misled” holds in s_1 . Suppose that, from s_1 to s_2 , a_1 informs a_2 that “misled” holds. The question is now, does “misled” hold in s_2 ? Well, if it does, then what was communicated over $[s_1, s_2]$ still holds in s_2 , so “misled” does not hold; but if it doesn’t, then what was communicated no longer holds, so “misled” does hold in s_2 .

The flaw in this argument is that it presupposes the independence axiom WRONG.1 that we rejected before. The argument presumes that if fluent $Q1 \neq Q2$, and $\text{do}(A1, \text{inform}(A2, Q1))$ occurs from s_1 to s_2 , then $\text{do}(A1, \text{inform}(A2, Q2))$ does not occur. (Our English description of the argument used the phrase “what was communicated between s_1 and s_2 ”, which presupposes that there was a unique content that was communicated.) But axiom I.4 asserts that many different fluents are communicated in the same act. Therefore, the argument collapses.

In particular, as we shall show, the clock time (in the sense of “the number of situations that have elapsed since the start of time”) is always common knowledge among all agents (Theorem 3, appendix A). Now, let q_1 be any fluent, and suppose that $\text{occurs}(\text{do}(a_1, \text{inform}(a_2, q_1)), s_1, s_2)$. Let $t_1 = \text{time}(q_1)$ and let q_2 be the fluent defined by the formula

$$\forall_S \text{holds}(S, q_2) \Leftrightarrow \text{holds}(S, q_1) \wedge \text{time}(S) = t_1.$$

By assumption, it is shared knowledge between a_1 and a_2 that $\text{holds}(s_1, q_2) \Leftrightarrow \text{holds}(s_1, q_1)$. Hence, by axiom I.4, $\text{occurs}(\text{do}(a_1, \text{inform}(a_2, q_2)), s_1, s_2)$. But by construction q_2 does not hold in s_1 ; hence the occurrence of $\text{do}(a_1, \text{inform}(a_2, q_2))$ from s_1 to s_2 is immediately obsolete. Therefore “misled” holds following *any* informative act.

Changing the definition of misled to use the universal quantifier, thus:

$$\begin{aligned} \text{holds}(S, \text{misled}) &\equiv \\ \forall_{Q, A1, A2} \text{occurs}(\text{do}(A1, \text{inform}(A2, Q)), S0, S) \wedge \neg \text{holds}(S, Q) \end{aligned}$$

does not rescue the contradiction. One need only change the definition of q_2 above to be

$$\forall_S \text{holds}(S, q_2) \Leftrightarrow \text{holds}(S, q_1) \vee \text{time}(S) \neq t_1.$$

Clearly, the new definition of “misled” *never* holds after any informative act.

Of course, if we extend the theory to include the underlying locutionary act, then this paradox may well return, as the locutionary act that occurs presumably is unique. However, as the content of a locutionary act is a quoted string, we can expect to have our hands full of paradoxes in that theory; this “misled” paradox will not be our biggest problem (Morgenstern, 1988).

6 Unexpected Hanging

The well-known paradox of the unexpected hanging (also known as the surprise examination) (Gardner, 1991; Quine, 1953) can be formally expressed in our theory; however, the paradox does not render the theory inconsistent. (The analysis below is certainly *not* a philosophically adequate solution to the paradox, merely an explanation of how our particular theory manages to side-step it.)

The paradox can be stated as follows:

A judge announces to a prisoner, “You will be hung at noon within 30 days; however, that morning you will not know that you will be hung that day.” The prisoner reasons to himself, “If they leave me alive until the 30th day, then I will know that morning that they will hang me that day. Therefore, they will have to kill me no later than the 29th day. So if I find myself alive on the morning of the 29th day, I can be sure that I will be hung that day. So they will have to kill me no later than the 28th day . . . So they can’t kill me at all!”

On the 17th day, they hung him at noon. He did not know that morning that he would be hung that day.

We can express the judge’s statement as follows:

$$\begin{aligned}
& \text{occurs}(\text{do}(\text{judge}, \text{inform}(\text{prisoner}, Q)), s_0, s_1) \wedge \\
& \forall_S \text{ holds}(S, Q) \Leftrightarrow \\
& \quad \forall_{SX} [S < SX \wedge \text{date}(SX) = \text{date}(S) + 31] \Rightarrow \\
& \quad \quad \exists_{SH, SM, SMA, SHA} \\
& \quad \quad \quad S < SM < SH < SX \wedge \text{hour}(SH) = \text{noon} \wedge \\
& \quad \quad \quad \text{holds}(SH, \text{hanging}) \wedge \text{hour}(SM) = 9\text{am} \wedge \\
& \quad \quad \quad \text{date}(SM) = \text{date}(SH) \wedge \\
& \quad \quad \quad \text{k_acc}(\text{prisoner}, SM, SMA) \wedge SMA < SHA \wedge \\
& \quad \quad \quad \text{hour}(SHA) = \text{noon} \wedge \text{date}(SM) = \text{date}(SH) \wedge \\
& \quad \quad \quad \neg \text{holds}(SHA, \text{hanging}).
\end{aligned}$$

That is: the content of the judge’s statement is the fluent defined by the following formula over S : On any timeline starting in S and going through some SX 31 days later, there is a situation SH at noon where you will be hung, but that morning SM you will not know you will be hung; that is, there is a SMA knowledge accessible from SM which is followed at noon by a situation SHA in which you are not hung.

Let UH^{lang} be the judge’s statement in English and let UH^{logic} be the fluent Q defined in the above formula. Let “kill(K)” be the proposition that the prisoner will be killed no later than the K th day, and let “kill_today” be the fluent that the prisoner will be killed today. It would appear that UH^{lang} is true; that the judge knows that in s_0 that it is true, and that UH^{logic} means the same as UH^{lang} . By axiom I.2, if the judge knows that UH^{logic} holds in s_0 , then he can inform the prisoner of it. How, then, does our theory avoid contradiction?

The first thing to note is that the prisoner *cannot* know UH^{logic} . There is simply no possible worlds structure in which the prisoner knows UH^{logic} . The proof is exactly isomorphic to the sequence of reasoning that prisoner goes through. Therefore, by Lemma 3.2 above, the judge cannot inform the prisoner of UH^{logic} ; if he did, the prisoner would know it to be true.

The critical point is that there is a subtle difference between UH^{lang} and UH^{logic} . The statement UH^{lang} asserts that the prisoner *will* not know kill_today — this means even *after* the judge finishes speaking. In our theory, however, one can only communicate properties of the situation at the beginning of the speech act and there is no way to refer to what *will* happens as distinguished from one *could* happen. So what UH^{logic} asserts is that the prisoner will not know kill_today *whatever* the judge decides to say or do in s_0 .

In fact, it is easily shown that either [the judge does not know in s_0 that UH^{logic} is true], or [UH^{logic} is false]. It depends on what the judge knows in s_0 . Let us suppose that in s_0 , it is inevitable that the prisoner will be killed on day 17 (the executioner has gotten irrevocable orders.) There are two main cases to consider.

- Case 1: All the judge knows $\text{kill}(K)$, for some $K > 17$. Then the most that the judge can tell the prisoner is $\text{kill}(K)$. In this case, UH^{logic} is in fact true in s_0 , but the judge does not know that it is true, because as far as the judge knows, it is possible that (a) he will tell the prisoner $\text{kill}(K)$ and (b) the prisoner will be left alive until the K th day, in which case the prisoner would know kill_today on the morning of the K th day.
- Case 2: The judge knows $\text{kill}(17)$. In that case, UH^{logic} is not even true in s_0 , because the judge has the option of telling the prisoner $\text{kill}(17)$, in which case the prisoner will know kill_today on the morning of the 17th day.

Again, we do not claim that this is an adequate solution to the philosophical problem, merely an explanation of how our formal theory manages to remain consistent and side-step the paradox. In fact, in the broader context the solution is not at all satisfying, for reasons that may well become serious when the theory is extended to be more powerful. There are two objections. First, the solution depends critically on the restriction that agents cannot talk about what *will* happen as opposed to what *can* happen; in talking about the future, they cannot take into account their own decisions or commitments about what they themselves are planning to do. One can extend the outer theory so as to be able to *represent* what will happen — in (Davis and Morgenstern, 2004), we essentially do this — but then the comprehension axiom I.5 must be restricted so as to exclude this from the scope of fluents that can be the content of an “inform” act. We do not see how this limitation can be overcome.

The second objection is that it depends on the possibility of the judge telling the prisoner $\text{kill}(17)$ if he knows this. Suppose that we eliminate this possibility? Consider the following scenario: The judge knows $\text{kill}(17)$, but he is unable to speak directly to the prisoner. Rather, he has the option of playing one of two tape recordings; one says “ $\text{kill}(30)$ ” and the other says UH^{logic} . Now the theory is indeed inconsistent. Since the prisoner cannot know UH^{logic} it follows that the judge cannot inform him of UH^{logic} ; therefore the only thing that the judge can say is “ $\text{kill}(30)$ ”. But in that case, the formula “ UH^{logic} ” is indeed true, and the judge knows it, so he should be able to push that button.

To axiomatize this situation we must change axiom I.2 to assert that that the only possible inform acts are $\text{kill}(30)$ and UH^{logic} .

Within the context of our theory, it seems to me that the correct answer is “So what?” Yes, you can set up a Rube Goldberg mechanism that creates this contradiction, but the problem is not with the theory, it is with the axiom that states that only these two inform acts are physically possible.

(Those readers, if any, who work through the proof of theorem 1 in appendix A may wonder what prevents this constraint from being incorporated into the construction of u-situations. After all, all that this amounts to is drastically restricting the class of “inform” acts that are added on. The answer is that which of the “inform” acts are allowed to exist now depends on the interpretation of a formula in the extended language, and that therefore the construction now involves a vicious cycle. See further the comments on Lemma 21.)

In a wider context, though, this answer will not serve. After all, it is physically possible to create this situation, and in a sufficiently rich theory of communication, it will be provable that you can create this situation. However, such a theory describing the physical reality of communication must include a theory of locutionary acts; i.e. sending signals of quoted strings. As mentioned above such a theory will run into *many* paradoxes; this one is probably not the most troublesome.

7 Consistency

Two paradoxes have come up, but the theory has side-stepped them both. How do we know that the next paradox won't uncover an actual inconsistency in the theory? We can eliminate all worry about paradoxes once and for all by proving that the theory is consistent. We do this by constructing a model satisfying the theory. More precisely, we construct a fairly broad class of models, establishing (informally) that the theory is not only consistent but does not necessitate any weird or highly restrictive consequences. (Just showing soundness with respect to a model or even completeness is not sufficient for this. For instance, if the theory were consistent only with a model in which every agent was always omniscient, and inform acts were therefore no-ops, then the theory would be *consistent* but not of any interest.)

As usual, establishing soundness has three steps: defining a model, defining an interpretation of the symbols in the model, and establishing that the axioms are true under the interpretation.

Our class of models is (apparently) more restrictive than the theory;⁸ that is, the theory is not complete with respect to this class of models. The major additional restrictions in our model are:

- I. Time must be discrete. We believe that this restriction can be lifted with minor modifications to the axioms, but this is beyond the scope of this paper. We hope to address it in future work.
- II. Time must have a starting point; it cannot extend infinitely far back. It would seem to be very difficult to modify our proof to remove this constraint; at the current time, it seems to depend on the existence of highly non-standard models of set theory.
- III. A knowledge accessibility link always connects two situations whose time is equal, where “time” measure the number of clock ticks since the start. In other words, all agents always have common knowledge of the time. In a discrete structure, this is a consequence of the axiom of memory. Therefore, it is not, strictly speaking, an additional restriction; rather, it is a non-obvious consequence of restriction (I). If we extend the construction to a non-discrete time line, some version of this restriction must be stated separately.

There are also more minor restrictions; for example, we will define shared knowledge to be the true transitive closure of knowledge, which is not expressible in a first-order language.

Theorem 1 below states that the axioms in this theory are consistent with essentially any physical theory that has a model over discrete time with a starting point state and physical actions.

Definition 1: A *physical language* is a first-order language containing the sorts “situations”, “agents”, “physical actionals”, “physical actions”, “physical fluents”, and “clock times”; containing the non-logical symbols, “<”, “do”, “occurs”, “holds”, “time”, and “communicate”; and excluding the symbols, “k_acc”, “inform”, and “sk_acc”.

Definition 2: (This is definition 6 of Appendix A). Let \mathcal{L} be a physical language, let \mathcal{T} be a theory over \mathcal{L} . \mathcal{T} is an *acceptable physical theory* (i.e. acceptable for use in theorem 1 below) if there exists a model \mathcal{M} and an interpretation \mathcal{I} of \mathcal{L} over \mathcal{M} such that the following conditions are satisfied:

1. \mathcal{I} maps the sort of clock times to the positive integers, and the relation $T1 < T2$ on clock times to the usual ordering on integers.

