

# $(\mathcal{C}+)^{++}$ : An Action Language for Modelling Norms and Institutions

Marek Sergot

Department of Computing  
Imperial College, London  
mjs@doc.ic.ac.uk

Technical Report 2004/8  
June 2004; revised April 2005

## Abstract

The language  $\mathcal{C}+$  of Giunchiglia, Lee, Lifschitz, McCain, and Turner (2004) is a formalism for specifying and reasoning about the effects of actions and the persistence ('inertia') of facts over time. An 'action description' in  $\mathcal{C}+$  is a set of  $\mathcal{C}+$  laws which define a labelled transition system of a certain kind. This document presents  $(\mathcal{C}+)^{++}$ , an extended form of  $\mathcal{C}+$  designed for representing norms of behaviour and institutional aspects of (human or computer) societies. There are two main extensions. The first is a means of expressing 'counts as' relations between actions, also referred to as 'conventional generation' of actions. The second is a way of specifying the permitted (acceptable, legal) states of a transition system and its permitted (acceptable, legal) transitions.

## 1 Introduction

The language  $\mathcal{C}+$  of Giunchiglia, Lee, Lifschitz, McCain, and Turner (2004) is a formalism for specifying and reasoning about the effects of actions and the persistence ('inertia') of facts over time. An 'action description' in  $\mathcal{C}+$  is a set of  $\mathcal{C}+$  laws which define a labelled transition system of a certain kind. Implementations supporting a wide range of querying and planning tasks are available, notably in the form of the 'Causal Calculator' CCALC<sup>1</sup>. This document presents  $(\mathcal{C}+)^{++}$ , an extended form of  $\mathcal{C}+$  designed for representing norms of behaviour and institutional aspects of (human or computer) societies. There are two main extensions. The first is a means of expressing 'counts as' relations between actions, also referred to as 'conventional generation' of actions. The second is a way of specifying the permitted (acceptable, legal) states of a transition system and its permitted (acceptable, legal) transitions.

There are two aims: to investigate how far one can get by building on transition systems, and to determine how far one can get by building on the language  $\mathcal{C}+$  in particular.

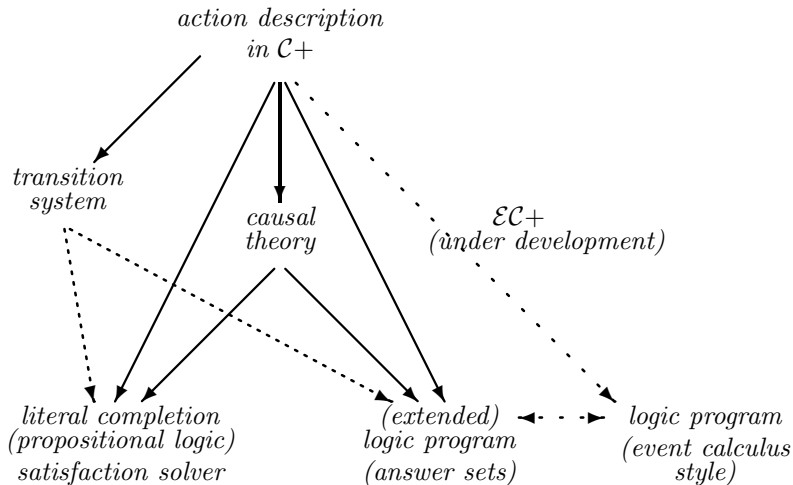
The language  $\mathcal{C}$  was introduced by Giunchiglia and Lifschitz (1998). It applies the ideas of 'causal theories' (Lifschitz, 1997; Turner, 1999) to reasoning about the effects of actions and the persistence ('inertia') of facts ('fluents'), building on earlier suggestions by McCain and Turner (1997). A summary of  $\mathcal{C}$  is included

---

<sup>1</sup><http://www.cs.utexas.edu/users/tag/cc>

in the survey article by Gelfond and Lifschitz (1998). (Giunchiglia and Lifschitz, 1999) discusses relationships to situation calculus.  $\mathcal{C}+$  extends  $\mathcal{C}$  by allowing multi-valued fluents as well as Boolean fluents (Giunchiglia et al., 2001) and generalises the form of rules in the language in various ways. The definitive presentation of  $\mathcal{C}+$ , and of the underlying formalism of ‘nonmonotonic causal theories’, is provided in (Giunchiglia et al., 2004).  $\mathcal{C}+$  and the Causal Calculator CCALC have been applied successfully to a number of non-trivial benchmark examples in the temporal reasoning literature (see e.g. (Akman et al., 2004) and the CCALC website). We have used it in our own work to construct executable specifications of agent societies (see e.g. (Artikis et al., 2003a,b)).

The language  $\mathcal{C}+$  provides a means of constructing a *transition system* with certain properties. A separate language is used for making assertions about this transition system—about what is true at various states in runs of the transition system—and for expressing queries about them. One implementation route is via the translation of a  $\mathcal{C}+$  action description into a causal theory, and thence into a set of formulas of (classical) propositional logic (its ‘literal completion’). This is the method used by CCALC: it performs the required translations and then invokes a standard propositional satisfaction solver to find models of the  $\mathcal{C}+$  action description. An alternative implementation route is provided by translations into extended logic programs (Lifschitz and Turner, 1999). We are developing another translation into logic programs with a different computational behaviour ( $\mathcal{EC}+$  in the diagram). That line of development is not covered in this document. See (Craven and Sergot, 2003).



In what follows, we give particular emphasis to the transition system semantics. We want to exploit the bridge this provides to other standard methods and tools in logic and computer science, such as temporal logics and model checking techniques.

Section 2 provides an overview of the two language extensions by way of motivation. Sections 3–7 present the language  $\mathcal{C}+$  and the formalism of nonmonotonic causal theories, and the transition system semantics of  $\mathcal{C}+$  on which the rest of the development builds. Section 8 presents the language  $(\mathcal{C}+)^+$  and the

treatment of ‘counts as’. Section 9 presents the language  $(\mathcal{C}+)^{++}$  and its treatment of permission. Section 10 summarises and gives a brief outline of further developments.

## 2 Outline

We want to be able to represent and reason about procedures and protocols—what actions are possible, permitted, valid at each stage; what rights, powers, obligations are created and destroyed as the protocol/procedure progresses; what kind of normative or institutional relations are created when the protocol/procedure reaches completion. We want to be able to determine from a recorded narrative of actions that have been performed what actions are now possible, permitted, valid, and what their effects will be; we should also like to be able to prove general properties of these protocols (termination, ‘fairness’, that such and such a protocol has a particular property, that two protocols have identical effects, and so on). Typical examples include auction protocols, rules of procedure, argumentation/dialogue games, negotiation protocols, contract formation, business rules and procedures, procedural law.

We have taken (labelled) transition systems as the basic representational structure. So we will have fluents (state variables, Boolean and non-Boolean, whose values vary from state to state), transition labels to represent actions, and in  $\mathcal{C}+$ , a way of specifying the effects of actions on fluents.  $\mathcal{C}+$  expressions of the form  $\alpha$  causes  $F$  if  $G$  express that  $F$  holds after any transition of type  $\alpha$  when  $G$  holds in the state immediately preceding the transition. Other language features of  $\mathcal{C}+$  will be presented later.  $\mathcal{C}+$  provides a treatment of default persistence (‘inertia’) of fluents, indirect effects of actions, and a form of concurrent execution of actions. It also allows us to express properties of transitions as a whole, such as saying that a transition is ‘permitted’ or ‘acceptable’ or ‘unusual’. In future developments we plan to introduce explicit representations of agents (including names for collective agents as well as for individuals), roles, and multiple institutions. These are not included yet.

The first additional feature of the extended language  $(\mathcal{C}+)^{++}$  is a relationship (‘counts as’) between actions: we want to be able to say that, in circumstances  $G$ , a transition of type  $\alpha$  counts as a transition of type  $\beta$ . The effects of such a transition are those that it has by virtue of being a transition of type  $\alpha$  and those it has by virtue of being a transition of type  $\beta$ . This is essentially the same idea as the ‘counts as’ conditional presented in (Jones and Sergot, 1996) but now transposed to a different treatment of action on a quite different semantical structure.

For example: suppose we wish to represent that birth of a person  $X$  initiates (or ‘causes’) that  $X$  is a citizen of the Republic (it does not matter which Republic—let it be the Republic of Great Britain) if either of  $X$ ’s parents is a citizen of the Republic at the time of  $X$ ’s birth.  $X$ ’s birth, moreover, makes the parents of  $X$  happy. Let  $birth(X, F, M)$  and  $acquires\_cit(X)$  be names of actions, and  $citizen(X)$ ,  $parent(X, Y)$ , and  $happy(X)$  be fluents, for  $X, Y, F$  and  $M$  ranging over some given set of person names. The expression  $birth(X, F, M)$  is to be read as representing the birth of a child  $X$  to a father  $F$  and a mother  $M$ . In

the language  $\mathcal{C}+$  we could write:

*birth* ( $X, F, M$ ) causes (*parent* ( $X, F$ )  $\wedge$  *parent* ( $X, M$ ))  
*birth* ( $X, F, M$ ) causes (*happy* ( $F$ )  $\wedge$  *happy* ( $M$ ))  
*birth* ( $X, F, M$ ) causes *citizen* ( $X$ ) if (*citizen* ( $F$ )  $\vee$  *citizen* ( $M$ ))

In the extended language  $(\mathcal{C}+)^{++}$  the same example could be represented instead as follows:

*birth* ( $X, F, M$ ) causes (*parent* ( $X, F$ )  $\wedge$  *parent* ( $X, M$ ))  
*birth* ( $X, F, M$ ) causes (*happy* ( $F$ )  $\wedge$  *happy* ( $M$ ))  
*birth* ( $X, F, M$ ) counts\_as *acquires\_cit* ( $X$ ) if (*citizen* ( $F$ )  $\vee$  *citizen* ( $M$ ))  
*acquires\_cit* ( $X$ ) causes *citizen* ( $X$ )

This  $(\mathcal{C}+)^{++}$  representation separates what some (e.g. (Searle, 1969)) have called a ‘brute fact’, the birth of  $X$ , from its institutional effects, the acquisition of citizenship by  $X$ , as recognised by the Republic. It separates the effects that the birth of  $X$  has by virtue of being a birth (the parents of  $X$  become happy, a brute fact), and possible effects it might have by virtue of being an acquisition of citizenship ( $X$  becomes a citizen, an institutional fact). This separation is essential, we believe, to give structure to the representation, though the benefits of that are difficult to demonstrate with such a small example. Larger and more realistic examples will be presented in due course. Furthermore, it allows us to specify with much greater precision what is permitted and what is not. For suppose that the Republic forbids the birth of  $X$  when the parents of  $X$  already have another child. In the language  $(\mathcal{C}+)^{++}$ , this can be expressed

not-permitted *birth* ( $X, F, M$ ) if *parent* ( $Y, F$ )  $\wedge$  *parent* ( $Y, M$ )  $\wedge$   $X \neq Y$

But does this imply that the acquisition of citizenship of  $X$  is also forbidden (not permitted) in these circumstances? No—there may be other means, such as naturalization, by which a person  $X$  can acquire citizenship, and those other means are not necessarily forbidden.

On the other hand, suppose that the children of Adam and Eve are not permitted to be citizens (for some reason or another). In the language  $(\mathcal{C}+)^{++}$ , we can write

not-permitted *citizen* ( $X$ )  $\wedge$  *parent* ( $X, Adam$ )  $\wedge$  *parent* ( $X, Eve$ )

In  $(\mathcal{C}+)^{++}$  it can be inferred from this (see Section 9), not only that *acquires\_cit* ( $X$ ) is not permitted for any child  $X$  of Adam and Eve, but also that *birth* ( $X, Adam, Eve$ ) is not permitted whenever *citizen* ( $Adam$ ) or *citizen* ( $Eve$ ) holds (because then *citizen* ( $X$ ) holds, and *citizen* ( $X$ ) is not permitted for this  $X$ ).

For the representation of protocols and procedures (such as the rules of an auction, or the formation of a contract) the ‘counts as’ relation provides a means—again, essential we believe—for expressing the distinction between unsuccessful attempts to perform a certain, institutionally significant, kind of action, and successful performance of that action. For example, in an offer-acceptance protocol as usually found in contract formation, we often say that  $X$  ‘made an offer’ to  $Y$  that such-and-such should be done even when the offer that  $X$  made was

in fact an invalid offer that has no possible effect. For example, when Jeremy aged 10 offers to sell his father’s car JYU-927W to Alex we still say Jeremy ‘made an offer’ to Alex even though clearly Jeremy, as a 10 year old child, has no power to offer to sell the car to anyone. Jeremy ‘makes an offer’ to Alex only in the sense that he performs some kind of communicative act that would have been regarded as a real, valid offer had Jeremy been empowered to offer to sell the car at the time. Since Jeremy was not so empowered, he performed the communicative act (he ‘made an offer’) but this act was not effective (he did not ‘make an offer’ in the other sense of that term).

In line with the account given in (Jones and Sergot, 1996), we will distinguish between the action that attempts, possibly unsuccessfully, to produce some institutionally significant fact (such as offering to sell the car) and the action that (by definition) successfully does produce that institutionally significant fact. The notation  $X:\alpha$  will represent that agent  $X$  performed (successfully) the action  $\alpha$ . The notation  $X \text{ signals } \alpha$  will represent that agent  $X$  made an attempt, possibly successful, possibly not, to perform the action  $\alpha$ . (The keyword ‘attempts’ or ‘indicates’ could have been chosen in place of ‘signals’. It is difficult to find a suitable term that works in all contexts.) How does  $X$  ‘signal’ or ‘indicate’ his attempt to offer? There are further rules of the ‘counts as’ type to specify that.

So we will have rules specifying when an attempted act is successful, i.e., rules of the form:

$$(X \text{ signals } \alpha) \text{ counts\_as } (X:\alpha) \text{ if } \dots \text{ conditions}$$

For convenience, we shall usually specify such rules using the abbreviation **pow**:

$$\text{pow}(X, \alpha) =_{\text{def}} (X \text{ signals } \alpha) \text{ counts\_as } (X:\alpha)$$

Example:

$$\text{pow}(X, \text{sell}(Y)) \text{ if } \text{owner}(X, Y)$$

(The example is for illustration. It does not pretend to be a representation of the concept of selling or of ownership.)

Here **pow** represents what Jones and Sergot (1996) termed ‘institutionalised power’, a generalisation to organisational structures in general (‘institutions’) of the concept that in legal theory is variously called ‘legal power’, ‘legal competence’, or ‘legal capacity’. Thus, in the legal context,  $\text{pow}(X, \alpha)$  represents that  $X$  has the legal competence/capacity to perform action  $\alpha$ .

*Permission* to perform certain kinds of acts is specified separately. For example, Jeremy’s father might want to tell his son that he (Jeremy) is not permitted to offer to sell car JYU-927W to anyone. He means by this:

$$\text{not-permitted } \text{Jeremy signals offer}(\text{Jeremy}, Y, \text{JYU-927W})$$

The fact that Jeremy, as a 10 year old child, is not empowered (not legally competent) to make an offer to sell the car is irrelevant—it is *attempts* by Jeremy to offer to sell the car that are not permitted.

Or to take another example, we might wish to say that in general, in a certain institution, acts of communicating offers without the requisite power/competence are regarded as anti-social and are not permitted:

$$\text{not-permitted } X \text{ signals offer}(X, Y, P) \text{ if } \neg \text{pow}(X, \text{offer}(X, Y, P))$$

In practice, we would probably want to say that making offers which are *prima facie* invalid is not permitted. That refinement is easily made.

As an *implementation* issue, we might wish to go on to specify that a system implementing a certain protocol/procedure should disable non-permitted actions. That is a separate consideration. The point is that we distinguish and specify separately:

- when attempts to perform actions ‘count as’ successful—these are the constitutive rules of the protocol;
- which kinds of actions (successful, unsuccessful) are permitted—these are normative-regulative rules that further classify runs of the protocol into permitted (acceptable, desirable) and not permitted (unacceptable, undesirable);
- how a system to implement/control/enforce the protocol should deal with prohibited (unacceptable, undesirable) runs of the protocols.

Examples of specifications of protocols and procedures from the area of multi-agent systems in this style, though not using the language  $(\mathcal{C}+)^{++}$ , may be found in (Artikis et al., 2002, 2003a,b). The aim of this report is to present the language  $(\mathcal{C}+)^{++}$  and its features.

### 3 Preliminaries

The earliest presentations of  $\mathcal{C}+$  gave the semantics of the language both in terms of transition systems and as translations to ‘causal theories’. The most recent presentation (Giunchiglia et al., 2004) downplays the transition systems aspect, presenting the language as a higher-level notation for causal theories of a certain kind with the relationship to transition systems presented as a secondary issue. In our case, however, one of the main aims is to explore the suitability of (labelled) transition systems as the semantic basis for a language to represent actions and norms, and for this reason we want to emphasise the transition systems aspects. The formal semantics of  $\mathcal{C}+$  and its extensions  $(\mathcal{C}+)^+$  and  $(\mathcal{C}+)^{++}$  will be given here in terms of transition systems, with translations to causal theories treated primarily as an implementation route.

The material in this section is quite standard, though the transition systems defined by  $\mathcal{C}+$  do have a novel feature in that they allow formulas to be evaluated on transition labels as well as states. The main purpose of the section is to introduce terminology and some notation.

#### 3.1 Multi-valued signatures

A *multi-valued propositional signature*  $\sigma$  (Giunchiglia et al., 2001, 2004) is a set of symbols called *constants*. For each constant  $c$  in  $\sigma$ , there is a non-empty set  $dom(c)$  of values called the *domain* of  $c$ . An *atom* of a signature  $\sigma$  is an expression of the form  $c=v$  where  $c$  is a constant in  $\sigma$  and  $v \in dom(c)$ . We will write  $Atoms(\sigma)$  for the set of atoms of signature  $\sigma$ . An atom  $c=v$  is *trivial* if the domain of  $c$  is a singleton, i.e., if  $dom(c) = \{v\}$  for some value  $v$ . We write  $Trivial(\sigma)$  for the set (possibly empty) of the trivial atoms of  $\sigma$ .

A *formula*  $\varphi$  of signature  $\sigma$  is any propositional compound of atoms of  $\sigma$ .

$$\varphi ::= \perp \mid \top \mid \text{any atom } c = v \mid \neg\varphi \mid \varphi \wedge \varphi \mid \varphi \vee \varphi.$$

We write  $\varphi \rightarrow \psi$  for truth functional ('material') implication ( $\varphi \rightarrow \psi =_{\text{def}} \neg\varphi \vee \psi$ ).  $c \neq v$  is shorthand for  $\neg(c = v)$ .  $\top$  (true) and  $\perp$  (false) are 0-ary connectives (rather than constants with a fixed interpretation).  $\text{Form}(\sigma)$  denotes the set of formulas of signature  $\sigma$ .

A *Boolean constant* is one whose domain is the set of truth values  $\{\mathbf{t}, \mathbf{f}\}$ . If  $p$  is a Boolean constant,  $p$  is shorthand for the atom  $p = \mathbf{t}$  and  $\neg p$  for the atom  $p = \mathbf{f}$ . Notice that, as defined here,  $\neg p$  is an *atom* when  $p$  is a Boolean constant.

An *interpretation* of  $\sigma$  is a function that maps every constant in  $\sigma$  to an element of its domain. An interpretation  $I$  *satisfies* an atom  $c = v$ , written  $I \models_{\sigma} c = v$ , if  $I(c) = v$ . We will often omit the subscript  $\sigma$  when context allows. The satisfaction relation  $\models_{\sigma}$  is extended from atoms to formulas in accordance with the standard truth tables for the propositional connectives. When  $X$  is a set of formulas we also write  $I \models_{\sigma} X$  to signify that  $I \models_{\sigma} \varphi$  for all formulas  $\varphi \in X$ .  $I$  is then a *model* for the set of formulas  $X$ . A set of formulas  $X$  is  $\sigma$ -*satisfiable* if it has a model, i.e., if  $I \models_{\sigma} X$  for some interpretation  $I$  of  $\sigma$ . Where  $\sigma$  is a signature, we will write  $\text{I}(\sigma)$  for the set of interpretations of  $\sigma$ .

An expression of the form  $c = d$ , where both  $c$  and  $d$  are constants, is understood as shorthand for the disjunction  $(c = v_1 \wedge d = v_1) \vee \dots \vee (c = v_k \wedge d = v_k)$  where  $v_1, \dots, v_k$  are the elements of  $\text{dom}(c) \cap \text{dom}(d)$ .  $c \neq d$  is shorthand for  $\neg(c = d)$ .

We write  $\models_{\sigma} \varphi$  to mean that  $I \models_{\sigma} \varphi$  for all interpretations  $I$  of  $\sigma$ . Where  $X$  is a set of formulas of signature  $\sigma$ ,  $X \models_{\sigma} \varphi$  denotes that  $I \models_{\sigma} \varphi$  for all interpretations  $I$  of  $\sigma$  such that  $I \models_{\sigma} X$ . When  $X'$  is a set of formulas of signature  $\sigma$ ,  $X \models_{\sigma} X'$  is shorthand for  $X \models_{\sigma} \varphi$  for all formulas  $\varphi \in X'$ . It is easy to check that (for finite  $X$  and  $X'$ )  $X \models_{\sigma} \varphi$  iff  $\models_{\sigma} \bigwedge(X) \rightarrow \varphi$ , and  $X \models_{\sigma} X'$  iff  $\models_{\sigma} \bigwedge(X) \rightarrow \bigwedge(X')$ .  $X \equiv_{\sigma} X'$  means that  $X \models_{\sigma} X'$  and  $X' \models_{\sigma} X$ .

**Proposition 1** *Let  $X$  and  $X'$  be sets of atoms of signature  $\sigma$ . Let  $\text{Trivial}(\sigma)$  denote the set (possibly empty) of trivial atoms of  $\sigma$ . If  $X$  is  $\sigma$ -satisfiable, then  $X \models_{\sigma} X'$  iff  $X' \subseteq X \cup \text{Trivial}(\sigma)$ .*

*Proof.* The notation  $X \models_{\sigma} X'$  is shorthand for  $X \models_{\sigma} \varphi$  for all  $\varphi \in X'$ , so it is sufficient to show that the result holds for all atoms of  $\sigma$ : we show that if  $X$  is  $\sigma$ -satisfiable, then  $X \models_{\sigma} c = v$  iff  $c = v \in X \cup \text{Trivial}(\sigma)$ .