⁸The only way to be sure that the theory is more general than the class of models is to prove that it is consistent with a broader class of models.

2. Axioms T.1 — T.9 in table ?? are true in \mathcal{M} under \mathcal{I} .
3. Theory \mathcal{T} is true in \mathcal{M} under \mathcal{I} .
4. The theory is consistent with the following constraint: In any situation S , if any communication act is feasible, then arbitrarily many physically indistinguishable communication acts are feasible.
5. If α is a predicate symbol in \mathcal{L} with more than one situational argument, then $\alpha(X_1 \dots X_k)$ holds only if all the situations among $X_1 \dots X_k$ are ordered with respect to $<$. (Note that this condition holds both when α is “ $<$ ” and α is “occurs”.) If $\beta(X_1 \dots X_k)$ is a function symbol, then the above condition holds for the relation $X_{k+1} = \beta(X_1 \dots X_k)$.

Condition (4) no doubt seems complex, strange, and restrictive. But in fact any physical model can be easily transformed into one satisfying this condition: take the original model and, wherever a communicative act occurs, make an infinite number of identical copies of the subtree following the branch where the act occurs. Moreover, most reasonable physical theories \mathcal{T} will accept this transformation, or can be straightforwardly transformed into theories that will accept this transformation. In fact, therefore, condition (4) is not a substantial restriction on \mathcal{T} .

For several reasons, it is unfortunate that condition (5) needs to be included:

- It was not included in the KR-2004 paper.
- I don’t know that it’s necessary; in fact, I suspect that the theorem is true even if this condition is dropped (certainly not true of the other conditions.)
- This condition is satisfied in most causal theories; generally a causal theory refers only to situations in a single time line, which is what is required here. However, it is hard to be sure that you will *never* encounter a causal theory where it would be natural to use a relation that violates this condition.

However, I have not found a proof for languages that violate this condition.

Of course, if \mathcal{L} contains a symbol α that violates condition 5, but that is defined in \mathcal{T} using a rule $\alpha \leftrightarrow \phi$ where ϕ contains only symbols that respect condition 5, then that is not a problem; we can simply replace α by ϕ throughout \mathcal{T} , and thus obtain a theory that respects condition 5. The problematic case is where there are symbols whose interpretation in \mathcal{I} violates condition 5, and that are not reducible to symbols that respect condition 5.

(The KR-2004 paper claims that condition (4) can be stated in a first-order axiom schema. This is in error. More precisely, I have not found any first-order axiom schema that can be used to instantiate condition 4 that I can prove to be sufficient for the theorem below.)

Theorem 1: Let \mathcal{T} be an acceptable physical theory, and let \mathcal{U} be \mathcal{T} together with axioms K.1 — K.8 and I.1 — I.5. Then \mathcal{U} is consistent.

It is possible to strengthen theorem 1 by adding in domain-specific axioms of knowledge acquisition and the associated frame axiom over accessibility relation, as described in section ??, plus conditions on the initial knowledge and ignorance of the agents. Specifically, we have the following theorem:

Theorem 2: Let \mathcal{T} be an acceptable physical theory, and let \mathcal{U} be the union of:

- A. \mathcal{T} ;

- B. Axioms K.1 — K.7 and I.1 — I.5.
- C. A collection of domain-specific knowledge acquisition axioms of the form specified in section ??.
- D. The frame axiom I.6 associated with the axioms in (C).
- E. Any set of axioms \mathcal{K} specifying the presence or absence of `k_acc` relations among situations at time 0 as long as:
 - i. The axioms in \mathcal{K} do not refer to any situations of time later than 0.
 - ii. The axioms in \mathcal{K} are consistent with \mathcal{T} , axioms K.1 — K.3, K.5 (as regards knowing the feasibility of actions at time 0); and the axioms in (C).

Then \mathcal{U} is consistent.

In appendix A, we sketch how the proof of theorem 1 is modified to give a proof of theorem 2.

8 Related Work

The theory presented here was originally developed as part of a larger theory of multi-agent planning (Davis and Morgenstern, 2004). That theory includes requests as speech acts as well as informative speech acts. However, our analysis of informative acts there was not as deep or as extensive in scope.

As far as we know, this is the first attempt to characterize the content of communication as a first-order property of possible worlds. Morgenstern (1988) develops a theory in which the content of communication is a string of characters. A number of BDI models incorporate various types of communication. The general BDI model was first proposed by Cohen and Perrault (1979); within that model, they formalized illocutionary acts such as “Request” and “Inform” and perlocutionary acts such as “Convince” using a STRIPS-like representation of preconditions and effects on the mental states of the speaker and hearer. Cohen and Levesque (1990) extend and generalize this work using an full modal logic of time and propositional attitudes. Here, speech acts are *defined* in terms of their effects; a request, for example, is any sequence of actions that achieves the specified effect in the mental state of the hearer.

Update logic (e.g. Plaza 1989; van Benthem 2003) combines dynamic logic with epistemic logic, introducing the dynamic operator $[A!]\phi$, meaning “ ϕ holds after A has been truthfully announced.”. The properties of this logic have been extensively studied. Baltag, Moss, and Solecki (2002) extend this logic to allow communication to a subset of agents, and to allow “suspicious” agents. Colombetti (1999) proposes a *timeless* modal language of communication, to deal with the interaction of intention and knowledge in communication. Parikh and Ramanujam (2003) present a theory of *messages* in which the meaning of a message is interpreted relative to a protocol.

There is a large literature on the applications of modal logics of knowledge to a multi-agent systems. For example, Sadek et al. (1997) present a first-order theory with two modal operators $B_i(\phi)$ and $I_i(\phi)$ meaning “Agent i believes that ϕ ” and “Agent i intends that ϕ ” respectively. An inference engine has been developed for this theory, and there is an application to automated telephone dialogue that uses the inference engine to choose appropriate responses to requests for information. However, the temporal language associated with this theory is both limited and awkward; it seems unlikely that the theory could be applied to problems involving multi-step planning. (The dialogue application requires only an immediate response to a query.)

The multi-agent communication languages KQML (Finin et al., 1993) and FIPA (FIPA, 2001) provide rich sets of communication “performatives”. KQML was never tightly defined (Woolridge 2002.) FIPA has a formal semantics defined in terms of the theory of (Sadek et al. 1997) discussed

above. However, the content of messages is unconstrained; thus, the semantics of the representation is not inherently connected with the semantics of the content, as in our theory.

Other modal theories of communication, mostly propositional rather than first-order, are discussed in (Wooldridge and Lomuscio, 2000; Lomuscio and Ryan, 2000; Rao, 1995).

The theory of runs and messages, presented in (Fagin et al. 1995) developed a constructive model of a system of agents. Each agent is characterized as a infinite sequence. A state of agent A at time T is the prefix of the first T element of the corresponding sequence. The global state of the system at time T is the tuple of the states of all the agents at time T . Two global system states $Q1$ and $Q2$ are knowledge accessible relative to A if the state of A is the same in $Q1$ and $Q2$. A message is an event that modifies the state of the sender when it is sent and the state of the recipient when received. There is a protocol that governs under what circumstances a sender may send a specified message. The meaning of the message can be identified with the knowledge gained by the recipient when it is received. If one identifies “possible world” with “global state of the system”, this gives a clear and simple semantics for knowledge and informative acts. Moreover, it has the remarkable advantage that, given a suitable interpretation of the symbols, axioms T.1–T.9, K.1–K.7, I.1, I.2, I.3, and I.6 are all consequences of the definition. However, it does not seem to be a quite adequate framework for our theory, since there doesn’t seem to be a way to achieve axioms I.4 and I.5. The reason that the (similar) semantics that we give in Appendix A does not quite fit within this framework is that this framework assumes a *fixed* set of messages, whereas in our semantics, a message sent at time K corresponds to a set of situations at time K ; and this kind of mutual recursion between system states and messages is not allowed within Fagin et al.’s definition of a system.

9 Conclusions and Open Problem

We have developed a theory of communications which allows the content of an informative act to include quantifiers and logical operators and to refer to physical states, events including other informative acts, and states of knowledge; all these in the past, present, or possible futures. We have proven that this theory is consistent, and compatible with a wide range of physical theories. We have examined how the theory avoids two potential paradoxes, and discussed how these paradoxes may pose a danger when these theories are extended. Elsewhere (Davis and Morgenstern, 2004) we have shown that the theory can be integrated with a similarly expressive theory of multi-agent planning.

The major technical problem that follows naturally on this work is to find ways to relax the limitations enumerated in section 1 while preserving the consistency of the theory. Let us discuss what is involved here a little.

The most irksome of the restrictions is that the sender AS knows when the communication is complete and that, when the communication is complete, the recipient AH knows when the communication was initiated. This rules out application to most mail-like communications, or any communication media with an unknown delay. The problem is to find a modified version of section I.4 which is suitable for this more relaxed setting. Untimed communication, especially where the recipient does not know the time when the communication was initiated, leads to complex and confusing possible worlds structures, and I have not yet managed to think my way through them. However, I would be surprised if there were any insurmountable problems here.

The restriction that the sender and recipient know each other is one that, in practice, is often enough violated, and it would certainly be interesting to relax this. If you relax this condition, then a timed communication (i.e. one satisfying I.4) gives rise to *anonymous shared knowledge*. That is,

the speaker and hearer know that they share the knowledge of the content; they just don't know who they are sharing the knowledge with. This is analogous to common knowledge among non-rigid sets (Fagin et al. 1995, section 6.4) but the different setting here raises different issues.

The restriction to discrete time obviously impedes the integration of this theory with physical theories that use continuous time. The problem is that the construction of the model in our consistency proof is inherently iterative over time, and there does not seem to be any easy way to modify this iterative structure. The proof will work if one makes strong assumptions about the discreteness of communicative acts; e.g. one posits that it is only physically possible to begin a communication in a situation whose clock time is a non-negative integer. It is conceivable that such a theory would suffice for most applications; one would have to look over examples of reasoning that integrate continuous physical reasoning with communication, which I have not yet done. I would conjecture that axioms K.1 — K.7 and I.1 — I.6 are in fact consistent with a continuous model of time, without modification, and without the need to impose strong conditions on the physics of communication, but I am certainly far from a proof.

Other, more far-reaching, problems include:

- The problem of characterizing a “maximally ignorant” initial state, discussed in section ??.
- Having defined the notion of a “generalized fluent”, an obvious next step is to define “ $\text{know}(A, Q)$ ” as a first-order function mapping agent A and fluent Q into the fluent of A knowing Q . The axiomatics of this representation would be interesting to study.
- Our work on integrating this work with a theory of planning (Davis and Morgenstern, 2004) involves some rather restrictive constraints on the protocols between agents. We would like to study how the theory can be modified to weaken these.
- To my mind, the brass ring in this field would be to integrate the above theory of illocutionary acts, which describes the content of communications, with a theory of locutionary acts, which would describe the form of communications. Achieving a theory that is both general and consistent would be a major accomplishment.

10 References

- Baltag, A., Moss, L. and Solecki, S. 2002. “The Logic of Public Announcements: Common Knowledge and Private Suspicions.”
- Benthem, J. van. 2003. “‘One is a Lonely Number’: on the logic of communication.” ILLC Tech Report 2003-07, Institute for Logic, Language and Computation, University of Amsterdam.
- Cohen, P.R. and Perrault, C.R. 1979. “Elements of a plan-based theory of speech acts.” *Cognitive Science*, vol. 3, no. 3, pp. 177-212.
- Cohen, P.R. and Levesque, H. 1990. “Intention is choice with commitment” *Artificial Intelligence*, vol. 42, nos. 2-3, pp. 213-261.
- Colombetti, M. 1999. “A Modal Logic of Intentional Communication.” *Mathematical Social Sciences*, vol. 38, pp. 171-196.
- Davis, E. 1988. “Inferring Ignorance from the Locality of Visual Perception.” *Proc. AAAI-88*, pp. 786-790
- Davis, E. “A First-Order Theory of Communicating First-Order Formulas,” *KR-04*.