Left-to-right: Suppose  $X$  is  $\sigma$ -satisfiable and  $c = v$  is an atom of  $\sigma$ . Suppose  $c = v$  is not trivial. We show that if  $c = v \notin X$  then  $X \not\models_{\sigma} c = v$ . Suppose  $c = v \notin X$ . There are two cases: either (i)  $c = v' \in X$  for some  $v \neq v'$ , or (ii) there is no  $c$  atom in  $X$ . For case (i),  $c = v' \in X$  implies  $X = X \cup \{c = v'\}$ , and since  $X$  is  $\sigma$ -satisfiable, so is  $X \cup \{c = v'\}$ . For case (ii), choose any interpretation  $I$  of  $\sigma$  which has  $I \models_{\sigma} X$  and  $I(c) = v'$  for some  $v' \neq v$ . There must be such a  $v' \in \text{dom}(c)$  because, by assumption,  $c = v$  is not trivial. So in both cases (i) and (ii), there is an interpretation  $I$  of  $\sigma$  such that  $I \models_{\sigma} X \cup \{c = v'\}$  for some  $v' \neq v$ . So now we have  $I \models_{\sigma} X$  but  $I \not\models_{\sigma} c = v$ , and so  $X \not\models_{\sigma} c = v$ .

Right-to-left: Suppose  $c = v \in X \cup \text{Trivial}(\sigma)$  and  $I \models_{\sigma} X$ . We show  $I \models_{\sigma} c = v$ . If  $c = v$  is trivial then  $I \models_{\sigma} c = v$  since clearly  $I \models_{\sigma} \text{Trivial}(\sigma)$  for any interpretation  $I$  of  $\sigma$ . If  $c = v$  is not trivial then  $c = v \in X$  and hence  $I \models_{\sigma} c = v$  because  $I \models_{\sigma} X$ .  $\square$

The  $\sigma$ -satisfiability condition is required in Proposition 1. For if  $X$  is not  $\sigma$ -satisfiable, then  $X \models_{\sigma} X'$  for any set  $X' \subseteq \text{Atoms}(\sigma)$ . And suppose  $c=v$  is a trivial atom of a signature  $\sigma$ . Then  $X \models_{\sigma} c=v$  for all  $X \subseteq \text{Atoms}(\sigma)$ , even when  $c=v \notin X$ .

### Reduction to Boolean signatures

It is easily seen that when the domain of every constant  $c$  of the signature is finite, all multi-valued constants and interpretations can be translated to Boolean constants and standard propositional interpretations (Giunchiglia et al., 2004): replace  $c$  by a family of Boolean constants, one for each value of  $\text{dom}(c)$ , together with the constraint that exactly one of these new constants must have value **t** in any interpretation.

In other words, given a signature  $\sigma$ , let  $\sigma_B$  be the signature which is identical to  $\sigma$  for all Boolean constants of  $\sigma$ , and in which there is a Boolean constant of the form  $c=v$  for every non-Boolean atom  $c=v$  of  $\sigma$ . When  $I$  is an interpretation of  $\sigma$ , let  $I_B$  be the interpretation of  $\sigma_B$  such that  $I_B(p) = I(p)$  for all Boolean constants  $p$  of  $\sigma$ , and  $I_B(c=v) = \mathbf{t}$  iff  $I(c) = v$  for all non-Boolean atoms  $c=v$  of  $\sigma$ .

An atom  $c=v$  can be viewed as a classical, propositional atom. Finding a model of a set  $X$  of formulas of a (multi-valued) signature  $\sigma$  can be reduced to finding a classical model of  $X$  together with the following set of additional formulas  $\text{Boolean}(\sigma)$ :

$$\bigvee_v (c=v) \wedge \bigwedge_{v \neq w} \neg(c=v \wedge c=w) \quad \text{for all } c \in \sigma$$

That is,  $I \models_{\sigma} \varphi$  iff  $I_B \models_{\sigma_B} \{\varphi\} \cup \text{Boolean}(\sigma)$ , and hence

$$X \models_{\sigma} \varphi \quad \text{iff} \quad X \cup \text{Boolean}(\sigma) \models_{\text{PL}} \varphi$$

where  $\models_{\text{PL}}$  denotes classical, truth-functional entailment.

## 3.2 Transition systems

A *labelled transition system* is a structure  $\langle S, \mathbf{A}, R \rangle$  in which

- $S$  is a (non-empty) set of ‘states’;
- $\mathbf{A}$  is a (non-empty) set of transition labels (also called ‘events’);
- $R$  is a set of transitions,  $R \subseteq S \times \mathbf{A} \times S$ .

It does not matter whether we think of the labelled transitions as a single three-place relation  $R$ , as here, or as a family of binary relations  $\{R_{\varepsilon}\}_{\varepsilon \in \mathbf{A}}$ . The former is chosen here for consistency with published accounts of the language  $\mathcal{C}+$ . A transition system can be depicted as a labelled directed graph. Every state  $s$  is a node of the graph. Labelled directed edges of the graph are the tuples  $(s, \varepsilon, s')$  of  $R$ .

We are free to interpret the labels on the transitions in various ways. The usual way is to see each label as corresponding to execution of an action or perhaps several actions concurrently. It is then usual to call the transition label an ‘event’. The triple  $(s, \varepsilon, s')$  represents execution of event  $\varepsilon$  in state  $s$  leading



(possibly non-deterministically) to the state  $s'$ . An event  $\varepsilon$  is *executable* in  $s$  when there is at least one tuple  $(s, \varepsilon, s')$  in  $R$ . An event  $\varepsilon$  is *deterministic* in  $s$  if there is at most one such  $s'$ .

**Paths, ‘runs’, or ‘histories’** A *run* or *trace* of a transition system is a finite or infinite ( $\omega$  length) path through the system. (One or other of the terms *run* or *trace* is often reserved to refer to infinite length paths. We will use ‘run’ and ‘path’ interchangeably and avoid the use of the term ‘trace’. The account of  $\mathcal{C}+$  in (Giunchiglia et al., 2004) uses the term ‘history’.)

Let  $\langle S, \mathbf{A}, R \rangle$  be a transition system. A *run* (or *path* or *history*) of length  $m$  is a sequence

$$s_0 \varepsilon_0 s_1 \cdots s_{m-1} \varepsilon_{m-1} s_m \quad (m \geq 0)$$

such that  $s_0, s_1, \dots, s_m \in S$ ,  $\varepsilon_0, \dots, \varepsilon_{m-1} \in \mathbf{A}$ , and  $(s_i, \varepsilon_i, s_{i+1}) \in R$  for  $0 \leq i < m$ .

Sometimes there is a distinguished set  $S_0 \subseteq S$  of *initial states*. All runs (or histories) are then defined so that their first state  $s_0 \in S_0$ . If there is a single initial state  $S_0 = \{s_0\}$  then the set of all runs of the transition system can be seen as a *tree* rooted in  $s_0$ .

### 3.3 Query languages

A wide variety of languages—we will call them query languages—can be interpreted on labelled transition systems. These include simple propositional languages, as well as temporal logics such as CTL and LTL widely used for expressing and verifying properties of transition systems.

We begin with states. Let  $\sigma^f$  be a multi-valued signature of constants called ‘state variables’, or more usually in AI terminology, *fluent constants*. Given a labelled transition system  $\langle S, \mathbf{A}, R \rangle$  we add a valuation function which specifies, for every fluent constant  $f \in \sigma^f$  and every state  $s \in S$ , a value in  $\text{dom}(f)$ . Throughout this paper we shall be dealing with the *special case* of transition systems in which

- each state  $s \in S$  is an interpretation of  $\sigma^f$ ,  $S \subseteq \text{I}(\sigma^f)$ .

The valuation function for fluent constants is then redundant: the state already specifies the value of each fluent constant in that state. It is often convenient to adopt the convention that an interpretation  $I$  of  $\sigma^f$  is represented by the set of atoms of  $\sigma^f$  that are satisfied by  $I$ . A state is then a (complete, and consistent) set of fluent atoms. We sometimes say a formula  $\varphi$  ‘holds in’ state  $s$  or ‘is true in’ state  $s$  as alternative ways of saying that  $s$  satisfies  $\varphi$ .

Although it is much less common, an idea employed in  $\mathcal{C}+$  and in the Causal Calculator CCALC is that another category of constants and formulas — *action formulas* — can be interpreted on the transition labels/events of a transition system. So, let  $\sigma^a$  be a multi-valued signature of constants called *action constants*, disjoint from  $\sigma^f$ . Given a labelled transition system  $\langle S, \mathbf{A}, R \rangle$  we add a valuation function for action constants which specifies, for every action constant  $a \in \sigma^a$  and every label/event  $\varepsilon \in \mathbf{A}$ , a value in  $\text{dom}(a)$ . Again, throughout the paper we shall be dealing with a special case, the case of labelled transition systems in which the set  $\mathbf{A}$  of labels/events is the set of interpretations of

$\sigma^a$ . In other words the transition systems of interest will be those of the form  $\langle \sigma^f, S, I(\sigma^a), R \rangle$ , on which we will interpret various query languages of signature  $\sigma^f \cup \sigma^a$ , or variations thereof.  $(\sigma^f, \sigma^a)$  is the ‘action signature’ of the transition system.

Note that since a transition label/event  $\varepsilon$  is an interpretation of  $\sigma^a$ , it is meaningful to say that  $\varepsilon$  satisfies an action formula  $\alpha$  ( $\varepsilon \models_{\sigma^a} \alpha$ ). When  $\varepsilon \models_{\sigma^a} \alpha$  we say that the event  $\varepsilon$  is of type  $\alpha$ . When  $\varepsilon \models_{\sigma^a} \alpha$  we also say that the transition  $(s, \varepsilon, s')$  is a transition of type  $\alpha$ .

Moreover, since a transition label is an interpretation of the action constants  $\sigma^a$ , it can also be represented by the set of atoms that it satisfies. The suggested reading of a transition label  $\{a_1 = v_1, a_2 = v_2, \dots, a_n = v_n\}$  for an action signature with action constants  $a_1, a_2, \dots, a_n$  is that it represents a composite action in which the elementary actions  $a_1 = v_1, a_2 = v_2, \dots, a_n = v_n$  are performed (or occur) concurrently. Where  $a$  is a Boolean action constant,  $\neg a$ , i.e.  $a = f$ , can be read as indicating that action  $a$  is not performed; and where all action constants are Boolean, the action  $\{a_1 = f, \dots, a_n = f\}$  can be read as representing the ‘null’ event.

For example: suppose there are three agents,  $a$ ,  $b$ , and  $c$  which can move in direction  $E$ ,  $W$ ,  $N$ , or  $S$ , or remain idle. Suppose (for the sake of an example) that they can also whistle as they move (they are trains, let us say). Let the action signature consist of action constants  $move(a)$ ,  $move(b)$ ,  $move(c)$  with domains  $\{E, W, N, S, idle\}$ , and Boolean action constants  $whistle(a)$ ,  $whistle(b)$ ,  $whistle(c)$ . Then one possible interpretation of the action signature, and therefore one possible transition label, is

$$\{move(a) = E, move(b) = N, move(c) = idle, whistle(a), \neg whistle(b), whistle(c)\}$$

Because of the way that action formulas are evaluated on a transition  $(s, \varepsilon, s')$ , an action formula can also be regarded as expressing a property of the transition  $(s, \varepsilon, s')$  as a whole. For example, we might include in the action signature a Boolean action symbol *unusual* which is intended to be satisfied by transition labels/events which are regarded as unusual for some reason or another. In later sections we will employ the action constant **trans** (possible values: **green** and **red**) to indicate whether a given transition is ‘permitted/acceptable’ (**green**) or ‘not permitted/unacceptable’ (**red**).

**Example (queries on states)** For an example of one possible query language, consider the propositional language of signature  $\sigma^f$  extended with expressions of the form  $[\alpha]\varphi$  where  $\alpha$  is any formula of signature  $\sigma^a$ .  $[\alpha]\varphi$  is intended to express that  $\varphi$  holds in every state following a transition of type  $\alpha$ . Let  $\mathcal{T} = \langle \sigma^f, S, I(\sigma^a), R \rangle$ . Define the satisfaction relation  $\models$  so that  $\mathcal{T}, s \models f = v$  iff  $s \models_{\sigma^f} f = v$  for fluent atoms  $f = v$ ,  $\mathcal{T}, s \models \neg\varphi$  iff  $\mathcal{T}, s \not\models \varphi$ ,  $\mathcal{T}, s \models (\varphi \wedge \varphi')$  iff  $\mathcal{T}, s \models \varphi$  and  $\mathcal{T}, s \models \varphi'$ , and so on as usual for the other truth functional connectives, and

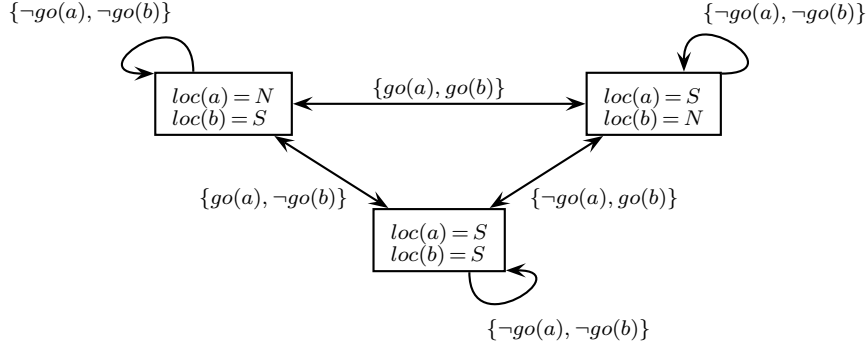
$$\mathcal{T}, s \models [\alpha]\varphi \text{ iff } \mathcal{T}, s' \models \varphi \text{ for every } s' \text{ such that } (s, \varepsilon, s') \in R \text{ and } \varepsilon \models_{\sigma^a} \alpha.$$

The notation  $[\alpha]\varphi$  is intended to be reminiscent of that of dynamic logic, except that in this language  $\alpha$  is an action *formula* not a transition *label* as in dynamic

logic. In particular, ‘action negation’, which is problematic in dynamic logic, has a straightforward meaning here:  $[\neg\alpha]\varphi$  simply says that  $\varphi$  holds in every state following a transition of type  $\neg\alpha$ .

We say (as usual) that a formula  $\varphi$  is *true* in the transition system  $\mathcal{T} = \langle \sigma^f, S, I(\sigma^a), R \rangle$ , written  $\mathcal{T} \models \varphi$ , when  $\mathcal{T}, s \models \varphi$  for every state  $s$  in  $\mathcal{T}$ .

**Example** Let  $\sigma^f$  be the set of fluent constants  $\{loc(a), loc(b)\}$  with possible values  $\{N, S\}$ , and let  $\sigma^a$  be the set of Boolean action constants  $\{go(a), go(b)\}$ . Consider the transition system  $\mathcal{T}$  depicted in the following diagram:



There is no state  $\{loc(a) = N, loc(b) = N\}$  in  $\mathcal{T}$  (for the sake of the example).

One can check by observation that, for example:

$$\begin{aligned} \mathcal{T} &\models loc(a) = N \rightarrow [go(a)] loc(a) = S \\ \mathcal{T} &\models loc(a) = N \rightarrow [go(a)] [go(a)] loc(a) = N \\ \mathcal{T} &\models loc(a) = N \rightarrow [\neg go(a)] loc(a) = N \\ \mathcal{T} &\models loc(a) = N \rightarrow [go(b)] loc(a) = S \\ \mathcal{T} &\models (loc(a) = S \wedge loc(b) = S) \rightarrow [go(a) \wedge go(b)] \perp \end{aligned}$$

**Example (time-stamped queries)** Query languages can also be interpreted on the paths (‘runs’) of a transition system. One candidate is the query language supported by the Causal Calculator CCALC. This uses propositional formulas of time-stamped fluent and action constants: the time-stamped fluent atom  $f[i] = v$  represents that fluent atom  $f = v$  holds at integer time  $i$ , or more precisely, that  $f = v$  is satisfied by the state  $s_i$  of a path  $s_0 \varepsilon_0 \cdots \varepsilon_{i-1} s_i, \cdots$  of the transition system; the time-stamped atom  $a[i] = v$  represents that action atom  $a = v$  is satisfied by the transition label  $\varepsilon_i$  of a path  $s_0 \varepsilon_0 \cdots s_i \varepsilon_i s_{i+1}, \cdots$ .

In general, given a multi-valued signature  $\sigma$  and a non-negative integer  $i$ , we write  $\sigma[i]$  for the signature consisting of all constants of the form  $c[i]$  where  $c$  is a constant of  $\sigma$ , with  $dom(c[i]) = dom(c)$ . For any non-negative integer  $m$ , we write  $\sigma_m$  for the signature  $\sigma[0] \cup \cdots \cup \sigma[m]$ . As usual,  $Form(\sigma_m)$  is the set of formulas of signature  $\sigma_m$ .

The time-stamped query language used (as in CCALC) to express properties of paths of length  $m$  of a transition system with action signature  $(\sigma^f, \sigma^a)$  is

the propositional language of signature  $\sigma_m^f \cup \sigma_{m-1}^a$ . In other words, the query language is  $Form(\sigma_m^f \cup \sigma_{m-1}^a)$ , whose formulas consist of:

- atoms  $f[i] = v$  where  $i \in 0..m$  and  $f = v$  is a fluent atom of  $\sigma^f$ ;
- atoms  $a[i] = v$  where  $i \in 0..m-1$  and  $a = v$  is an action atom of  $\sigma^a$ ;
- all truth-functional compounds of the above.

Let  $\pi = s_0 \varepsilon_0 s_1 \cdots s_{m-1} \varepsilon_{m-1} s_m$  be a path of length  $m$  of a transition system  $\mathcal{T}$  of action signature  $(\sigma^f, \sigma^a)$ . An atom  $f[i] = v$  for any fluent constant  $f$  of  $\sigma^f$  and  $0 \leq i \leq m$  is *true on* path  $\pi$  (or ‘holds on’ path  $\pi$ , or ‘is satisfied by’ path  $\pi$ ), written  $\mathcal{T}, \pi \models_m f[i] = v$ , when  $s_i \models_{\sigma^f} f = v$ ; for action constants  $a$  of  $\sigma^a$  and  $0 \leq i < m$ ,  $\mathcal{T}, \pi \models_m a[i] = v$  when  $\varepsilon_i \models_{\sigma^a} a = v$ ; and  $\models_m$  is extended to formulas  $\varphi$  of signature  $\sigma_m = \sigma_m^f \cup \sigma_{m-1}^a$  by the usual truth tables for the propositional connectives.

We will say  $\varphi$  is true on paths of length  $m$  of  $\mathcal{T}$ , written  $\mathcal{T} \models_m \varphi$  when  $\mathcal{T}, \pi \models_m \varphi$  for all paths of length  $m$  of  $\mathcal{T}$ .

The definition of the satisfaction relation  $\models_m$  can be stated more concisely, as follows. First some notation (which will also be useful later).

Let  $\pi[i]$  denote the  $i$ th component of a path  $\pi$ : that is, when  $\pi = s_0 \varepsilon_0 s_1 \cdots s_i \varepsilon_i s_{i+1} \cdots$ , let  $\pi[i] = s_i \cup \varepsilon_i$ . Clearly,  $\pi[i]$  is an interpretation of  $\sigma^f \cup \sigma^a$  when  $\pi$  is a path of a transition system with action signature  $(\sigma^f, \sigma^a)$ .

For any formula  $\psi$  of signature  $\sigma^f \cup \sigma^a$ , let  $\psi[i]$  stand for the formula of signature  $\sigma^f[i] \cup \sigma^a[i]$  obtained by time-stamping every constant in  $\psi$  with  $i$ , that is, replacing every constant  $c$  in  $\psi$  by the constant  $c[i]$ . Clearly, every formula  $\varphi$  of signature  $\sigma_m = \sigma_m^f \cup \sigma_{m-1}^a$  is a truth-functional compound of formulas of the form  $\psi[i]$  where  $0 \leq i \leq m$  and  $\psi$  is a formula of signature  $\sigma = \sigma^f \cup \sigma^a$ .

Now, for any path  $\pi$  of length  $m$  of a transition system  $\mathcal{T}$  of action signature  $(\sigma^f, \sigma^a)$ , we have

$$\mathcal{T}, \pi \models_m \psi[i] \text{ iff } \pi[i] \models_{\sigma} \psi$$

**Example (contd)** Consider again the transition system  $\mathcal{T}$  depicted in the figure above. We have, amongst other things:

$$\mathcal{T} \models_1 (loc(a)[0] = N \wedge go(a)[0]) \rightarrow loc(a)[1] = S$$

$$\mathcal{T} \models_2 (loc(a)[0] = N \wedge go(a)[0] \wedge go(a)[1]) \rightarrow loc(a)[2] = N$$

$$\mathcal{T} \models_1 (loc(a)[0] = N \wedge \neg go(a)[0]) \rightarrow loc(a)[1] = N$$

$$\mathcal{T} \models_2 (loc(a)[0] = N \wedge go(b)[0] \wedge go(a)[1]) \rightarrow loc(a)[2] = N$$

$$\mathcal{T} \models_m (loc(a)[i] = N \wedge loc(a)[i+2] = N) \rightarrow (go(a)[i] \wedge go(a)[i+1]) \text{ for all } 0 \leq i \leq m-2$$

$$\mathcal{T} \models_m (loc(a)[i] = S \wedge loc(b)[i] = S) \rightarrow (\neg go(a)[i] \wedge \neg go(b)[i]) \text{ for all } 0 \leq i \leq m-1$$

A wide variety of other languages can be interpreted on transition systems. In Section 9.6 we will show how the temporal logic CTL can be used to express properties of a transition system defined by an action description of  $(\mathcal{C}^+)^{++}$ .

## 4 The Action Description Language $\mathcal{C}+$

The language  $\mathcal{C}+$  has evolved through several versions. Here we follow the (definitive) presentation in (Giunchiglia et al., 2004) though with more emphasis on the transition system semantics and with some small notational and terminological differences. These differences are recorded in the text.

### 4.1 Syntax

An *action signature*  $(\sigma^f, \sigma^a)$  is a (non-empty) set  $\sigma^f$  of *fluent constants* and a (non-empty) set  $\sigma^a$  of *action constants*. Fluent constants are partitioned into *simple* fluent constants and *statically determined* fluent constants. *Simple fluent constants* are related to actions by *dynamic laws*, which specify how the values of simple fluents change in transitions from one state to another. The values of *statically determined fluents* can also vary from state to state but they are defined by *static laws* relating their values in a state to the values of other fluents in that state. Static laws can also express constraints on the values of simple fluent constants in a state.

Action constants are partitioned into *exogenous* action constants, which are used to name kinds of actions, and *transition attributes* which are used to represent attributes of actions and other properties of transitions. (This partitioning of action constants is not part of standard presentations of  $\mathcal{C}+$ . There, the declaration of exogenous actions is done by means of a special kind of  $\mathcal{C}+$  law and not, as here, as part of the signature of the language. There is no practical difference between these two treatments, as will be explained later. We prefer to specify exogenous actions in the signature for uniformity with the treatment of fluent constants, but nothing of significance turns on this.)