- Davis, E. and Morgenstern, L. 2004. "A First-Order Theory of Communication and Multi-Agent Plans." Submitted to *Journal of Logic and Computation*.
- Fagin, R., J. Halpern, Y. Moses, and M. Vardi. 1995. *Reasoning about Knowledge*. MIT Press.
- Finin et al. 1993. "Specification of the KQML agent communication language." DARPA knowledge sharing initiative external interfaces working group.
- FIPA 2001. "The foundation for intelligent physical agents." <http://www.fipa.org/>
- Gardner, M. 1991. *The Unexpected Hanging and Other Mathematical Diversions*. Chicago University Press.
- Lomuscio, A. and Ryan, M. 2000. "A spectrum of modes of knowledge sharing between agents." *Intelligent Agents VI: Agent Theories, Architectures, and Languages*. Lecture Notes in Artificial Intelligence 1757, Springer-Verlag, pp. 13-26.
- McDermott, D. 1982. "A Temporal Logic for Reasoning about Processes and Plans." *Cognitive Science*, Vol. 6, pp. 101-155.
- Moore, R. 1980. "Reasoning about Knowledge and Action." Tech. Note 191, SRI International, Menlo Park, CA.
- Moore, R. 1985a. "A Formal Theory of Knowledge and Action." In Jerry Hobbs and Robert Moore, (eds.) *Formal Theories of the Commonsense World*. ABLEX Publishing, Norwood, New Jersey, pp. 319-358.
- Moore, R. 1985b. "Semantical Considerations on Nonmonotonic Logic." *Artificial Intelligence*. vol. 25 p. 75-94.
- Morgenstern, L. 1987. "Foundations of a Logic of Knowledge, Action, and Communication." NYU Ph.D. Thesis.
- Parikh, R. and Ramanujam, R., 2003. "A Knowledge-Based Semantics of Messages." *Journal of Logic, Language, and Information*. vol. 12 no. 4.
- Plaza, J. 1989. "Logics of Public Announcements." *Proc. 4th International Symposium on Methodologies for Intelligence Systems*.
- Quine, W.V.O. 1953. "On a So-Called Paradox." *Mind* vol. 62, pp. 65-67.
- Rao, A.S. 1995. "Decision Procedures for Propositional Linear Time Belief-Desire-Intention Logics." *Intelligent Agents II: Agent Theories, Architectures, and Languages*. Lecture Notes in Artificial Intelligence 1037, Springer-Verlag, pp. 33-48.
- Reiter, R. 2001. *Knowledge in Action: Logical Foundations for Specifying and Implementing Dynamical Systems*. MIT Press.
- Sadek, M.D., Bretier, P. and Panaget, F. 1997. "ARTIMIS: Natural dialogue meets rational agency." *Proc. IJCAI-97*, pp. 1030-1035.
- Scherl, R. and Levesque, H. 1993. "The Frame Problem and Knowledge Producing Actions." *Proc. AAAI-93*, pp. 689-695.
- Scherl, R. and Levesque, H. 2003. "Knowledge, action, and the frame problem." *Artificial Intelligence*, vol. 144 no. 1, pp. 1-39.
- Woolridge, M. 2002. *An Introduction to MultiAgent Systems*. John Wiley and Sons.
- Wooldrige, M. and Lomuscio, A. 2000. "Reasoning about Visibility, Perception, and Knowledge." *Intelligent Agents VI: Agent Theories, Architectures, and Languages*. Lecture Notes in Artificial

Appendix A: Proof of Theorem 1

This appendix contains a proof of theorem 1. Specifically, we prove that if \mathcal{T} is a physical theory over integer-valued time satisfying a few, not very restrictive, constraints, then \mathcal{T} is consistent with our axioms of knowledge and of communication.

Outline of paper: In section A.1 we give a formal definition of what we mean by a physical theory. In section A.2, we show how a model of a physical theory can be extended to incorporate knowledge relations and informative actions. In section A.3, we define the interpretation of our theory over the new model. In section A.4, we prove that this interpretation over this model satisfies both the original physical theory and the axioms of knowledge and communication.

A.1: A Physical Theory

A physical theory is a set of constraints on actions and fluents. A communicative action may have physical preconditions, effects, or other constraints, but these may not depend on the *content* of the communication. That is, from the physical point of view, communicative actions are distinguished only by the identity of the speaker and hearer, not the content. Physical theories do not refer to knowledge states.

Our objective here is to prove that any reasonable physical theory is consistent with our theory of knowledge and communication. To do this, we have to ensure that the two theories “join up”, so to speak; specifically, that the physical theory does not impose any constraints that are incompatible with the epistemic theory. There are three potential sources of trouble.

- Axiom I.1, I.2, and I.4 together imply that, if AS can communicate with AH then, in general, there are a large number of different possible communicative acts that AS can perform. Specifically, in any situation S , if $Q1$ and $Q2$ are fluents such that (a) AS knows that both $Q1$ and $Q2$ hold; but (b) it is not shared knowledge between AS and AH that $Q1 \Leftrightarrow Q2$, then the act of AS informing AH that $Q1$ different from the act of AS informing AH that $Q2$. The physical theory could make this impossible by asserting that only a small number of different communicative acts are feasible in S . For instance, the statement that only two different communicative acts are feasible in $s0$ could be stated in the formula

$$\begin{aligned} & \exists_{S1a, S1b} \text{occurs}(\text{do}(\text{as}, \text{communicate}(\text{ah})), s0, S1a) \wedge \\ & \text{occurs}(\text{do}(\text{as}, \text{communicate}(\text{ah})), s0, S1b) \wedge S1a \neq S1b \wedge \\ & \forall_{S1} \text{occurs}(\text{do}(\text{as}, \text{communicate}(\text{ah})), s0, S1) \Rightarrow [S1 = S1a \vee S1 = S1b] \end{aligned}$$

To block this, we impose condition (4) in definition 6 below: A physical theory must be consistent with the constraint that, if any communicative action is feasible in a situation, then infinitely many physically indistinguishable actions are feasible in that situation.

- Axiom I.5 asserts the existence of a large number of fluents. The physical theory could assert that only a limited class of fluents exist. E.g. the following axiom asserts that the only fluents have the form “on(A, B)” where A and B are blocks.

$$\forall_Q \exists_{A, B} \text{block}(A) \wedge \text{block}(B) \wedge Q = \text{on}(A, B).$$

This is not at all far-fetched; one approach to the frame problem is to assert “The only fluents changed by action A are $Q1 \dots Qk$,” which leads to the same kind of problem.

We get around this problem by distinguishing between *physical fluents* and *general fluents*, and requiring that a physical theory can only refer to physical fluents.

- Similarly, the theory of communication requires the existence of actionals “inform(AH, Q)” and of actions “do($AS, \text{inform}(AH, Q)$).” We have to make sure that the physical theory does not simply prohibit these; e.g. assert that the only possible actionals have the form “communicate(AH)” and “puton(A, B)”. To insure this, we require that the physical theory can only refer to physical actions and actionals.

Definition 1: A *physical language* is a first-order language containing the sorts “situations”, “agents”, “physical actionals”, “physical actions”, “physical fluents”, and “clock times”; containing the non-logical symbols, “<”, “do”, “occurs”, “holds”, “time”, and “communicate”; and excluding the symbols, “k_acc”, “inform”, and “sk_acc”. (The language may or may not contain any sort or non-logical symbol other than those mentioned above.)

Definition 2: Let \mathcal{L} be a physical language. Let \mathcal{M} be a model and let \mathcal{I} be an interpretation of \mathcal{L} in \mathcal{M} . Let s_0 and s_1 be situations in \mathcal{M} . Situation s_1 is a *successor* of s_0 if $s_0 < s_1$ and there is no situation s_m such that $s_0 < s_m < s_1$

Here, and in subsequent definitions, we implicitly use \mathcal{I} to apply nomenclature from \mathcal{L} to entities in \mathcal{M} . More formal statements of the condition “ $s_0 < s_1$ ” above would be, “The pair $\langle s_0, s_1 \rangle \in \mathcal{I}(\text{'<'})$ ” or “The open formula $SA < SB$ is true in \mathcal{M} under \mathcal{I} under the valuation $SA \rightarrow s_0, SB \rightarrow s_1$.”. We will use the shorter form when it is clear; when necessary, we will be more precise.

Definition 3: Let $\mathcal{L}, \mathcal{M}, \mathcal{I}$ be as above. Let s_0, s_1 be situations in \mathcal{M} . We say that s_1 is a *communication successor* of s_0 if s_1 is a successor of s_0 and there exist agents as, ah and a situation sz such that $s_1 \leq sz$ and occurs(do($as, \text{communicate}(ah)$), s_0, sz).

Definition 4: Let $\mathcal{L}, \mathcal{M}, \mathcal{I}$ be as above. Let τ be a function from \mathcal{M} to itself which is one-to-one and onto. The function τ is said to be a *situational automorphism* if the following conditions hold:

1. If X is not a situation, then $\tau(X) = X$.
2. Let α be a predicate symbol in \mathcal{L} with k arguments or a function symbol with $k - 1$ arguments. Note that, under standard Tarskian semantics, $\mathcal{I}(\alpha)$ is a set of k -tuples of elements of \mathcal{M} . A tuple $\langle x_1 \dots x_k \rangle \in \mathcal{I}(\alpha)$ if and only if $\langle \tau(x_1) \dots \tau(x_k) \rangle \in \mathcal{I}(\alpha)$.

Definition 5: Two situations SA and SB are *indistinguishable* if the following holds: Let SSA be the part of the time structure following SA and SSB be the part of the time structure following SB .

$$\begin{aligned} SSA &= \{ S \in \mathcal{M} \mid SA \leq S \} \\ SSB &= \{ S \in \mathcal{M} \mid SB \leq S \} \end{aligned}$$

Then there exists a situational automorphism τ over \mathcal{M} such that $\tau(SSA) = SSB$, $\tau(SSB) = SSA$, and for any situation S which is not in SSA and SSB , $\tau(S) = S$.

Definition 6: Let \mathcal{L} be a physical language, and let \mathcal{T} be a theory over \mathcal{L} . \mathcal{T} is an *acceptable physical theory* (i.e. acceptable for our discussion here) if there exists a model \mathcal{M} and an interpretation \mathcal{I} of \mathcal{L} over \mathcal{M} such that the following conditions are satisfied:

1. \mathcal{I} maps the sort of clock times to the positive integers, and the relation $T1 < T2$ on clock times to the usual ordering on integers.
2. \mathcal{M} satisfies Axioms T.1 — T.9 in table ?? under \mathcal{T} , where T.8 and T.9 are restricted to physical actions.
3. \mathcal{M} satisfies theory \mathcal{T} under \mathcal{I} .
4. For any situations s_0, s_1 and agents a_s, a_h in \mathcal{M} , if s_1 is a communication successor of s_0 , then there are infinitely many successors of s_0 that are physically indistinguishable from s_1 .
5. If α is a predicate symbol in \mathcal{L} with more than one situational argument, then $\alpha(X_1 \dots X_k)$ holds only if all the situations among $X_1 \dots X_k$ are ordered with respect to $<$. (Note that this condition holds both when α is “ $<$ ” and α is “occurs”.) If $\beta(X_1 \dots X_k)$ is a function symbol, then the above condition holds for the relation $X_{k+1} = \beta(X_1 \dots X_k)$.

We can now state precisely the theorem that is the objective of this appendix.

Theorem 1:

Let \mathcal{T} be an acceptable physical theory, and let \mathcal{A} be \mathcal{T} together with axioms K.1 — K.7 and I.1 — I.5, and with T.8 and T.9 extended to arbitrary actions. Then \mathcal{A} is consistent.

Sections A.2-A.4 give the proof of this theorem.

A.2: Model construction

Sketch of model construction

The main sticking point of the proof is as follows: In order to satisfy the comprehension axiom, we must define a fluent to be any set of situations. However, if Q is a fluent, then the act of AS informing AH of Q in $S1$ generates a new situation; and if we generate a separate “inform” act for each fluent, then we would have a unsolvable vicious circularity.

We are rescued here by axiom I.4 together with the theorem, proven in theorem 2 below, that, in a discrete time structure satisfying the axiom of memory (K.4), knowledge accessibility relations can only connect situations of the same time, and therefore the current time is always common knowledge between all agents. Let q_1 be any fluent that holds in situation s_1 . By axiom I.4, if AS informs AH of q_1 over the interval $[s_1, s_2]$ and AS and AH have shared knowledge that $q_1 \Leftrightarrow q_2$ in s_1 , then the same act can be characterized as AS informing AH of q_2 . Let $t_1 = \text{time}(s_1)$. Let q_2 be the fluent such that $\text{holds}(S, q_2) \Leftrightarrow \text{holds}(S, q_1) \wedge \text{time}(S) = t_1$. Then AS and AH have shared knowledge in s_1 that q_1 is equivalent to q_2 . Therefore, it suffices to generate an occurrence of an inform act starting in $S1$ only for fluents like q_2 that specify the current time, and such a fluent can be identified with a set of situations of the same time as $S1$. This limitation allow us to break the circularity in the construction of situations and informative acts: the content of informative acts starting at time K is a subset of the situations whose time is K ; informative acts starting in time K generate situations whose time is $K + 1$.

Therefore, we can use the “algorithm” shown in table ?? to construct a model of the theory \mathcal{A} . The main difference between the model \mathcal{M} of theory \mathcal{T} and the model \mathcal{U} of \mathcal{A} is that \mathcal{U} contains many more situations. To avoid confusion, we will call the situations of \mathcal{M} “p-situations” and call the situations of \mathcal{U} “u-situations”. Each u-situation US has a corresponding p-situation, denoted $\text{PHYS}(US)$, which is physically indistinguishable from US . The difference is that US may associate specific contents with some of the communication actions that precede $\text{PHYS}(US)$.