A *fluent formula* is any truth-functional compound of fluent atoms (i.e., a formula of signature  $\sigma^f$ ). An *action formula* is any formula of signature  $\sigma^a$  that contains at least one action constant. Notice that since  $\top$  and  $\perp$  are treated as 0-ary connectives, they are not action formulas, though  $a \wedge \neg a$  and  $a \vee \neg a$  are for any Boolean action constant  $a$ . The language also allows formulas of signature  $\sigma^f \cup \sigma^a$ .

The sentences of  $\mathcal{C}+$  are of three kinds.

- A *static law* (or rule) is an expression of the form

$$\text{caused } F \text{ if } G \tag{1}$$

where  $F$  and  $G$  are fluent formulas (i.e., formulas of signature  $\sigma^f$ ).  $F$  is the *head* of the rule and  $G$  is the *body*. Static laws are used to express constraints that hold in all states and to define statically determined fluents.

- A *fluent dynamic law* (or rule) is an expression of the form

$$\text{caused } F \text{ if } G \text{ after } \psi \tag{2}$$

where  $F$  and  $G$  are fluent formulas, and  $\psi$  is a formula of signature  $\sigma^f \cup \sigma^a$ , with the restriction that statically determined fluent constants may not appear in the head of the rule,  $F$ . We will say  $G$  *after*  $\psi$  is the body of the rule.

Informally, in a transition  $(s, \varepsilon, s')$ , formulas  $F$  and  $G$  are evaluated at  $s'$  (the resulting state), fluent atoms in  $\psi$  are evaluated at  $s$  (i.e., in the state immediately before the transition), and action atoms in  $\psi$  are evaluated on the transition label  $\varepsilon$  itself. It may be helpful to note that any set of fluent dynamic laws can be written equivalently as a set of fluent dynamic laws of the form

$$\text{caused } F \text{ if } G \text{ after } H \wedge \alpha$$

where  $H$  is a fluent formula and  $\alpha$  is an action formula, and this is the form that will appear most frequently in subsequent sections. Fluent dynamic laws are primarily used to express how the values of fluents are affected by different kinds of actions, and to specify which of the fluents are ‘inertial’, that is, which fluents are such that their values persist by default from one state to the next. The expression

$$\alpha \text{ causes } F \text{ if } H$$

is shorthand for a fluent dynamic law of the form:  $\text{caused } F \text{ if } \top \text{ after } H \wedge \alpha$ .

- An *action dynamic law* (or rule) is an expression of the form

$$\text{caused } \alpha \text{ if } \psi \tag{3}$$

where  $\alpha$  is an action formula (i.e., a formula of signature  $\sigma^a$ ) and  $\psi$  is any formula of signature  $\sigma^f \cup \sigma^a$ .  $\alpha$  is the head of the rule and  $\psi$  is the body.

Action dynamic laws are used to express, among other things, that any transition of type  $\alpha$  must also be of type  $\alpha'$  ( $\text{caused } \alpha' \text{ if } \alpha$ ), or that any transition from a state satisfying fluent formula  $G$  must be of type  $\beta$  ( $\text{caused } \beta \text{ if } G$ ). Examples will be provided in later sections.

An *action description* is a set of static and dynamic laws.

In the remainder of this presentation we will omit the keyword `caused` when writing causal laws. This is simply to shorten expressions and avoid linebreaks.

**Example** The effects of toggling a switch between on and off can be represented by a simple Boolean fluent constant *on* and an (exogenous) Boolean action constant *toggle* and the following pair of laws:

$$\begin{aligned} \text{toggle causes } on & \text{ if } \neg on \\ \text{toggle causes } \neg on & \text{ if } on \end{aligned}$$

which, recall, are shorthand for the following fluent dynamic laws:

$$\begin{aligned} on & \text{ if } \top \text{ after } toggle \wedge \neg on \\ \neg on & \text{ if } \top \text{ after } toggle \wedge on \end{aligned}$$

Note that default persistence (‘inertia’) of fluent values is not a built-in feature of the  $\mathcal{C}+$  language. One specifies explicitly which fluents are ‘inertial’ by means of a  $\mathcal{C}+$  law of the form

$$\text{inertial } f$$

This is shorthand for the set of fluent dynamic laws of the form

$$f = v \text{ if } f = v \text{ after } f = v, \quad \text{for every } v \in \text{dom}(f)$$

How this law works to express default persistence of  $f = v$  will become clearer when we look at the semantics of  $\mathcal{C}+$  laws. Notice that a statically determined fluent constant cannot be inertial, since inertial declarations are fluent dynamic laws, and no statically determined fluent can be in the head of a fluent dynamic law.

Early versions of the language also included a separate category of *rigid constants* in the action signature, used to represent fluents whose values are constant and do not vary from state to state. In the latest presentations (Giunchiglia et al., 2004), rigid constants are not features of an action signature but of an action description. A fluent constant  $c$  is rigid (relative to an action description  $D$ ) if the value of  $c$  is the same in every state of the transition system defined by  $D$ . Trivial fluent constants are clearly rigid, but there may also be rigid constants which are not trivial.

### Definite action descriptions

Of particular interest are *definite* action descriptions.

An action description  $D$  is *definite* when, for all static laws, fluent dynamic laws, and action dynamic laws in  $D$ :

- the head of every law is either a fluent atom (resp., action atom) or the symbol  $\perp$ , and
- no atom is the head of infinitely many laws of  $D$ .

(Note that, as defined above, a Boolean fluent constant  $p$  and its negation  $\neg p$ , and a Boolean action constant  $a$  and its negation  $\neg a$  are atoms, and hence are included in the definition.)

## 4.2 Semantics

Consider an action description  $D$  of  $\mathcal{C}+$  of signature  $(\sigma^f, \sigma^a)$ .

- A *state* is an interpretation of the fluent constants  $\sigma^f$  that is closed under the static laws of  $D$ , that is to say, that satisfies  $G \rightarrow F$  for every static law  $F$  if  $G$  in  $D$ . A state must further satisfy certain restrictions concerning statically determined fluents which we will identify presently.
- A *transition label* or *event* is an interpretation of the action constants  $\sigma^a$ . A transition label must satisfy certain constraints concerning the exogenous action constants, as identified below.
- A *transition* is a triple  $(s, \varepsilon, s')$  where  $s$  and  $s'$  are states and  $\varepsilon$  is a transition label/event. A transition *defined by* an action description  $D$  must also satisfy the fluent dynamic laws and the action dynamic laws of  $D$ , in a sense to be defined below.

The *transition system* defined by an action description  $D$  is the transition system  $\langle \sigma^f, S, I(\sigma^a), R \rangle$  which has the states of  $D$  as its states  $S$ , and which has  $(s, \varepsilon, s') \in R$  when  $(s, \varepsilon, s')$  is a transition of  $D$ .

## Semantics: Definitions

Let  $T_{static}(s)$  stand for the heads of all static laws in  $D$  whose bodies are satisfied by  $s$ ; let  $E(s, \varepsilon, s')$  stand for the heads of all fluent dynamic laws in  $D$  whose bodies are satisfied by the transition  $(s, \varepsilon, s')$ ; and let  $A(\varepsilon, s)$  stand for the heads of all action dynamic laws whose bodies are satisfied by the transition  $(s, \varepsilon, s')$ . More precisely:

$$\begin{aligned} T_{static}(s) &=_{\text{def}} \{F \mid F \text{ if } G \text{ is in } D, s \models G\} \\ E(s, \varepsilon, s') &=_{\text{def}} \{F \mid F \text{ if } G \text{ after } \psi \text{ is in } D, s' \models G, s \cup \varepsilon \models \psi\} \\ A(\varepsilon, s) &=_{\text{def}} \{\alpha \mid \alpha \text{ if } \psi \text{ is in } D, s \cup \varepsilon \models \psi\} \end{aligned}$$

(Throughout this section we omit subscripts on  $\models$  whenever context allows. Notice that since  $\sigma^f$  and  $\sigma^a$  are disjoint,  $s \models_{\sigma^f} F$  iff  $s \models_{\sigma^f \cup \sigma^a} F$  for any fluent formula  $F$ , and  $\varepsilon \models_{\sigma^a} \alpha$  iff  $\varepsilon \models_{\sigma^f \cup \sigma^a} \alpha$  for any action formula  $\alpha$ .)

In the absence of statically determined fluents, an interpretation  $s$  of  $\sigma^f$  is a state of  $D$  iff it is closed under all the static laws, i.e., iff  $s \models T_{static}(s)$ . In the presence of statically determined fluents, the characterisation of what constitutes a state is a little more complicated.

An interpretation  $s$  of  $\sigma^f$  is a *state* of  $D$  when

- $s \models T_{static}(s)$
- there is no other interpretation  $s'$  of  $\sigma^f$  which agrees with  $s$  on the interpretation of all simple fluents and is such that  $s' \models T_{static}(s)$ .

Alternatively: let  $Simple(s)$  denote the simple fluent atoms that are satisfied by an interpretation  $s$  of  $\sigma^f$ . Clearly  $s \models Simple(s)$  for any  $s$ . So then we can say the following.

**Definition 2** *Let  $D$  be an action description with fluent signature  $\sigma^f$ . An interpretation  $s$  of  $\sigma^f$  is a state of  $D$  when*

$$\{s' \in \mathbf{I}(\sigma^f) \mid s' \models T_{static}(s) \cup Simple(s)\} = \{s\}$$

*In words: when  $s \models T_{static}(s) \cup Simple(s)$  and there is no other interpretation  $s'$  of  $\sigma^f$  such that  $s' \models T_{static}(s) \cup Simple(s)$ .*

The rationale behind these definitions, especially that of the uniqueness condition, is perhaps far from obvious. It is a direct application of the methods of ‘nonmonotonic causal theories’ (see Section 5 below) in which the language  $\mathcal{C}+$  has its origins. For convenience we will write  $X \models \varphi$  to denote that  $X$  is the only interpretation that satisfies a formula (or set of formulas)  $\varphi$ . In this notation, an interpretation  $s$  of  $\sigma^f$  is a state of  $D$  when

$$s \models T_{static}(s) \cup Simple(s)$$

A transition  $(s, \varepsilon, s')$  satisfies the action dynamic laws of  $D$  when:

- $\varepsilon \models A(\varepsilon, s)$
- there is no other interpretation  $\varepsilon'$  of  $\sigma^a$  which agrees with  $\varepsilon$  on the interpretation of all exogenous action constants and is such that  $\varepsilon' \models A(\varepsilon, s)$ .



Again, let  $Exog(\varepsilon)$  denote the exogenous atoms of  $\sigma^a$  that are satisfied by an interpretation  $\varepsilon$  of  $\sigma^a$ . Clearly,  $\varepsilon \models Exog(\varepsilon)$ . The transition  $(s, \varepsilon, s')$  satisfies the action dynamic laws of  $D$  when

$$\varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon)$$

that is, when  $\{\varepsilon' \in I(\sigma^a) \mid \varepsilon' \models A(\varepsilon, s) \cup Exog(\varepsilon)\} = \{\varepsilon\}$ .

$(s, \varepsilon, s')$  is a *transition defined by* an action description  $D$  (or simply: is a transition of  $D$ ) if  $s$  is a state of  $D$ , the transition  $(s, \varepsilon, s')$  satisfies the action dynamic laws of  $D$ , and the resulting state  $s'$  is the only interpretation of  $\sigma^f$  that satisfies all formulas  $T_{static}(s') \cup E(s, \varepsilon, s')$ . (Early presentations of  $\mathcal{C}+$ , e.g. (Giunchiglia and Lifschitz, 1998; Giunchiglia et al., 2001) referred to this as a *causally explained* transition.)

**Definition 3** Let  $D$  be an action description of  $\mathcal{C}+$  of signature  $(\sigma^f, \sigma^a)$ . Let  $s$  and  $s'$  be interpretations of  $\sigma^f$ , and  $\varepsilon$  an interpretation of  $\sigma^a$ .  $(s, \varepsilon, s')$  is a transition of  $D$  iff

- $s \models T_{static}(s) \cup Simple(s)$
- $\varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon)$
- $s' \models T_{static}(s') \cup E(s, \varepsilon, s')$

In words:  $(s, \varepsilon, s')$  is a transition of  $D$  iff  $s$  is a state of  $D$ ,  $\varepsilon$  is the only model of  $A(\varepsilon, s) \cup Exog(\varepsilon)$ , and  $s'$  is the only model of  $E(s, \varepsilon, s')$  that also satisfies the static laws of  $D$ .

The uniqueness condition in the above definition does *not* imply that all events are deterministic, as will be illustrated in examples later.

The *transition system* defined by an action description  $D$  is the transition system  $\langle \sigma^f, S, I(\sigma^a), R \rangle$  which has the states of  $D$  as its states  $S$ , and which has  $(s, \varepsilon, s') \in R$  when  $(s, \varepsilon, s')$  is a transition of  $D$ . We will say that  $\pi$  is a path of  $D$  when  $\pi$  is a path of the transition system defined by  $D$ .

Notice that the last condition of Definition 3 makes no mention of  $Simple(s')$ . It is not obvious therefore that the definition is well formed, that is, that the resulting state  $s'$  of every transition  $(s, \varepsilon, s')$  defined by  $D$  is guaranteed to be a state of  $D$ . It is easy to check that this is so. It relies on the fact that formulas in  $E(s, \varepsilon, s')$  contain only simple and no statically determined fluent constants. (It might appear there is a simple argument relying on the observation that  $E(s, \varepsilon, s') \subseteq Simple(s')$ . But that is not so, because in the general case  $E(s, \varepsilon, s')$  is a set of *formulas* not a set of *atoms*.)

**Proposition 4** If  $(s, \varepsilon, s')$  is a transition defined by an action description  $D$  then  $s'$  is a state of  $D$ .

*Proof.* We show that if  $s'$  is the only model of  $T_{static}(s') \cup E(s, \varepsilon, s')$  then  $s'$  is the only model of  $T_{static}(s') \cup Simple(s')$ .

Clearly  $s'$  is a model of  $T_{static}(s') \cup Simple(s')$  because  $s' \models T_{static}(s')$  (from the assumption) and  $s' \models Simple(s')$  (by definition). So it remains to show that  $s'$  is the only such model.

Notice first that a state  $s'$  can be represented in the form  $s'_s \cup \delta'$  where  $s'_s$  is an interpretation of the simple fluent constants and  $\delta'$  is an interpretation of the statically determined fluent constants. We then have  $\text{Simple}(s'_s \cup \delta') = s'_s$ . So suppose that  $s'_s \cup \delta'$  is the only model of  $T_{\text{static}}(s'_s \cup \delta') \cup E(s, \varepsilon, s'_s \cup \delta')$ . We need to show that  $s'_s \cup \delta'$  is the only model of  $T_{\text{static}}(s'_s \cup \delta') \cup \text{Simple}(s'_s \cup \delta')$ , that is, that there is no other interpretation  $s''_s \cup \delta''$  such that  $s''_s \cup \delta'' \models T_{\text{static}}(s'_s \cup \delta') \cup \text{Simple}(s'_s \cup \delta')$ .

For any  $s''_s \cup \delta''$ ,  $s''_s \cup \delta'' \models \text{Simple}(s'_s \cup \delta')$  if and only if  $s''_s = s'_s$ , because  $\text{Simple}(s'_s \cup \delta') = s'_s$ .

So it remains to show that there is no other interpretation  $\delta'' \neq \delta'$  of the statically determined fluent constants such that  $s'_s \cup \delta'' \models T_{\text{static}}(s'_s \cup \delta')$ .

Suppose there were such a  $\delta''$ .  $s'_s \cup \delta'' \models E(s, \varepsilon, s'_s \cup \delta')$ , and  $E(s, \varepsilon, s'_s \cup \delta')$  contains no statically determined fluents, so  $s'_s \cup \delta'' \models E(s, \varepsilon, s'_s \cup \delta')$  also. Now we have  $s'_s \cup \delta'' \models T_{\text{static}}(s'_s \cup \delta') \cup E(s, \varepsilon, s'_s \cup \delta')$ , which contradicts the assumption that  $s'_s \cup \delta'$  is the only such interpretation.  $\square$

In future sections, (in particular in Section 6), we will be primarily concerned with paths/runs of the transition system defined by  $D$ . Proposition 4 can also be applied to simplify the the characterisation of paths of  $D$ , as follows.  $(s_{i-1}, \varepsilon_{i-1}, s_i)$  and  $(s_i, \varepsilon_i, s_{i+1})$  are both transitions of  $D$  when

$$\begin{aligned} s_{i-1} & \not\models T_{\text{static}}(s_{i-1}) \cup \text{Simple}(s_{i-1}), \\ \varepsilon_{i-1} & \not\models A(\varepsilon_{i-1}, s_{i-1}) \cup \text{Exog}(\varepsilon_{i-1}), \\ s_i & \not\models T_{\text{static}}(s_i) \cup E(s_{i-1}, \varepsilon_{i-1}, s_i), \end{aligned}$$

and

$$\begin{aligned} s_i & \not\models T_{\text{static}}(s_0) \cup \text{Simple}(s_i), \\ \varepsilon_i & \not\models A(\varepsilon_i, s_i) \cup \text{Exog}(\varepsilon_i), \\ s_{i+1} & \not\models T_{\text{static}}(s_{i+1}) \cup E(s_i, \varepsilon_i, s_{i+1}). \end{aligned}$$

But by Proposition 4,  $s_i \not\models T_{\text{static}}(s_i) \cup E(s_{i-1}, \varepsilon_{i-1}, s_i)$  already implies  $s_i \not\models T_{\text{static}}(s_0) \cup \text{Simple}(s_i)$  ( $s_i$  is a state of  $D$ ), and so this condition can be dropped. We thus have the following.

**Proposition 5** *Let  $D$  be a  $C+$  action description of signature  $(\sigma^f, \sigma^a)$ .*

*$s_0 \varepsilon_0 s_1 \cdots \varepsilon_{m-1} s_m$  is a path of length  $m \geq 0$  of the transition system defined by  $D$  iff  $s_0, s_1, \dots, s_m$  are interpretations of  $\sigma^f$  and  $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{m-1}$  are interpretations of  $\sigma^a$  such that*

$$\begin{aligned} s_0 & \not\models T_{\text{static}}(s_0) \cup \text{Simple}(s_0) \\ \varepsilon_0 & \not\models A(\varepsilon_0, s_0) \cup \text{Exog}(\varepsilon_0) \\ s_1 & \not\models T_{\text{static}}(s_1) \cup E(s_0, \varepsilon_0, s_1) \\ & \vdots \end{aligned}$$

$$\begin{aligned}
\varepsilon_i & \models A(\varepsilon_i, s_i) \cup Exog(\varepsilon_i) \\
s_{i+1} & \models T_{static}(s_{i+1}) \cup E(s_i, \varepsilon_i, s_{i+1}) \\
& \vdots \\
\varepsilon_{m-1} & \models A(\varepsilon_{m-1}, s_{m-1}) \cup Exog(\varepsilon_{m-1}) \\
s_m & \models T_{static}(s_m) \cup E(s_{m-1}, \varepsilon_{m-1}, s_m)
\end{aligned}$$

*Proof.* In the preceding discussion. The cases  $m = 0$  and  $m = 1$  are trivial.  $\square$

### Semantics: Definite action descriptions

When  $D$  is a definite action description, the conditions under which an interpretation  $s$  of  $\sigma^f$  is a state, and the conditions under which  $(s, \varepsilon, s')$  is a transition defined by  $D$ , can be simplified very considerably. For then  $T_{static}(s)$ ,  $E(s, \varepsilon, s')$ , and  $A(\varepsilon, s)$  are all sets of *atoms*, possibly also containing the element  $\perp$ . And if we represent an interpretation  $s$  of  $\sigma^f$  by the set of atoms of  $\sigma^f$  that are satisfied by  $s$ , the condition that  $s$  is the *only* interpretation of  $\sigma^f$  that satisfies all formulas  $T_{static}(s) \cup Simple(s)$  is equivalently stated as  $s = T_{static}(s) \cup Simple(s) \cup Trivial(\sigma^f)$ , where  $Trivial(\sigma^f)$  denotes the set (possibly empty) of the trivial atoms of  $\sigma^f$ . Trivial atoms complicate the account slightly but since they can be used to represent *rigid* fluent constants we will keep open the possibility that they are present.

Informally: suppose first that  $s = T_{static}(s) \cup Simple(s)$ . Any interpretation  $s'$  of  $\sigma^f$  satisfying  $T_{static}(s) \cup Simple(s)$  must also be such that  $T_{static}(s) \cup Simple(s) \subseteq s'$ ; every interpretation  $s'$  of  $\sigma^f$  must also satisfy  $Trivial(\sigma^f)$ . However,  $s'$  cannot contain any additional atoms of  $\sigma^f$  not in  $T_{static}(s) \cup Simple(s) \cup Trivial(\sigma^f)$  because  $s$  is already a complete set of atoms and  $s'$  would not then be an interpretation of  $\sigma^f$ . So  $s$  is unique. We can generalise the above informal argument, as follows.

**Proposition 6** *Let  $X$  be a set of atoms of signature  $\sigma$ . Let  $I$  be an interpretation of  $\sigma$ , represented as a set of atoms. Let  $Trivial(\sigma)$  denote the set (possibly empty) of trivial atoms of  $\sigma$ . Then  $I$  is the unique interpretation of  $\sigma$  satisfying  $X$  iff  $I = X \cup Trivial(\sigma)$ .*

*Proof.* Left-to-right: Suppose the set of atoms  $I$  is the only interpretation of  $\sigma$  such that  $I \models X$ . Then it follows that  $X \models_{\sigma} I$  (i.e., atoms  $I$  are all satisfied in all models of  $X$ ): there is only one interpretation of  $\sigma$  that satisfies  $X$ , that interpretation is  $I$ , and clearly  $I \models I$ .  $I \models X$  so  $X$  is  $\sigma$ -satisfiable, and by Proposition 1,  $X \models_{\sigma} I$  implies  $I \subseteq X \cup Trivial(\sigma)$ . Further, since  $I \models X$  we have  $X \subseteq I$ . And clearly  $Trivial(\sigma) \subseteq I$ . So we have  $X \cup Trivial(\sigma) \subseteq I \subseteq X \cup Trivial(\sigma)$ .

Right-to-left: Suppose  $I = X \cup Trivial(\sigma)$ . Clearly  $I$  satisfies  $X$  because  $X \subseteq I$ . To show  $I$  is unique, suppose  $I' \models X$  for some other interpretation  $I'$  of  $\sigma$ , also represented as a set of atoms. We have  $X \subseteq I'$ , and so  $I \subseteq X \subseteq I'$ , and hence  $I \subseteq I'$ . But for interpretations  $I$  and  $I'$ ,  $I \subseteq I'$  iff  $I' = I$ .  $\square$

Applying Proposition 6 to the definition of a state of  $D$ , we obtain the following concise definition for the case where  $D$  is a definite action description.