Constructing a model

```

procedure model_construct(in  $\mathcal{T}$  : an acceptable physical theory;
                         $\mathcal{M}$  : a model of theory  $\mathcal{T}$ )
    return a structure of u-situations over which we will define
            a model of the extended theory.
for each p-situation  $PS$  in  $\mathcal{M}$ , construct a u-situation  $US$ .
    Label  $\text{PHYS}(US) = PS$ ,  $\text{time}(US)=0$ .
for (each agent  $A$ ), define the relation  $\text{K\_ACC}(A, \cdot, \cdot)$ 
    to be some equivalence relation over the u-situations constructed above.
for ( $K=0$  to  $\infty$ ) do {
    for (each u-situation  $S$  of time  $K$ ) do {
        for (each p-situation  $PS$  following  $\text{PHYS}(S)$  in  $\mathcal{M}$ )
            construct a new u-situation  $S1$  and mark  $\text{PHYS}(S1)=PS$ ;
        for (each pair of agents  $AS, AH$ ) do {
            if (in  $\mathcal{M}$  there is an act starting in  $S$  of  $AS$  communicating to  $AH$ )
                then {
                     $\text{SSL} :=$  the set of u-situations knowledge-accessible from  $S$ 
                        relative to the knowledge of  $AS$ ;
                     $\text{SSU} :=$  the set of u-situations knowledge-accessible from  $S$ 
                        relative to the shared knowledge of  $AS$  and  $AH$ ;
                    for (each set  $SS$  that is a subset of  $\text{SSU}$  and a superset of  $\text{SSL}$ ) do {
                        construct an action “do( $AS, \text{inform}(AH, SS)$ )” starting in  $S$ ;
                        construct a successor  $S1$  of  $S$  corresponding to the execution of this action;
                        label  $\text{PHYS}(S1)$  to be a u-situation in  $\mathcal{M}$  following a communicate action in  $\text{PHYS}(S)$ ;
                    }
                }
        }
    }
    use the axioms of knowledge to construct a valid set of
        knowledge accessibility relations over the new u-situations
} return (the set of u-situations plus the set of knowledge accessibility relations)

```

Table 4: Construction of a model

Theorem 3: If the set of clocktimes is equal to the positive integers, then for any situations SA, SB , if $\text{k_acc}(A, SA, SB)$ then $\text{time}(SA)=\text{time}(SB)$.

Proof: Suppose that $\text{time}(SA) < \text{time}(SB)=k$. By axioms T.7, T.6 and T.5, there exist situations $SB_0 < SB_1 < \dots < SB_{k-1} < SB$ such that $\text{time}(SB_i)=i$. By axiom K.4 there exist $SA_0 \dots SA_{k-1}$ such that $\text{k_acc}(A, SA_i, SB_i)$, $SA_{i-1} < SA_i$ and $SA_{k-1} < SA$; but this is impossible, since $\text{time}(SA) < k$.

Formal construction of the model

The definitions in this section essentially amount to a formalized re-statement of the “algorithm” in table ??.

Let \mathcal{L} be a physical language. Let \mathcal{T} be an acceptable physical theory over \mathcal{L} . Let \mathcal{M} be a model and let \mathcal{I} be an interpretation of \mathcal{L} satisfying the conditions of definition 6.

The remaining definitions in this section are relative to a fixed choice of \mathcal{L} , \mathcal{T} , \mathcal{M} , and \mathcal{I} .

For convenience, for each symbol τ in \mathcal{T} , including sorts, we use the same symbol in block

capitals to denote the image of τ under \mathcal{I} ; this is an individual, a subset, a mapping, or a relation over \mathcal{M} . Thus, for example, AGENTS is the image under \mathcal{I} of the sort “agents”; TIME is the image under \mathcal{I} of the function symbol “time” and so on.

We now proceed to building up the set of u-situations. This construction is recursive over time. Naturally, the base case is at time 0.

The most important and complex part of the construction is the wider class of situations that we will need. In general a u-situation US is a pair $\langle S1, MM \rangle$ where:

- S1 is a p-situation. We will write $S1 = \text{PHYS}(US)$.
- MM is a set of 4-tuples $\langle AS, AH, USSQ, SX \rangle$. AS and AH are agents; USSQ is a set of u-situations; and SX is a p-situation such that $SX < \text{PHYS}(US)$ and such that $\text{OCCURS}(\text{DO}(AS, \text{COMMUNICATE}(AH)), SX, SZ)$. Such a tuple asserts that an action of AS informing AH of content USSQ began in a u-situation $USX < US$. We write $MM = MM(US)$.

It will be convenient to posit the existence of an atomic entity **INFORM**, which is not in \mathcal{M} , and of an entity **DO**.

Definition 7: Let PS be a p-situation such that $\text{TIME}(PS) = 0$. A *u-situation at time 0* is a pair of the form $US = \langle PS, \emptyset \rangle$. The function $\text{ANCESTOR}(US)$ maps a u-situation US to a set of u-situations, the ancestors of US in the time structure.

Definition 8: A *time structure of depth 0* TS is a pair:

- The set of u-situations $U_SITS = \{ \langle PS, \emptyset \rangle \mid PS \in \text{SITUATIONS}, \text{TIME}(PS) = 0 \}$ with one u-situation for each p-situation at time 0.
- A function K_ACC mapping any agent $A \in \text{AGENTS}$ to an equivalence relation over U_SITS.

Definitions 9 through 15 are mutually recursive, successively building up the model forward in time.

Definition 9: Let TS be a time structure of depth K. Let US be a u-situation of time K in TS. Let $S1 = \text{PHYS}(US)$. Let MM be a collection of 4-tuples as described above. Let S2 be a successor to S1. The *simple successor to US parallel to S2* is the pair $\langle S2, MM \rangle$.

Definition 10: Let $TS = \langle U_SITS, K_ACC \rangle$, US, S1, MM be as above. Let AS and AH be agents. A *possible communicative content* from AS to AH is a set of u-situations USSQ of time K in U_SITS satisfying the following: Let USSL be the set of u-situations USA in TS such that $\langle US1, USA \rangle \in K_ACC(AS)$. Let USSU be the set of u-situations USA in USSL such that there exist $US_0 = US, US_1, US_2 \dots US_N = USA$, such that for each J, $\langle US_J, US_{J+1} \rangle$ is either in $K_ACC(AS)$ or in $K_ACC(AH)$. Then $USSL \subseteq USSQ \subseteq USSU$.

The 4-tuple $\langle AS, AH, USSQ, S1 \rangle$ is called an *inform indicator* starting in S1.

Definition 11: Let TS, US, S1, MM be as above. Let S2 be a successor of S1. Let $I = \langle AS, AH, USSQ, S1 \rangle$ be an inform indicator starting in S1. I *possibly leads toward* S2 if there exists $SZ \geq S2$ such that $\text{OCCURS}(\text{DO}(AS, \text{COMMUNICATE}(AH)), S1, SZ)$. An *informative sheaf* in US toward S2 is a set

MMX of inform indicators in US toward S2 such that no two elements of MMX have the same speaker and the same hearer. An *informative successor* to US toward S2 is a pair $\langle S2, MM2 \rangle$ where MM2 is the union of MM with some informative sheaf in US toward S2.

Definition 12: Let TS, US, S1, S2 be as above. A *u-successor set* for US toward S2 is the union of

- The simple successor to US,S2.
- A set USS of informative successors to US,S2 with the following property: If M is any inform indicator in S1, then there exists an element $\langle S2, MM \rangle \in USS$ such that $M \in MM$. That is, every inform indicator is attached to at least one successor of US.

A u-successor of a u-situation at time K is a u-situation at time $K + 1$. If US1 is a u-successor of US then $ANCESTOR(US1) = ANCESTORS(US) \cup \{ US \}$.

Definition 13: Let TS be a time-structure of depth K . A *u-situation successor space* for TS is the union over [all u-situations US of depth K in TS] and [all successors S2 of $PHYS(US)$] of some u-successor set for US,S2.

Definition 14: Let $TS = \langle U_SITS, K_ACC \rangle$ be a time-structure of depth K . Let USA and USB be u-situations of depth K in TS. Let US1A be a u-successor of USA and let US1B be a u-successor of USB. Let A be an agent. Then US1B is *possibly knowledge accessible* from US1A relative to A if the following conditions hold:

- $\langle USA, USB \rangle \in K_ACC(A)$.
- For any actional Z and p-situations SXA, SYA, if $OCCURS(DO(A, Z), SXA, SYA)$ and $SXA \leq PHYS(USA)$, then
 - If $SYA < PHYS(USA)$, then there exist SXB, SYB such that $OCCURS(DO(A, Z), SXB, SYB)$ and $SYB < PHYS(USB)$.
 - If $SYA = PHYS(USA)$, then there exists SXB such that $OCCURS(DO(A, Z), SXB, PHYS(USB))$.
 - If $SXA < PHYS(USA) < SYA$, then there exist SXB, SYB such that $OCCURS(DO(A, Z), SXB, SYB)$ and $SXB < PHYS(USB) < SYB$.
 - If $SXA = PHYS(USA) < SYA$, then there exists SYB such that $OCCURS(DO(A, Z), PHYS(USB), SYB)$.
- If there exists a tuple $\langle AS, A, USSQ, SX \rangle$ in $MM(USA)$ and $OCCURS(DO(AS, COMMUNICATE(AH)), SX, PHYS(USA))$ then there exists a p-situation SXB and a tuple $\langle AS, A, USSQ, SXB \rangle$ in $MM(USB)$ and $OCCURS(DO(AS, COMMUNICATE(AH)), SXB, PHYS(USB))$. (That is, if AS has completed informing A of USSQ, then A knows that AS has completed informing him of USSQ.)

Definition 15: Let TS be a time-structure of depth K . A *possible successor* to TS is a pair $TS1 = \langle U_SITS1, K_ACC1 \rangle$ where

- U_SITS1 is a u-situation successor space for TS.

- for each agent $A \in \text{AGENTS}$, $\text{K_ACC1}(A)$ is an equivalence relation over U_SITS1 , which is a subset of the relation, “USB is possibly knowledge accessible from USA.” (Note that, since all the conditions on “possibly knowledge accessible” have the form “Some property holds on US1A iff the corresponding property holds on US1B,” the relation “possibly knowledge accessible relative to (A)” is itself always an equivalence relation.)

TS1 is said to be of depth $K+1$.

Finally, we let this construction go from time 0 to infinity.

Definition 16: Let $\text{TS}_0 = \langle \text{U_SITS}_0, \text{K_ACC}_0 \rangle$, $\text{TS}_1 = \langle \text{U_SITS}_1, \text{K_ACC}_1 \rangle$, \dots be a sequence such that TS_0 is a time structure of depth 0 and for each i , TS_{i+1} is a possible successor for TS_i . Then the pair $\text{TS}_\infty = \langle \text{U_SITS}_\infty, \text{K_ACC}_\infty \rangle = \langle \cup_i \text{U_SITS}_i, \cup_j \text{K_ACC}_j \rangle$ is a *communicative model extension* of \mathcal{M}, \mathcal{I} .

A.3: Interpretation

Let \mathcal{L} , \mathcal{M} and \mathcal{I} be as in the previous section. Let \mathcal{W} be the language \mathcal{L} combined with the following additional elements:

- The sorts “fluent”, “actional” and “actions”, which are super-categories of the sort “physical fluent”, “physical action”, and “physical actional”, respectively.
- The symbols “k_acc”, “sk_acc”, and “inform”.

Let $\text{TS}_\infty = \langle \text{U_SITS}_\infty, \text{K_ACC}_\infty \rangle$ be a communicative model extension of \mathcal{M}, \mathcal{I} .

In this section, we define an interpretation \mathcal{J} of \mathcal{W} in terms of constructions over TS_∞ and \mathcal{M} . For notational convenience, we will write the image of a symbol under \mathcal{J} by writing it in lower-case boldface; thus, for example, $\text{sk_acc} = \mathcal{J}(\text{“sk_acc”})$. We will use ordinary Roman font where symbols are used in prefix notation and are interpreted under \mathcal{J} . For example, if we write “occurs($E, S1, S2$)” we mean the interpretation of “occurs” under \mathcal{J} . Note that, if a symbol is in \mathcal{L} , then its interpretation under \mathcal{I} may be different than its interpretation under \mathcal{J} .

We will first discuss the construction of \mathcal{J} informally and then proceed to the formal definition.