**Proposition 7** *Let  $D$  be a definite action description of signature  $(\sigma^f, \sigma^a)$ . Let  $s$  be a set of atoms of  $\sigma^f$  representing an interpretation of  $\sigma^f$ .  $s$  is a state of  $D$  iff*

$$s = T_{static}(s) \cup Simple(s) \cup Trivial(\sigma^f)$$

□

Notice that the above characterisation of a state does not hold for an *arbitrary* set  $s$  of atoms of  $\sigma^f$  —  $s$  must be an interpretation (a consistent and complete set of atoms) of  $\sigma^f$ .

It is perhaps helpful to see that the conditions for an interpretation  $s$  of  $\sigma^f$  to be a state of  $D$  can be equivalently stated as follows. First, in the case where there are no trivial atoms in  $\sigma^f$ :

- $T_{static}(s) \subseteq s$
- $s - Simple(s) = T_{static}(s) - Simple(s)$

In other words, the state  $s$  must satisfy the static laws of  $D$ , and the set of statically determined atoms in  $s$ ,  $s - Simple(s)$ , must be completely determined by the static laws of  $D$ . When  $Trivial(\sigma^f) \neq \emptyset$ :

- $T_{static}(s) \subseteq s$
- $s - Simple(s) = (T_{static}(s) \cup Trivial(\sigma^f)) - Simple(s)$

By similar applications of Proposition 6, for any definite action description  $D$ , the condition that a set  $s'$  of atoms represents the only interpretation of  $\sigma^f$  satisfying  $T_{static}(s') \cup E(s, \varepsilon, s')$  is equivalently stated as  $s' = T_{static}(s') \cup E(s, \varepsilon, s') \cup Trivial(\sigma^f)$ ; and the condition that a set  $\varepsilon$  of action atoms represents the only interpretation of  $\sigma^a$  satisfying  $A(\varepsilon, s) \cup Exog(\varepsilon)$  is equivalently stated as  $\varepsilon = A(\varepsilon, s) \cup Exog(\varepsilon) \cup Trivial(\sigma^a)$ . So then:

**Proposition 8** *Let  $D$  be a definite action description of signature  $(\sigma^f, \sigma^a)$ .  $(s, \varepsilon, s')$  is a transition of  $D$  iff  $s$  and  $s'$  are interpretations of  $\sigma^f$  and  $\varepsilon$  is an interpretation of  $\sigma^a$  such that:*

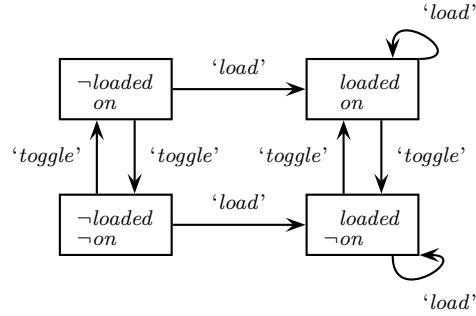
- $s = T_{static}(s) \cup Simple(s) \cup Trivial(\sigma^f)$  ( $s$  is a state of  $D$ )
- $\varepsilon = A(\varepsilon, s) \cup Exog(\varepsilon) \cup Trivial(\sigma^a)$
- $s' = T_{static}(s') \cup E(s, \varepsilon, s') \cup Trivial(\sigma^f)$

### 4.3 Example

Signature: simple Boolean fluent constants *loaded*, *on*; exogenous Boolean action constants *load*, *toggle*.

There are no static laws in this first example.

inertial *loaded*  
 inertial *on*  
*load* causes *loaded*  
*toggle* causes *on* if  $\neg$ *on*  
*toggle* causes  $\neg$ *on* if *on*



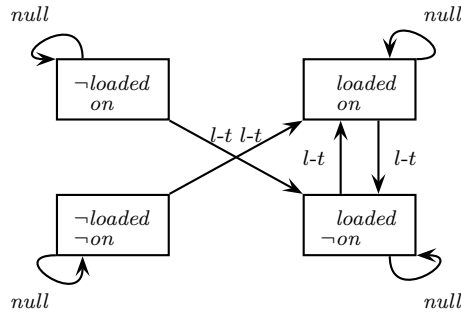
(Action *load* is supposed to mean something like ‘ensure that loaded’. Otherwise we would change the action description to *load* causes *loaded* if  $\neg$ *loaded*.)

In the diagram, transition labels ‘*load*’ and ‘*toggle*’ are shorthand for  $\{load, \neg toggle\}$  and  $\{\neg load, toggle\}$ , respectively.

Let’s consider state  $\{\neg loaded, on\}$  and event  $\{load, \neg toggle\}$ . There are no static laws so we can ignore  $T_{static}$  in this example, and there are no action dynamic laws to consider.

- Is  $(\{\neg loaded, on\}, \{load, \neg toggle\}, \{loaded, on\})$  a transition defined by this action description?  
 Yes, because  $E(\{\neg loaded, on\}, \{load, \neg toggle\}, \{loaded, on\}) = \{loaded, on\}$ .
- Is  $(\{\neg loaded, on\}, \{load, \neg toggle\}, \{loaded, \neg on\})$  a transition defined by this action description??  
 No, because  $E(\{\neg loaded, on\}, \{load, \neg toggle\}, \{loaded, \neg on\}) = \{loaded\} \neq \{loaded, \neg on\}$ .

There are two other kinds of transitions, not shown in the diagram above. They are transitions with labels  $\{load, toggle\}$  and  $\{\neg load, \neg toggle\}$ .



The label  $l-t$  in the diagram is shorthand for  $\{load, toggle\}$  and  $null$  is shorthand for  $\{\neg load, \neg toggle\}$ .

If we wanted to eliminate the  $null$  events, we could add the following law to the action description:

$$\perp \text{ if } \top \text{ after } \neg load \wedge \neg toggle$$

for which there is a standard abbreviation in  $\mathcal{C}+$ :

$$\text{nonexecutable } \neg load \wedge \neg toggle$$

**Remark: the reading of causes** The choice of the abbreviation  $\text{causes}$  in  $\mathcal{C}+$  can be misleading. For example, we might observe that whenever Jack leaves work, he is tired. We can express this by means of a fluent dynamic law

$$\text{tired}(Jack) \text{ if } \top \text{ after } \text{leaves\_work}(Jack)$$

In abbreviated form this is

$$\text{leaves\_work}(Jack) \text{ causes } \text{tired}(Jack)$$

But we would surely not want to say that Jack's action of leaving work *causes* him to be tired.

We will retain the abbreviation  $\text{causes}$  in what follows since most other accounts and examples of  $\mathcal{C}+$  make extensive use of it. However, we emphasise that the reading of  $\text{causes}$  in examples to follow should not be given any strong of reading of causation: it is the unabbreviated fluent dynamic law that is intended.

## 4.4 Abbreviations

The language  $\mathcal{C}+$  provides various abbreviations, of which  $\text{causes}$ ,  $\text{inertial}$ , and  $\text{nonexecutable}$  are the most common. Here is the full list, as presented in (Giunchiglia et al., 2004, Appendix B).

In the following,  $F$  and  $G$  are fluent formulas,  $\alpha$  and  $\alpha'$  are action formulas,  $\psi$  is a formula,  $f$  is a fluent constant, and  $a$  is an action constant.

$F$	$F \text{ if } \top$
$\alpha$	$\alpha \text{ if } \top$
$F \text{ after } \psi$	$F \text{ if } \top \text{ after } \psi$

default $F$	$F \text{ if } F$
default $F \text{ if } G$	$F \text{ if } F \wedge G$
inertial $f$	$f = v \text{ if } f = v \text{ after } f = v, \text{ for all } v \in \text{dom}(f)$
inertial $f \text{ if } \psi$	$f = v \text{ if } f = v \text{ after } f = v \wedge \psi, \text{ for all } v \in \text{dom}(f)$
$\alpha \text{ causes } F$	$F \text{ if } \top \text{ after } \alpha$
$\alpha \text{ causes } F \text{ if } \psi$	$F \text{ if } \top \text{ after } \alpha \wedge \psi$
nonexecutable $\alpha$	$\perp \text{ if } \top \text{ after } \alpha \quad (\text{or: } \alpha \text{ causes } \perp)$
nonexecutable $\alpha \text{ if } \psi$	$\perp \text{ if } \top \text{ after } \alpha \wedge \psi \quad (\text{or: } \alpha \text{ causes } \perp \text{ if } \psi)$
$\alpha \text{ may cause } F$	$F \text{ if } F \text{ after } \alpha$
$\alpha \text{ may cause } F \text{ if } \psi$	$F \text{ if } F \text{ after } \alpha \wedge \psi$

The last of these, *may cause*, is used for specifying the effects of non-deterministic actions. An example will follow later.

The form *nonexecutable  $\alpha$*  is, strictly, not allowed according to (Giunchiglia et al., 2004) but is used in that paper.

We also have, where  $a$  is an action constant and  $f$  a fluent constant:

exogenous $a$	$a = v$ if $a = v$ , for all $v \in dom(a)$ (or: default $a = v$ , for all $v \in dom(a)$ )
exogenous $a$ if $\psi$	$a = v$ if $a = v \wedge \psi$ , for all $v \in dom(a)$ (or: default $a = v$ if $\psi$ , for all $v \in dom(a)$ )
exogenous $f$	$f = v$ if $f = v$ , for all $v \in dom(f)$ (or: default $f = v$ , for all $v \in dom(f)$ )
exogenous $f$ if $G$	$f = v$ if $f = v \wedge G$ , for all $v \in dom(f)$ (or: default $f = v$ if $G$ , for all $v \in dom(f)$ )

Exogenous fluent constants are those whose values can vary from one state to another, but not under the effects of actions in the action description. Whether it is raining could be represented by an exogenous (Boolean) fluent constant, for example.

For reference, here are the remaining abbreviations given in (Giunchiglia et al., 2004, Appendix B). They are not used in this paper.

$\alpha$ causes $\alpha'$ if $\psi$	$\alpha'$ if $\alpha \wedge \psi$
$\alpha$ may cause $\alpha'$ if $\psi$	$\alpha'$ if $\alpha' \wedge \alpha \wedge \psi$
default $\alpha$	$\alpha$ if $\alpha$
default $\alpha$ if $\psi$	$\alpha$ if $\alpha \wedge \psi$
default $F$ if $G$ after $\psi$	$F$ if $F \wedge G$ after $\psi$
default $F$ after $\psi$	default $F$ if $\top$ after $\psi$ (or: $F$ if $F$ after $\psi$ )
constraint $F$	$\perp$ if $\neg F$
constraint $F$ after $\psi$	$\perp$ if $\neg F$ after $\psi$
rigid $f$	$\perp$ if $f \neq v$ after $f = v$ , for all $v \in dom(f)$ (or: constraint $f = v$ after $f = v$ , for all $v \in dom(f)$ )
always $\psi$	$\perp$ after $\neg\psi$

For  $f$  a *statically determined* Boolean fluent constant, and  $a$  a Boolean action constant:

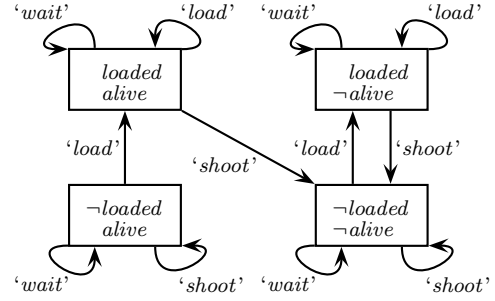
$F$ if $G$ unless $f$	$F$ if $G \wedge \neg f$ default $\neg f$
$F$ if $G$ after $\psi$ unless $a$	$F$ if $G$ after $\psi \wedge \neg a$ default $\neg a$
$\alpha$ if $\psi$ unless $a$	$\alpha$ if $\psi \wedge \neg a$ default $\neg a$

Note that many of the above abbreviations are ambiguous as given in (Giunchiglia et al., 2004, Appendix B), in the special case that components are  $\top$  or  $\perp$ .