The first issue is fluents. On the one hand, axiom I.5 asserts that every property of situations $\alpha(S)$ has an associated fluent Q_α such that Q_α holds in just those situations satisfying S . The usual extensionalizing trick, therefore, is to identify Q_α with the set of u-situations satisfying α ; generally, to identify fluents with sets of situations. On the other hand, to extend the theory \mathcal{T} to the new model, we must make sure that every physical fluent in \mathcal{T} is still a fluent in the new theory. Moreover it is possible that \mathcal{T} involves the existence of two different fluents that are in fact coextensional in terms of the situations where they hold, but differ in terms of some other property of interest to \mathcal{T} . Therefore, we define a general fluent as a pair of a label and a set of u-situations. For a physical fluent that is, so to speak, grandfathered from \mathcal{T} , the label is just the physical fluent; for all other fluents, the label is immaterial. A physical fluent Q holds in u-situation S just if Q holds in $\text{PHYS}(S)$.

The second issue is the occurrence of actions. For physical actions, as for physical fluents, we use the “PHYS” mapping to guide us; a physical action E occurs from US1 to US2 if E occurs from $\text{PHYS}(\text{US1})$ to $\text{PHYS}(\text{US2})$. For informative events, there are two steps. First, axiom I.4 asserts that “do(as,inform(ah,q1))” and “do(as,inform(ah,q2))” co-occur from us1 to us2 if the intersection of q1 with the set of u-situations that are sk-accessible relative to as,ah from us1 is the

same as the intersection of us2 with that set. Second, the occurrence from us1 to us2 of the act “do(as,inform(ah,q0))” where q0 is a subset of the sk-accessible situations is indicated in the second (MM) field of the u-situation us1.

Finally for simplicity we assume that there are no “pointless coincidences” between \mathcal{M} and the constructions we will use in \mathcal{J} . That is to say: It is conceivable that \mathcal{M} itself happens to contain, as an entity, some tuple that we will want to define as an entity in the denotation of \mathcal{J} . Such a coincidence would cause propositions to be true and false in ways that we don’t intend. One could block this by modifying definition 17 below as follows: Wherever the definition constructs an tuple, add an additional element that is not an element of \mathcal{M} (e.g. \mathcal{M} itself.) That will block any such coincidences. For the sake of readability, I have omitted these.

Otherwise, the definition is pretty much straightforward.

Definition 17: A *general fluent* is a pair $\langle \text{LABEL}, \text{USS} \rangle$ where LABEL is either a physical fluent or 0, and USS is a set of u-situations.

Definition 18: For any PF in PHYSICAL-FLUENTS, define PF_MAP(PF) to be the pair $\langle \text{PF}, \{ \text{US} \mid \text{US} \in \text{U-SITUATIONS} \wedge \text{HOLDS}(\text{PHYS}(\text{US}), \text{PF}) \} \rangle$. Define PF_IMAGES = $\{ \text{PF_MAP}(\text{PF}) \mid \text{PF} \in \text{PHYSICAL-FLUENTS} \}$

Definition 19: We define a general mapping “U2P_MAP” from constructions over TS_∞ to entities in \mathcal{M} as follows:

- If U is a u-situation, then $\text{U2P_MAP}(\text{U}) = \text{PHYS}(\text{U})$.
- If $\text{U} = \langle \text{PF}, \text{USS} \rangle \in \text{PF_IMAGES}$ then $\text{U2P_MAP}(\text{U}) = \text{PF}$.
- If $\text{U} \in \mathcal{M}$ then $\text{U2P_MAP}(\text{U}) = \text{U}$.
- Else $\text{U2P_MAP}(\text{U})$ is undefined.

In reading definition 20 below, keep in mind that, in the standard Tarskian semantics for first-order logic, the denotation of a function or a predicate symbol is a set of tuples. Similarly, we take the denotation of a sort to be a set of entities.

Definition 20: (Long) Let $\mathcal{L}, \mathcal{M}, \mathcal{I}, \mathcal{W}, \mathcal{U}$ be as above. We define the function \mathcal{J} over the sorts and symbols of \mathcal{W} as follows:

Sorts:

$\mathcal{J}(\text{the sort “clock time”}) = \text{the non-negative integers.}$

$\mathcal{J}(\text{the sort “agent”}) = \mathcal{I}(\text{“agent”}).$

$\mathcal{J}(\text{the sort “situation”}) = \text{the set of u-situations in } \mathcal{U}.$

$\mathcal{J}(\text{the sort “fluent”}) = \text{the set of general fluents.}$

$\mathcal{J}(\text{the sort “physical fluent”}) = \text{PF_IMAGES.}$

$\mathcal{J}(\text{the sort “physical actional”}) = \mathcal{I}(\text{“physical actional”})$

$\mathcal{J}(\text{the sort “physical action”}) = \mathcal{I}(\text{“physical action”})$

Let **informative_actionals** $\equiv \{ \langle \text{INFORM}, \text{AH}, \text{Q} \rangle \mid \text{A} \in \text{agent} \wedge \text{Q} \in \text{fluent} \}.$

Let **informative_actions** $\equiv \{ \langle \mathbf{DO}, A, Z \rangle \mid Z \in \mathbf{informative_actionals} \}$

$\mathcal{J}(\text{the sort "actional"}) = \mathcal{I}(\text{"physical actional"}) \cup \mathbf{informative_actionals}$.

$\mathcal{J}(\text{the sort "action'}) = \mathcal{I}(\text{"physical action"}) \cup \mathbf{informative_actions}$.

If σ is any other sort used in \mathcal{L} , then $\mathcal{J}\sigma = \mathcal{I}(\sigma)$.

Non-logical symbols:

$\mathcal{J}("<")$ (as a predicate on clock times) = the usual ordering on integers.

$\mathcal{J}("<")$ (as a predicate on situations) = $\{ \langle S1, S2 \rangle \mid S1, S2 \in \mathbf{situation} \text{ and } S1 \text{ is an ancestor of } S2. \}$

$\mathcal{J}(\text{"holds"}) = \{ \langle S, Q \rangle \mid S \in \mathbf{situation}, Q = \langle PF, USS \rangle \in \mathbf{fluent} \text{ and } S \in USS. \}$

$\mathcal{J}(\text{"time"}) = \{ \langle S, T \rangle \mid S \in \mathbf{situation}, T \in \mathbf{clocktime} \text{ and } S \text{ is of time } T \}$.

$\mathcal{J}(\text{"communicate"}) = \mathcal{I}(\text{"communicate"})$

$\mathcal{J}(\text{"do"}) = \mathcal{I}(\text{"do"}) \cup \{ \langle A, Z, \langle \mathbf{DO}, A, Z \rangle \rangle \mid A \in \mathbf{agent} \text{ and } Z \in \mathbf{informative_actionals} \}$

$\mathcal{J}(\text{"inform"}) = \{ \langle AH, Q, \langle \mathbf{INFORM}, AH, Q \rangle \rangle \mid AH \in \mathbf{agent} \text{ and } Q \in \mathbf{fluent} \}$

$\mathcal{J}(\text{"k_acc"}) = \{ \langle A, S1, S2 \rangle \mid A \in \mathbf{agents} \text{ and } \langle S1, S2 \rangle \in \mathbf{K_ACC}_\infty(A). \}$

$\mathcal{J}(\text{"sk_acc"}) =$
 $\{ \langle AS, AH, SA, SB \rangle \mid$
 exists($S_0 = SA, S_1 \dots S_k = SB$) such that
 for ($i = 1 \dots k$) either **k_acc**(AS, S_{i-1}, S_i) or **k_acc**(AH, S_{i-1}, S_i)
 $\}$.

$\mathcal{J}(\text{"occurs"}) =$
 $\{ \langle E, US1, US2 \rangle \mid$
 $E \in \mathbf{action}$ and $US1, US2 \in \mathbf{situation}$ and
 either [$E \in \mathcal{I}(\text{"physical action"})$ and $\mathbf{OCCURS}(E, \mathbf{PHYS}(US1), \mathbf{PHYS}(US2))$] or
 [there exist ($A, AH \in \mathbf{agent}; Q1, Q2 \in \mathbf{fluent}; USS1, USS2$) such that
 $E = \langle \mathbf{DO}, AS, \langle \mathbf{INFORM}, AH, Q1 \rangle \rangle$ and
 $Q1 = \langle PF1, USS1 \rangle, Q2 = \langle PF2, USS2 \rangle;$
 $USS2 = \{ US \in USS1 \mid \langle AS, AH, US1, US \rangle \in \mathbf{sk_acc} \},$
 $\mathbf{OCCURS}(\mathbf{DO}(AS, \mathbf{COMMUNICATE}(AH)), \mathbf{PHYS}(US1), \mathbf{PHYS}(US2))$ and
 $\langle AS, AH, USS2, \mathbf{PHYS}(US1) \rangle \in \mathbf{MM}(US2)$
 $]$
 $\}$

Let α be any symbol in \mathcal{L} other than those enumerated above. $\mathcal{I}(\alpha)$ is a set of tuples of entities in \mathcal{M} . A tuple T' is a replacement for tuple T if, for each index I , $\mathbf{U2P_MAP}(T'[I]) = T[I]$. Then $\mathcal{J}(\alpha)$ is the set of all replacements R for the tuples in $\mathcal{I}(\alpha)$, such that any two situations in R are ordered under $\mathcal{J}("<")$.

End of definition 20.

Definition 21: The model \mathcal{U} is the union of **clocktime**, **agent**, **situation**, **fluent**, **actional**, **action** and \mathcal{M} .

Note that the function $\text{U2P_MAP}(X)$ is defined for exactly those entities X which are in $\mathcal{J}(\sigma)$ where σ is one of the sorts in the physical language (clock times, situations, agents, physical fluents, physical actionals, physical actions, and other sorts in \mathcal{L}).

A.4: Soundness

Throughout this section: Let \mathcal{L} be a physical language. Let \mathcal{T} be an acceptable physical theory over \mathcal{L} . Let \mathcal{M} be a model and let \mathcal{I} be an interpretation of \mathcal{L} in \mathcal{M} that satisfies \mathcal{T} . Let \mathcal{U} and \mathcal{J} be defined as above.

We will assume that \mathcal{L} is strongly sorted; in particular, that every variable in \mathcal{L} is labelled with its sort. A valuation over variables in \mathcal{L} is required to respect the sort constraint. That is, if μ_i is a variable of sort σ_i , and V is a valuation of μ_i in \mathcal{M} then $V(\mu_i) \in \mathcal{I}(\sigma_i)$. If W is a valuation of μ_i in \mathcal{U} then $W(\mu_i) \in \mathcal{J}(\sigma_i)$.

Lemma 1: For every p-situation PS in \mathcal{M} there exists a u-situation US in \mathcal{U} such that $\text{PHYS}(US)=PS$.

Proof by induction on $\text{TIME}(PS)$. If $\text{TIME}(PS)=0$ then there exists a corresponding u-situation by definition 8. Suppose the statement is true for all PS such that $\text{TIME}(PS)=k$. Let $PS1$ be a p-situation such that $\text{TIME}(PS1)=k+1$. By axiom T.7, $PS1$ is the successor of some situation $PS0$ such that $\text{TIME}(PS0) = k$. By the induction hypothesis, there is a situation $US0$ such that $\text{PHYS}(US0)=PS0$. By definition 9 there is a simple successor $US1$ of $US0$ such that $\text{PHYS}(US1)=PS1$.

Lemma 2: For any u-situation U , $\text{TIME}(\text{PHYS}(U)) = \text{TIME}(U)$. For any two u-situations $U1, U2$ if $U1 < U2$ then $\text{PHYS}(U1) < \text{PHYS}(U2)$.

Proof: Immediate from the definition of $\mathcal{J}("<")$ in definition 20 and the definition of "ANCESTORS" in definitions 7 and 12.

Lemma 3: Let $\mu_1 \dots \mu_k$ be variables in \mathcal{L} . Let V be a valuation mapping each variable μ_i into $\mathcal{I}(\sigma_i)$. Then there exists a valuation W into \mathcal{U} such that $\text{U2P_MAP}(W(\mu_i)) = V(\mu_i)$.

Proof: Immediate from Lemma 1 together with the construction of U2P_MAP and the fact that, for each sort σ , U2P_MAP maps an element of $\mathcal{J}(\sigma)$ to an element of $\mathcal{I}(\sigma)$.

Lemma 4: Let $\alpha(\mu_1 \dots \mu_k)$ be a predicate symbol in \mathcal{L} , including equality. Let W be a valuation from μ_i into \mathcal{U} . Define $V(\mu_i) = \text{U2P_MAP}(W(\mu_i))$. Then $\alpha(\mu_1 \dots \mu_k)$ holds in \mathcal{U} under \mathcal{J} , W if and only if (a) $\alpha(\mu_1 \dots \mu_k)$ holds in \mathcal{M} under \mathcal{I}, V and (b) any two situations $W(\mu_i)$ and $W(\mu_j)$ are ordered under $\mathcal{J}("<")$.

Proof: We must consider separately the cases where α is (A) equality over non-situations; (B) equality over situations; (C) the symbol "<" over clock times; (D) the symbol "<" over situations; (E) the symbol "occurs"; (F) the symbol "holds"; (G) any other predicate symbol in \mathcal{L} .

(A) Equality over non-situations: from definitions 19 and 20.

(B) Equality over situations: Following definitions 19 and 20, this amounts to the claim that $US1=US2$ if and only if $\text{PHYS}(US1)=\text{PHYS}(US2)$ and $US1$ and $US2$ are ordered. The implication from left to right is trivial. For the implication from right to left, consider that, if $US1$ and $US2$ are ordered but $US1 \neq US2$, then either $US1 < US2$ or $US2 < US1$. If $US1 < US2$, then $\text{time}(US1) < \text{time}(US2)$

so by lemma 2 $\text{PHYS}(\text{US1}) \neq \text{PHYS}(\text{US2})$; and likewise if $\text{US2} < \text{US1}$.

(C) The symbol “ $<$ ” over clock times: From the fact that the interpretation is the same under \mathcal{J} as under \mathcal{I} (Definition 20).

(D) The symbol “ $<$ ” over situation: Analogous to (B) above.

(E) The symbol “occurs”. By definition 20, if E is a physical action then $\text{occurs}(E, S1, S2)$ occurs under \mathcal{J} if and only if $\text{occurs}(E, \text{PHYS}(S1), \text{PHYS}(S2))$ under \mathcal{I} .

(F) Let μ_1, μ_2 be variables of sorts “situation” and “physical fluent” respectively. Let $\text{PF} = \text{V}(\mu_1)$. Since $\text{U2P_MAP}(\text{W}(\mu_2)) = \text{V}(\mu_2) = \text{PF}$, by definition 19 $\text{W}(\mu_2) \in \text{PF_IMAGES}$, which, by definition 18 and 19, means that $\text{W}(\mu_2) = \langle \text{PF}, \{ \text{US} \in \text{U-SITUATIONS} \mid \text{HOLDS}(\text{PHYS}(\text{US}), \text{PF}) \} \rangle$. By definition 20 it follows that $\langle \text{W}(\mu_1), \text{W}(\mu_2) \rangle \in \mathcal{J}$ (“holds”) if and only if $\langle \text{V}(\mu_1), \text{PF} \rangle \in \mathcal{I}$ (“holds”)

(G) α is any other predicate symbol in \mathcal{L} . Immediate from definition 20.

Lemma 5: Let $\beta(\mu_1 \dots \mu_k)$ be a function symbol in \mathcal{L} . Let W be a valuation from μ_i into \mathcal{U} such that, for any two situational variables μ_p and μ_q , $\text{W}(\mu_p)$ and $\text{W}(\mu_q)$ are ordered with respect to \mathcal{J} (“ $<$ ”). Define $\text{V}(\mu_i) = \text{U2P_MAP}(\text{W}(\mu_i))$. Then the value of $\beta(\mu_1 \dots \mu_k)$ in \mathcal{M} under \mathcal{I} , V is the image under U2P_MAP of the value of $\beta(\mu_1 \dots \mu_k)$ in \mathcal{U} under \mathcal{J}, W .

Proof: As in the proof of lemma 4, we must consider separately the cases where β is (A) the function symbol “do”; (B) the function symbol “time”; (C) any other function symbol in \mathcal{L} .

(A) By definitions 19 and 20, if A is an agent and Z is a physical action then $\text{U2P_MAP}(\mathcal{J}(\text{do}(\text{A}, \text{Z}))) = \mathcal{J}(\text{do}(\text{A}, \text{Z})) = \mathcal{I}(\text{do}(\text{A}, \text{Z})) = \mathcal{I}(\text{do}(\text{U2P_MAP}(\text{A}), \text{U2P_MAP}(\text{Z})))$. (Again, we are mildly abusing notation.)

(B) By definitions 19 and 20, if US is a u-situation then $\text{U2P_MAP}(\mathcal{J}(\text{time}(\text{US}))) = \mathcal{J}(\text{time}(\text{US})) = \mathcal{I}(\text{time}(\text{PHYS}(\text{US}))) = \mathcal{I}(\text{time}(\text{U2P_MAP}(\text{US})))$.

(C) Let β be any other function symbol. Let $\langle x_1 \dots x_k, y \rangle$ be any tuple where the x_i and y are entities in the image under \mathcal{J} of the sorts in \mathcal{L} . Then by the last part of definition 20,

$\langle x_1 \dots x_k, y \rangle \in \mathcal{J}(\beta)$ iff $\langle \text{U2P_MAP}(x_1) \dots \text{U2P_MAP}(x_k), \text{U2P_MAP}(y) \rangle \in \mathcal{I}(\beta)$.

But for any terms $\gamma_1 \dots \gamma_k$ and any valuation W from the variables in the γ 's to \mathcal{U} , the denotation of $\beta(\gamma_1 \dots \gamma_k)$ under \mathcal{J}, W is equal to y just if the tuple $\langle \mathcal{J}(\gamma_1) \dots \mathcal{J}(\gamma_k), y \rangle$ is in $\mathcal{J}(\beta)$; and likewise for \mathcal{I} .

Unfortunately, U2P_MAP does not preserve truth-values of predicates over unordered u-situations; it is possible that $\text{U2P_MAP}(\text{US1}) = \text{U2P_MAP}(\text{US2})$ even though $\text{US1} \neq \text{US2}$. or that $\text{U2P_MAP}(\text{US1}) < \text{U2P_MAP}(\text{US2})$ even if US1 and US2 are unordered. There is, moreover, in general no way to modify U2P_MAP to preserve inequality, since the cardinality of the set of u-situations may be larger than the cardinality of p-situations. Therefore, in establishing below that if an open formula with inequalities or orderings is satisfiable in \mathcal{J} then it is also satisfiable in \mathcal{I} , it is necessary to continuously “patch” the mapping U2P_MAP by mapping a u-situation US into some p-situation that is physically indistinguishable from $\text{U2P_MAP}(\text{US})$. Fortunately, we had the foresight to provide ourselves with plenty of these. Stating this exactly is a little involved; definitions 22-24 and corollary 6 through lemma 9 accomplish this.

Definition 22: Let τ be a function from \mathcal{U} to itself which is one-to-one and onto. The function τ is said to be a *physical automorphism* over \mathcal{U} if the following conditions hold:

1. If X is not a u-situation, then $\tau(X) = X$.
2. Let $\alpha(\mu_1 \dots \mu_k)$ be any atomic formula in \mathcal{L} with free variables $\mu_1 \dots \mu_k$. Let W and Y be

valuations from μ_i to \mathcal{U} such that $Y(\mu_i) = \tau(W(\mu_i))$. Then Y satisfies α only if W satisfies α .

Note that condition (2) only applies to formulas in the *physical* language \mathcal{L} , not in the broader language.

Definition 23: Let S_1, S_2 be either two p-situations or two u-situations. Situation S is the *latest common ancestor* (LCA) of S_1 and S_2 , if $S \leq S_1, S \leq S_2$ and S is the latest situation with that property. Since the order relation on situations is a forest of trees, any two situations have at most one latest common ancestor

Definition 24: Let $\langle \mu_1 \dots \mu_k \rangle$ be a k-tuple of variables. Let W be a valuation of the μ 's to \mathcal{U} and let V be a valuation of the μ 's to \mathcal{M} . V is said to be an *image* of W if the following conditions hold:

- If μ is not a situational variable, then $V(\mu) = \text{U2P_MAP}(W(\mu))$.
- There exists a physical automorphism τ over \mathcal{U} such that, for each pair of situational variables μ_i, μ_j , if S is the latest common ancestor of $W(\mu_i), W(\mu_j)$ then $\text{PHYS}(\tau(S))$ is the LCA of $V(\mu_i), V(\mu_j)$; and if $W(\mu_i)$ and $W(\mu_j)$ have no common ancestor, then $V(\mu_i)$ and $V(\mu_j)$ have no common ancestor.

We say that the automorphism τ *establishes the correspondence* between W and V .

Corollary 6: Let $\mu_1 \dots \mu_k, W, V$, and τ be as in definition 24. For each i , $V(\mu_i) = \text{U2P_MAP}(\tau(W(\mu_i)))$.

If μ_i is a situational variable, then applying definition 24 and choosing $j = i$, since $W(\mu_i)$ is the LCA of $W(\mu_i)$ and itself, we have $\text{U2P_MAP}(\tau(W(\mu_i))) = \text{PHYS}(\tau(W(\mu_i))) = \text{LCA}(V(\mu_i), V(\mu_i)) = V(\mu_i)$. If μ_i is not a situational variable, then the result is immediate.

Lemma 7: Let μ_1 and μ_2 be situational variables in \mathcal{L} . Let W and V be valuations of μ_1, μ_2 to \mathcal{U} and \mathcal{M} respectively, and let V be an image of W . Then $W(\mu_1) = W(\mu_2)$ if and only if $V(\mu_1) = V(\mu_2)$ and $W(\mu_1) < W(\mu_2)$ if and only if $V(\mu_1) < V(\mu_2)$

Proof: Let τ be an automorphism that establishes the correspondence between W and V . If $W(\mu_1) = W(\mu_2)$ then $V(\mu_1) = V(\mu_2)$, since $V(\mu) = \text{PHYS}(\tau(W(\mu)))$ and is thus a function of $W(\mu)$. If $W(\mu_1) < W(\mu_2)$ then by lemma 2, $V(\mu_1) < V(\mu_2)$.

Suppose that $V(\mu_1) = V(\mu_2)$. Thus, $\text{LCA}(V(\mu_1), V(\mu_2)) = V(\mu_1) = V(\mu_2)$. By definition 24 $\text{LCA}(W(\mu_1), W(\mu_2)) = W(\mu_1) = W(\mu_2)$.

Suppose that $V(\mu_1) < V(\mu_2)$. Thus, $\text{LCA}(V(\mu_1), V(\mu_2)) = V(\mu_1)$. By definition 24, $\text{LCA}(W(\mu_1), W(\mu_2)) = W(\mu_1)$. Therefore $W(\mu_1) \leq W(\mu_2)$. Since $V(\mu_1) \neq V(\mu_2)$, it follows from the earlier part of this lemma that $W(\mu_1) \neq W(\mu_2)$; hence $W(\mu_1) < W(\mu_2)$.

Lemma 8: Let $\alpha(\mu_1 \dots \mu_k)$ be a predicate symbol in \mathcal{L} . Let W be a valuation of the μ 's to \mathcal{U} and let V be an image of W . Then α holds in \mathcal{U} under W if and only if α holds in \mathcal{M} under V .

Proof: Let τ be an automorphism that establishes the correspondence between W and V . Let $Q(\mu_i) = \tau(W(\mu_i))$. By definition 22, $\alpha(\mu_1 \dots \mu_k)$ holds under \mathcal{J}, W if and only if it holds under \mathcal{J}, Q . By lemma 4, $\alpha(\mu_1 \dots \mu_k)$ holds under \mathcal{J}, Q if and only if it holds under \mathcal{I}, V and for any two situational variables μ_a, μ_b , $Q(\mu_a)$ and $Q(\mu_b)$ are ordered. By lemma 7, $Q(\mu_a)$ and $Q(\mu_b)$ are ordered if and only if $V(\mu_a)$ and $V(\mu_b)$ are ordered; and by condition 5 of definition 6, $\alpha(\mu_1 \dots \mu_k)$ holds under

\mathcal{I}, V only if $V(\mu_a)$ and $V(\mu_b)$ are ordered. Putting these together, it follows that $\alpha(\mu_1 \dots \mu_k)$ holds under \mathcal{J}, W if and only if it holds under \mathcal{I}, V .

Lemma 9: Let $\beta(\mu_1 \dots \mu_k)$ be a function symbol in \mathcal{L} , and let μ_{k+1} be another variable. Let W be a valuation of the μ 's to \mathcal{U} and let V be an image of W . Then the equation $\mu_{k+1} = \beta(\mu_1 \dots \mu_k)$ holds in \mathcal{U} under W if and only if it holds in \mathcal{M} under V .

Proof: Exactly analogous to the proof of lemma 8, substituting lemma 5 for lemma 4.

Lemma 10: Let $\alpha(\mu_1 \dots \mu_k)$ be a quantifier-free formula in \mathcal{L} . Let W be a valuation of the μ 's to \mathcal{U} and let V be an image of W . Then α holds in \mathcal{U} under W if and only if α holds in \mathcal{M} under V .

Proof: Straightforward structural induction over the form of α , using lemmas 8 and 9.

Lemma 11: Let $\mu_1 \dots \mu_k$ be variables whose sorts are in \mathcal{L} . Let W be a valuation from variables $\mu_1 \dots \mu_k$ to \mathcal{U} and let V be an image of W . (We will include here the case where $k = 0$; in that case, W and V are null valuations.) Let μ_{k+1} be a new variable of sort σ_{k+1} .

1. Let A be an entity in $\mathcal{J}(\sigma_{k+1})$. Let $W' = W \cup \{ \mu_{k+1} \rightarrow A \}$. Then there exists B in \mathcal{M} such that $V' = V \cup \{ \mu_{k+1} \rightarrow B \}$ is an image of W' .
2. Let B be an entity in $\mathcal{I}(\sigma_{k+1})$. Let $V' = V \cup \{ \mu_{k+1} \rightarrow B \}$. Then there exists A in \mathcal{U} such that V' is an image of $W' = W \cup \{ \mu_{k+1} \rightarrow A \}$.

Proof:

Let τ be a physical automorphism over \mathcal{U} that establishes the correspondence of W and V . If the sort of μ_{k+1} is not a situation, then both (1) and (2) are trivial; one can take $A=B$, leave the automorphism τ unchanged, and the result is immediate from the definitions. Therefore, we may assume that the sort of μ_{k+1} is a situation, and therefore A is a u-situation and B is a p-situation. Without loss of generality, renumber the variables $\mu_1 \dots \mu_k$ so that $\mu_1 \dots \mu_m$ are situational variables and the rest are not situational variables.

In both halves of the lemma, in order to show that W' is an image of V' we must exhibit an automorphism τ' that establishes this correspondence.

Let us write $PT(S) = PHYS(\tau(S))$, and $S_i = W(\mu_i)$ for $i = 1 \dots m$.

Part 1. There are three cases:

Case A. $m = 0$. In this case, one can choose $B=PHYS(A)$, and τ' to be the identity automorphism.

Case B. Suppose that $A \leq S_i$ for some i . Let $\tau' = \tau$, and let $B=PT(A)$. For any j , let S be the LCA of S_j and A . There are four cases:

- B.i. $S_j \leq A$. In this case $S = S_j$. Since W is an image of V under τ , $PT(S) = PT(S_j) = V(\mu_j)$.
- B.ii. $A \leq S_j$. In this case $S=A$. Since τ is an automorphism, $\tau(S) \leq \tau(S_j)$. By lemma 2, $PT(S) = PT(A) \leq PT(S_j)$ so $PT(S)$ is the LCA of $PT(A)$ and $PT(S_j)$.
- B.iii. A and S_j are unordered but have LCA S . Then S is the LCA of S_i and S_j , so $PT(S)$ is the LCA of $PT(S_i)$ and $PT(S_j)$. Since $PT(S) < PT(A) \leq PT(S_i)$, it follows that $PT(S)$ is the LCA of $PT(A)$ and $PT(S_j)$.

B.iv. A and S_j have no common ancestor. Hence S_i and S_j have no common ancestor. Hence $\text{PT}(S_i)$ and $\text{PT}(S_j)$ have no common ancestor. Hence $\text{PT}(A) < \text{PT}(S_i)$ and $\text{PT}(S_j)$ have no common ancestor.

Case C. Suppose that A does not precede any of the S_j . Consider the set $\text{LL} = \text{LCA}(A, S_1) \dots \text{LCA}(A, S_m)$. If LL is non-empty, let S be the latest situation in LL. We have three cases:

C.i. LL is empty; that is, none of the S_j are ordered with respect to A. Then none of the values of $\tau(S_j)$ are ordered with respect to $\tau(A)$, so by lemma 4, none of the values of $\text{PT}(S_j)$ are ordered with respect to $\text{PT}(A)$. Hence, we may choose $\tau' = \tau$ and $\text{B} = \text{PT}(A)$.

C.ii. S is equal to one of the S_i . Then for each S_j , $\text{LCA}(A, S_j) = \text{LCA}(S_i, S_j)$. Thus, again, we may choose $\tau' = \tau$ and $\text{B} = \text{PHYS}(\tau(A))$.

C.iii. S is not equal to any of the S_i . Note that at least there must be one of the $S_j > S$; call this S_x . Let Q be the successor of S such that $Q \leq A$. There are two cases:

C.iii.a. Q is not a communicative successor of S. Then $\tau(Q)$ is not a communicative successor of $\tau(S)$. For any S_j , if $S < S_j$, let Q_j be the successor of S such that $Q_j \leq S_j$. By the construction in definitions 9-12, it follows that $\text{PT}(Q)$ is not equal to $\text{PT}(Q_j)$. Therefore $\text{PT}(S)$ is the LCA of $\text{PT}(A)$ and $\text{PT}(S_j)$. If S_j is not ordered with respect to S, then the LCA of S_j and A is the same as the LCA of S_j and S_x (or neither of these LCA's exists), so again $\text{LCA}(\text{PT}(A), \text{PT}(S_j)) = \text{LCA}(\text{PT}(S_x), \text{PT}(S_j)) = \text{PHYS}(\text{LCA}(\tau(S_x), \tau(S_j))) = \text{PHYS}(\text{LCA}(\tau(A), \tau(S_j)))$. Therefore we can choose $\tau' = \tau$ and $\text{B} = \text{PHYS}(\tau(A))$.

C.iii.b. Q is a communicative successor of S. Here, finally, is the case where τ may need to be modified. Let $Q_1 \dots Q_p$ be all the successors of S that precede one of the S_i . By property (4) of definition 6, there are infinitely many successors of $\text{PT}(S)$ that are physically indistinguishable from $\text{PT}(Q)$. Let C be one such that is not equal to $\text{PT}(Q_i)$ for any i . Let ω be the automorphism of \mathcal{M} that interchanges the subtree of p-situations following C with the subtree of p-situations following $\text{PT}(Q)$ and leaves the rest of \mathcal{M} the same (see definition 5). Let $\tau' = \tau \circ \omega$. Let $\text{B} = \text{PHYS}(\tau'(A))$. Now, suppose $S_j > S$. Then the LCA of S_j and $A = S$. Since $\text{PHYS}(\tau'(A))$ is a descendant of C, which is a successor of $\text{PHYS}(\tau(S))$ and $\text{PHYS}(\tau'(S_j))$ is a descendant of $\text{PHYS}(\tau(Q_j))$ which is a different successor of $\text{PHYS}(\tau(S))$, it follows that the $\text{LCA}(\text{PHYS}(\tau'(A)), \text{PHYS}(\tau'(S_j))) = \text{PHYS}(\tau(S))$. Alternatively, if S_j is not ordered with respect to S, then we still have $\text{LCA}(\text{PHYS}(\tau(A)), \text{PHYS}(\tau(S_j))) = \text{PHYS}(\text{LCA}(\tau(A), \tau(S_j)))$, by exactly the same argument as in case C.iii.a.

Part 2. The proof of part 2 is exactly analogous to that of part 1, but going in the opposite direction. ■

Lemma 12: Let α be a prenex formula in \mathcal{L} with m quantifiers and k free variables $\mu_1 \dots \mu_k$. Let W be a valuation from variables $\mu_1 \dots \mu_k$ to \mathcal{U} and let V be an image of W. Then α is true under \mathcal{J}, W if and only if it is true under \mathcal{I}, V .

Proof by induction on m , the number of quantifiers.

If $m = 0$, then the statement is just lemma 10.

Suppose the statement is true for all formulas with m quantifiers. Let α be a formula with $m + 1$ quantifiers. There are four cases:

Case 1: α is true under \mathcal{J}, W and α has the form " $\exists X \beta(X)$ ", where β is a formula with m quantifiers and $k + 1$ free variables. Since α is true, there exists an entity $A \in \mathcal{U}$ and a valuation $W' =$

$W \cup \{X \rightarrow A\}$ such that β is true under \mathcal{J}, W' . By lemma 11 there exists a valuation V' that is an image of W' . By the inductive hypothesis, β is true under \mathcal{I}, V' . Hence α (that is, $\exists_X \beta$) is true under \mathcal{I}, V .

Case 2: α is true under \mathcal{I}, V and α has the form “ $\exists_X \beta(X)$ ”. Since α is true, there exists an entity $B \in \mathcal{M}$ and a valuation $V' = V \cup \{X \rightarrow B\}$ such that β is true under \mathcal{I}, V' . By lemma 11 there exists a valuation W' such that V' is an image of W' . By the inductive hypothesis, β is true under \mathcal{J}, W' . Hence α is true under \mathcal{J}, W .

Case 3: α is true under \mathcal{J}, W and α has the form “ $\forall_X \beta(X)$ ”. Let γ be the transformation into prenex form of $\neg\alpha$. Then γ is false under \mathcal{J}, W , and γ has the form “ $\exists_X \delta$ ” where δ is the prenex form of $\neg\beta$. By the contrapositive to case 2 above, γ is false under \mathcal{I}, V ; hence α is true under \mathcal{I}, V .

Case 4: α is true under \mathcal{I}, V and α has the form “ $\forall_X \beta(X)$ ”. Exactly analogous to case (4), but using the contrapositive to case 1.

Corollary 13: All the physical axioms of \mathcal{T} , axioms T.1-T.7, and axioms T.8 and T.9 restricted to physical actions are true in \mathcal{U} under interpretation \mathcal{J} .

Proof: Immediate from lemma 12, taking $k = 0$. and using the fact that the axioms in \mathcal{T} and axioms T.1-T.9 are true in \mathcal{M} (by definition of \mathcal{M}).

Lemma 14: If $\text{PS1}=\text{PHYS}(\text{US1})$ and PS1 and PSZ are ordered, then there exists USZ such that US1 and USZ are ordered, and $\text{PSZ}=\text{PHYS}(\text{USZ})$.

Proof: If $\text{PS1}=\text{PSZ}$ then $\text{USZ}=\text{US1}$.

If $\text{PSZ} < \text{PS1}$, then let USZ be the ancestor of US1 at time $\text{TIME}(\text{PSZ})$.

If $\text{PS1} < \text{PSZ}$, then let $s_1 = \text{PS1}, s_2 \dots s_k = \text{PSZ}$ be p-situations such that s_{i+1} is a successor of s_i . Using definition 9 iteratively, let US_2 be the simple successor to US1 parallel to PS2 , let US_3 be the simple successor to US_2 parallel to PS3 , and so on. Then US_k satisfies the desired conditions on USZ .

Lemma 15: Axioms T.8 extended to general actions and K.1—K.8 are true in \mathcal{U} under \mathcal{J} . (I’m just bunching together the axioms whose proof is easy.)

Proof:

T.8 Immediate from the definition of \mathcal{J} (“occurs”) (Definition 20).

K.1–K.3. Immediate from definition 15, which requires $\text{K-ACC}(A)$ to be an equivalence relation on u-situations.

K.4–K.6 Immediate from definition 14, which restricts the “possibly accessible” on situations that hold on the left-hand side of each of these relations to those that satisfy the conditions on the right-hand side of these implications; plus definition 15, which states that the actual knowledge accessibility relation are a subset of the possibly accessible relations.

K.7, K.8. Immediate from the definition of \mathcal{J} (“sk_acc”) in definition 20.

Lemma 16: Axiom I.1 is true in \mathcal{U} under \mathcal{J} .

Proof:

By the definition of \mathcal{J} (“occurs”) in definition 20, if $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1}, \text{US2})$ then there exist $\text{QA}, \text{PF1}, \text{USS1}, \text{PFA}, \text{USSA}$ such that $\text{Q1} = \langle \text{PF1}, \text{USS1} \rangle$, $\text{QA} = \langle \text{PFA}, \text{USS2} \rangle$, $\text{USSA} = \{ \text{US} \in \text{USS1} \mid \langle \text{AS}, \text{AH}, \text{US1}, \text{US} \rangle \in \mathbf{k_acc} \}$, and $\langle \text{AS}, \text{AH}, \text{USS2}, \text{S1}, \text{PHYS}(\text{US2}) \rangle \in \text{MM}(\text{US2})$. Let USQ be the successor of US1 such that $\text{USQ} \leq \text{US2}$. By definition 9, $\langle \text{AS}, \text{AH}, \text{USS2}, \text{S1}, \text{PHYS}(\text{US2}) \rangle \in \text{MM}(\text{USQ})$. By definition 10 and 11, $\text{OCCURS}(\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH})), \text{PHYS}(\text{US1}), \text{PHYS}(\text{US2}))$. By definition 20, $\text{occurs}(\text{do}(\text{AS}, \text{communicate}(\text{AH})), \text{US1}, \text{US2})$.

Lemma 17: Axiom T.9 extended to general actions is true in \mathcal{U} under \mathcal{J} .

Proof:

Let $\text{US1}, \text{US2}, \text{USX}$, and USY be u-situations and E an event such that $\text{occurs}(\text{E}, \text{US1}, \text{US2})$, $\text{US1} < \text{USX} < \text{US2}$ and $\text{USX} < \text{USY}$. Let $\text{S1} = \text{PHYS}(\text{US1})$ and $\text{S2} = \text{PHYS}(\text{US2})$. By definition 20, E is either a physical action or an informative action. The case where E is a physical action is covered in corollary 13. Suppose that E is an informative action; let $\text{E} = \langle \text{DO}, \text{AS}, \langle \text{INFORM}, \text{AH}, \text{Q1} \rangle \rangle$. By definition 20 there exist $\text{QA}, \text{PF1}, \text{USS1}, \text{PFA}, \text{USSA}$ such that $\text{Q1} = \langle \text{PF1}, \text{USS1} \rangle$, $\text{QA} = \langle \text{PFA}, \text{USS2} \rangle$, $\text{USSA} = \{ \text{US} \in \text{USS1} \mid \langle \text{AS}, \text{AH}, \text{US1}, \text{US} \rangle \in \mathbf{k_acc} \}$, and $\langle \text{AS}, \text{AH}, \text{USSA}, \text{S1} \rangle \in \text{MM}(\text{US2})$. By definition 9, $\langle \text{AS}, \text{AH}, \text{USSA}, \text{S1} \rangle \in \text{MM}(\text{USX})$. By axiom T.9 applied to the action $\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH}))$ there exists a situation SZ such that $\text{ordered}(\text{SZ}, \text{PHYS}(\text{SY}))$, $\text{SZ} > \text{SX}$, and $\text{OCCURS}(\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH}, \text{S1}, \text{SZ}))$. By lemma 14, there exists USZ such that $\text{PHYS}(\text{USZ}) = \text{SZ}$ and USZ is ordered with respect to USY . It follows that $\text{USZ} > \text{USX}$ and that $\langle \text{AS}, \text{AH}, \text{USS2}, \text{S1} \rangle \in \text{MM}(\text{USZ})$. By definition 20, $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1}, \text{USZ})$.

Lemma 18: Axiom I.2 is true in \mathcal{U} under \mathcal{J} .

Proof:

Let AS, AH be agents, let $\text{US1}, \text{US2}$ be u-situations, and let $\text{Q} = \langle \text{PF}, \text{USSQ} \rangle$ be a general fluent. Let $\text{US1ACC} = \{ \text{USA} \mid \langle \text{AS}, \text{AH}, \text{US1}, \text{USA} \rangle \in \mathbf{sk_acc} \}$, the set of situations accessible from US1 in the shared knowledge of AS and AH . Let $\text{USSA} = \text{USSQ} \cap \text{US1ACC}$, the set of situations satisfying Q that are knowledge accessible from S1 , relative to the shared knowledge of AS and AH . Let $\text{S1} = \text{PHYS}(\text{US1})$.

Suppose that $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1}, \text{US2})$. By definition 20 (denotation of “occurs”), the tuple $\langle \text{AS}, \text{AH}, \text{USSA}, \text{S1} \rangle \in \text{MM}(\text{US2})$. Let USY be the successor of US1 that is an ancestor of US2 . By definitions 9, 11, and 12 it follows that $\text{MM}(\text{USY})$ contains the tuple $\langle \text{AS}, \text{AH}, \text{USSA}, \text{S1} \rangle$. By definition 10, USSA is a possible communicative content for S1 from AS to AH ; hence, by definition 10, every situation that is knowledge accessible from US1 relative to AS is an element of USSA and therefore an element of $\text{USSQ} \supset \text{USSA}$. By definition 20 (“holds”) Q holds in every situation accessible from US1 .

Conversely, if Q holds in every situation accessible from S1 , then USSA is a possible communicative content from AS to AH . Suppose that $\text{OCCURS}(\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH})), \text{S1}, \text{S2})$. Let SY be the successor of S1 such that $\text{SY} \leq \text{S2}$. By definition 12, there exists an informative successor USY of US1 such that $\langle \text{AS}, \text{AH}, \text{USSA}, \text{S1} \rangle \in \text{MM}(\text{USY})$. By axiom T.9 there exists a situation $\text{USZ} \geq \text{USY}$ such that $\text{OCCURS}(\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH})), \text{US1}, \text{USZ})$. By definitions 9, 11, 12 $\langle \text{AS}, \text{AH}, \text{USSA}, \text{S1} \rangle \in \text{MM}(\text{USZ})$. By definition 20, $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{S1}, \text{SZ})$.

Lemma 19: Axiom I.3 is true in \mathcal{U} under \mathcal{J} .

Proof: Assume that $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1}, \text{US2})$ and that $\text{k_acc}(\text{AH}, \text{US2}, \text{US2A})$. We need to prove that there exists a situation US1A such that $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1A}, \text{US2A})$ and $\text{k_acc}(\text{AH}, \text{US1}, \text{US1A})$.

Define USSA as in the proof of lemma 18. By definition 20 (denotation of “occurs”) since $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1}, \text{US2})$ it follows that the tuple $\langle \text{AS}, \text{AH}, \text{USSA}, \text{PHYS}(\text{US1}) \rangle \in \text{MM}(\text{US2})$ and $\text{OCCURS}(\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH})), \text{PHYS}(\text{US1}), \text{PHYS}(\text{US2}))$. By definition 15, since $\text{k_acc}(\text{AH}, \text{US2}, \text{US2A})$, US2A is possibly knowledge accessible from US2 relative to AH . By definition 14, the tuple $\langle \text{AS}, \text{AH}, \text{USSA}, \text{PS1A} \rangle \in \text{MM}(\text{US2A})$ for some p-situation $\text{PS1A} < \text{PHYS}(\text{US2A})$, and $\text{OCCURS}(\text{DO}(\text{AS}, \text{COMMUNICATE}(\text{AH})), \text{PS1A}, \text{PHYS}(\text{US2A}))$. By theorem 3 and axiom K.8, any two situations that are sk_acc are at the same time. Hence, all the situations in USSA are at the same time, and by definition 10 this time must be equal to $\text{TIME}(\text{US1})$ and to $\text{TIME}(\text{US1A})$. Hence $\text{TIME}(\text{US1}) = \text{TIME}(\text{US1A})$. By axiom A.4, since $\text{k_acc}(\text{AH}, \text{US2}, \text{US2A})$, $\text{US1} < \text{US2}$, $\text{US1A} < \text{US2A}$ and $\text{TIME}(\text{US1}) = \text{TIME}(\text{US1A})$ it follows that $\text{k_acc}(\text{AH}, \text{US1}, \text{US1A})$. Hence, the set of situations that are accessible relative to the shared knowledge of AS and AH is the same starting from US1 as starting from US1A . Hence the act of AS informing AH of Q starting in US1A uses the tuple $\langle \text{AS}, \text{AH}, \text{USSA}, \text{PS1A} \rangle$. Thus by definition 20, $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{Q})), \text{US1A}, \text{US2A})$.

Lemma 20: Axiom I.4 is true in \mathcal{U} under \mathcal{J} .

Proof: Suppose that $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{QX})), \text{US1}, \text{US2})$ and that QY is a fluent. Let US3 be the successor of US1 such that $\text{US3} \leq \text{US2}$. Let $\text{QX} = \langle \text{PFX}, \text{USSQX} \rangle$; $\text{QY} = \langle \text{PFY}, \text{USSQY} \rangle$; $\text{QXA} = \text{USSQX} \cap \{ \text{USA} \mid \text{sk_acc}(\text{AS}, \text{AH}, \text{US1}, \text{USA}) \}$, and $\text{QYA} = \text{USSQY} \cap \{ \text{USA} \mid \text{sk_acc}(\text{AS}, \text{AH}, \text{US1}, \text{USA}) \}$. By definition 20, $\langle \text{AS}, \text{AH}, \text{QXA}, \text{PHYS}(\text{US1}) \rangle \in \text{MM}(\text{US2})$. By definitions 9, 11, 12, $\langle \text{AS}, \text{AH}, \text{QXA}, \text{PHYS}(\text{US1}) \rangle \in \text{MM}(\text{US3})$.

I. (Left to right in the two-way implication.) Suppose that $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{QY})), \text{US1}, \text{US2})$. By the same argument as above $\langle \text{AS}, \text{AH}, \text{QYA}, \text{PHYS}(\text{US1}) \rangle \in \text{MM}(\text{US3})$. But by definition 11, US3 contains at most one inform indicator with starting point $\text{PHYS}(\text{US1})$, speaker AS , and hearer AH . Hence $\text{QXA} = \text{QYA}$. That is, if situation USA is accessible from US1 relative to the shared knowledge of AS and AH , then QX holds in USA iff QY holds in USA .

II. (Right to left in the two-way implication.) If it is the case that $\forall_{S1A} \text{sk_acc}(\text{AS}, \text{AH}, S1, S1A) \Rightarrow [\text{holds}(S1A, \text{QX}) \Leftrightarrow \text{holds}(S1A, \text{QY})]$ then $\text{QXA} = \text{QYA}$, so by definition 20, $\text{occurs}(\text{do}(\text{AS}, \text{inform}(\text{AH}, \text{QY})), \text{US1}, \text{US2})$.

Lemma 21: Axiom I.5 (the comprehension axiom) is true in \mathcal{U} under \mathcal{J} .

Proof: Immediate from definitions 17 and 20, using the comprehension axiom of set theory.

Considering how problematic the comprehension axiom would seem to be it may be surprising that it has a one-line proof. In fact, one might say that the whole construction we went through in section A.3 is precisely tailored so that the comprehension axioms *should* have a one-line proof. Nonetheless the reader may well have legitimate worries about such a powerful axiom, that are hardly assuaged by the above proof. Let me therefore discuss further how this whole construction works.

The key point is this: There is *no circularity whatever* in the whole structure of definitions given in section 3. The structure of u-situations is built up iteratively forward in time. The label on an “inform” action A is a set of u-situations contemporaneous with the start of A ; it gives rise to a new u-situations at the next point in time. Iterating from 1 to infinity gives us a well-defined and fixed set \mathcal{U} of all u-situations. Definition 17 defines a fluent as a subset of \mathcal{U} . Definition 20 defines the occurrence of an inform action in terms of these fluents and of the labels on the actions.

More generally, definition 20 defines the denotation of every symbol in \mathcal{W} extensionally, in terms of structures over \mathcal{U} and \mathcal{M} and the interpretation \mathcal{I} ; no aspect of \mathcal{J} is defined in terms of \mathcal{J} itself (except as a convenient abbreviation.) Having adopted definition 20, \mathcal{J} is now fixed, and it is fixed which fluents satisfy which formulas under \mathcal{J} .

But isn't it inherently circular to say, for example,

q1 is the fluent such that
 $\forall_S \text{ holds}(S, \text{q1}) \Leftrightarrow \exists_{AS, AH, S2, Q} \text{ occurs}(\text{do}(AS, \text{inform}(AH, Q)), S, S2)$

considering that the quantification over Q contains q1 itself? Not at all, no more than saying

0 is the number such that, $\forall_X, X + 0 = X$

when the quantification over X includes 0 itself. The formula above is just a description of q1, and the axioms are sufficient to guarantee that a q1 satisfying this definition exists.

Theorem 1:

Let \mathcal{T} be an acceptable physical theory, and let \mathcal{A} be \mathcal{T} together with axioms K.1 — K.8 and I.1 — I.5, and with T.8 and T.9 extended to arbitrary actions. Then \mathcal{A} is consistent.

Proof: We have shown that a model and an interpretation satisfying \mathcal{A} can be constructed. ■

Theorem 2: Let \mathcal{T} be an acceptable physical theory, and let \mathcal{U} be the union of:

- A. \mathcal{T} ;
- B. Axioms K.1 — K.7 and I.1 — I.5.
- C. A collection of domain-specific knowledge acquisition axioms of the form specified in section ??.
- D. The frame axiom I.6 associated with the axioms in (C).
- E. Any set of axioms \mathcal{K} specifying the presence or absence of k_acc relations among situations at time 0 as long as:
 - i. The axioms in \mathcal{K} do not refer to any situations of time later than 0.
 - ii. The axioms in \mathcal{K} are consistent with \mathcal{T} , axioms K.1 — K.3, K.5 (as regards knowing the feasibility of actions at time 0); and the axioms in (C).

Then \mathcal{U} is consistent.

Sketch of Proof: The proof of theorem 1 needs to be modified as follows:

- In definition 8, initialize the K_ACC function at time 0 to satisfy the union of the axioms in (E) with the axioms enumerated in E.ii.
- In definition 14, add to the conditions on US1B being possibly knowledge accessible from US1A:

For each axiom in (C) of the form “ A always knows whether $\Phi_i(A, S)$,” the condition $\Phi(\text{US1B}) \Leftrightarrow \Phi(\text{US1A})$ must hold.

- Modify the second bullet in definition 15 to read, “For each agent A , $K_ACC1(A)$ is the relation over u -situations, ‘ $US1B$ is knowledge accessible from $US1A$ relative to A .’”

The proof that the additional axioms enumerated in theorem 2 are satisfied is then straightforward. ■