#### 4.5 Example (‘Yale Shooting Problem’)

Signature: simple Boolean fluent constants *loaded*, *alive*; exogenous Boolean action constants *load*, *shoot*, *wait*.

inertial *loaded*  
 inertial *alive*  
*load* causes *loaded*  
*shoot* causes  $\neg$ *alive* if *loaded*  
*shoot* causes  $\neg$ *loaded*  
 nonexecutable *shoot*  $\wedge$  *load*  
 nonexecutable *wait*  $\wedge$  *shoot*  
 nonexecutable *wait*  $\wedge$  *load*  
 nonexecutable  $\neg$ *wait*  $\wedge$   $\neg$ *shoot*  $\wedge$   $\neg$ *load*



The diagram uses the same shorthand convention for transition labels as used earlier.

It is not possible to load and shoot a gun at the same time: *shoot*  $\wedge$  *load* events are eliminated by the first of the **nonexecutable** laws.

Alternatively, we could dispense with the action constant *wait* and represent it instead by the ‘null’ event  $\{\neg$ *shoot*,  $\neg$ *load* $\}$ . The last three lines of the action description could then be deleted.

Consider now the query languages introduced in Section 3.3. Let  $D_{\text{YSP}}$  be the action description above. We write  $D_{\text{YSP}} \models \varphi$  (resp.  $D_{\text{YSP}} \models_m \psi$ ) when  $\mathcal{T}_{\text{YSP}}$  is the transition system defined by  $D_{\text{YSP}}$  and  $\mathcal{T}_{\text{YSP}} \models \varphi$  (resp.  $\mathcal{T}_{\text{YSP}} \models_m \psi$ ). One can see from the diagram above that we have, for example:

$$D_{\text{YSP}} \models [\textit{shoot}] [\textit{load}] \neg \textit{alive}$$

$$D_{\text{YSP}} \models_3 (\textit{load}[0] \wedge \textit{shoot}[2]) \rightarrow \neg \textit{alive}[3]$$

and indeed

$$D_{\text{YSP}} \models_m (\textit{load}[i] \wedge \textit{shoot}[j]) \rightarrow \neg \textit{alive}[k] \quad \text{for any } 0 \leq i < j < k \leq m$$

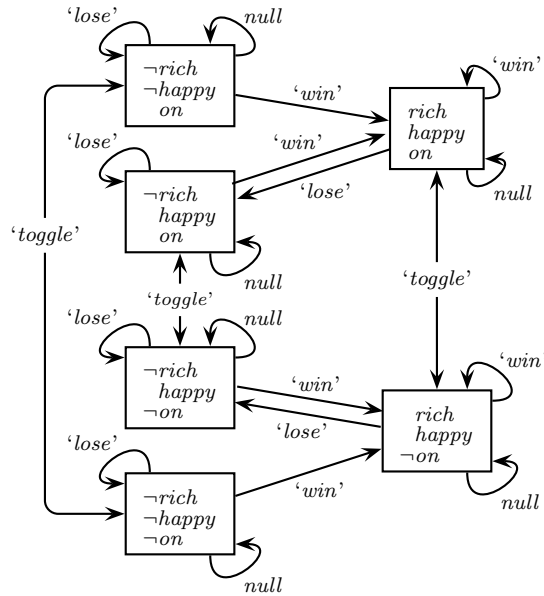
#### 4.6 Example

This is a completely artificial example, just for the sake of illustrating how static laws work as constraints on simple fluents.

Signature: simple Boolean fluent constants *rich*, *happy*, *on*; exogenous Boolean action constants *win*, *lose*, *toggle*.

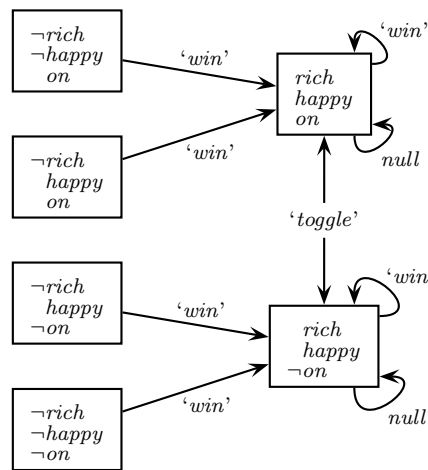


inertial *rich*  
 inertial *on*  
 inertial *happy*  
*win* causes *rich*  
*lose* causes  $\neg$ *rich*  
*toggle* causes *on* if  $\neg$ *on*  
*toggle* causes  $\neg$ *on* if *on*  
  
*happy* if *rich* (a static law)  
  
 nonexecutable  $\textit{lose} \wedge \textit{win}$



Because of the static law, there are only 6 states not  $2^3 = 8$ . The diagram does not show the transitions with labels  $\{\textit{toggle}, \textit{win}, \neg\textit{lose}\}$  and  $\{\textit{toggle}, \neg\textit{win}, \textit{lose}\}$ . This is just to reduce clutter.

What if we drop the law that specifies *happy* is inertial? In that case, since there is nothing to 'cause'  $\neg$ *happy* ( $\neg$ *happy*  $\notin T_{static}(s)$  for any  $s$  and  $\neg$ *happy*  $\notin E(s, \varepsilon, s')$  for any transition  $(s, \varepsilon, s')$ ), there can be no transitions to any state where  $\neg$ *happy* holds, and by a similar argument, no transitions to *happy* states where *rich* does not hold. There are thus no 'lose' transitions, and some 'null' and 'toggle' transitions are eliminated also, as shown in the following diagram.



without inertial *happy*

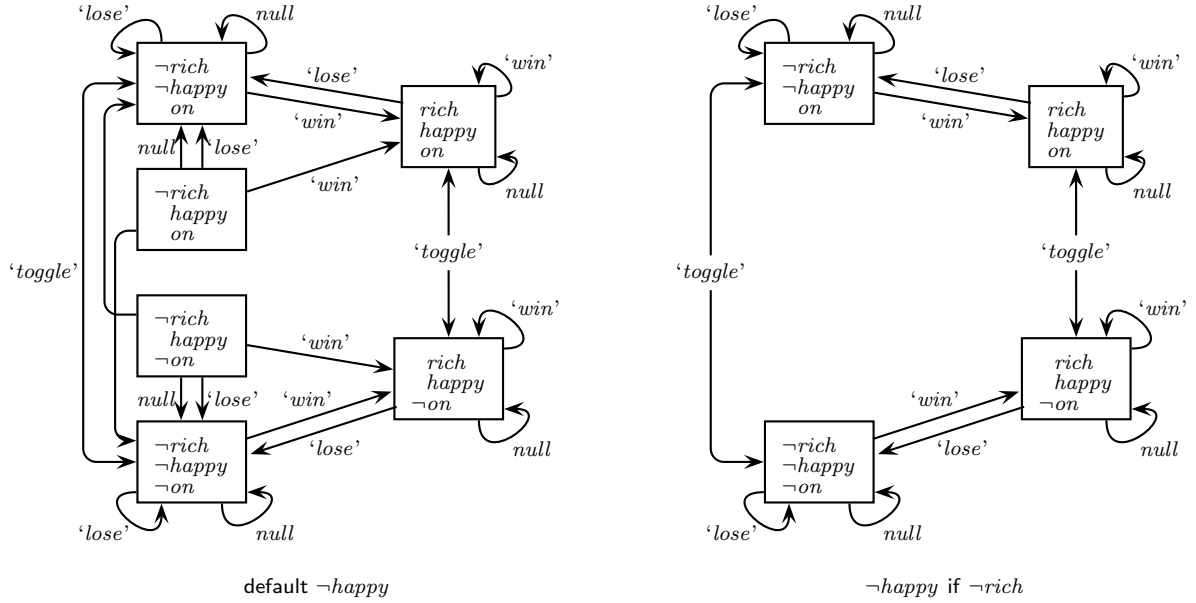
Suppose instead of inertial *happy* we add another static law

$$\text{default } \neg\text{happy} \quad (\text{i.e., } \neg\text{happy} \text{ if } \neg\text{happy})$$

Now we get transitions to  $\neg\text{happy}$  states, but still no transitions to *happy* states where *rich* does not hold. Also, some *null* and *'lose'* transitions have different effects, since *happy* goes to  $\neg\text{happy}$  by default. This is illustrated in the leftmost of the two diagrams below.

Suppose we make *happy* exogenous (in effect adding another law **default** *happy* to **default**  $\neg\text{happy}$ )? Now the static law *happy* if *rich* is strong enough to eliminate unwanted states but the *'toggle'* and *null* transitions, which were deterministic, now become non-deterministic. (We omit the diagram for this case.)

Finally, suppose instead of **default** *happy*, we add the static law  $\neg\text{happy}$  if  $\neg\text{rich}$ . Then there are only four states instead of six: states with  $\text{happy} \wedge \neg\text{rich}$  are eliminated. This is shown in the rightmost diagram below.



Now suppose we make *happy* a statically determined fluent constant. Then there are only two states in the transition system:  $\{\text{rich}, \text{happy}, \text{on}\}$  and  $\{\text{rich}, \text{happy}, \neg\text{on}\}$ . This is easy to check using the characterisation of a state,  $s = T_{\text{static}}(s) \cup \text{Simple}(s)$ :

$s$	$= T_{\text{static}}(s) \cup \text{Simple}(s)$	
$\{\text{rich}, \text{happy}, \text{on}\}$	$\{\text{happy}\} \cup \{\text{rich}, \text{on}\}$	✓
$\{\text{rich}, \text{happy}, \neg\text{on}\}$	$\{\text{happy}\} \cup \{\text{rich}, \neg\text{on}\}$	✓
$\{\neg\text{rich}, \text{happy}, \text{on}\}$	$\{\}$ $\cup$ $\{\neg\text{rich}, \text{on}\}$	×
$\{\neg\text{rich}, \text{happy}, \neg\text{on}\}$	$\{\}$ $\cup$ $\{\neg\text{rich}, \neg\text{on}\}$	×
$\{\neg\text{rich}, \neg\text{happy}, \text{on}\}$	$\{\}$ $\cup$ $\{\neg\text{rich}, \text{on}\}$	×
$\{\neg\text{rich}, \neg\text{happy}, \neg\text{on}\}$	$\{\}$ $\cup$ $\{\neg\text{rich}, \neg\text{on}\}$	×

(We omit the two interpretations satisfying  $rich \wedge \neg happy$  since they are obviously not states.)

Now add another static law: **default**  $\neg happy$  (i.e.,  $\neg happy$  if  $\neg happy$ ):

$s$	$=$	$T_{static}(s) \cup Simple(s)$	
$\{ rich, happy, on \}$	$=$	$\{ happy \} \cup \{ rich, on \}$	$\checkmark$
$\{ rich, happy, \neg on \}$	$=$	$\{ happy \} \cup \{ rich, \neg on \}$	$\checkmark$
$\{ \neg rich, happy, on \}$	$\neq$	$\{ \} \cup \{ \neg rich, on \}$	$\times$
$\{ \neg rich, happy, \neg on \}$	$\neq$	$\{ \} \cup \{ \neg rich, \neg on \}$	$\times$
$\{ \neg rich, \neg happy, on \}$	$=$	$\{ \neg happy \} \cup \{ \neg rich, on \}$	$\checkmark$
$\{ \neg rich, \neg happy, \neg on \}$	$=$	$\{ \neg happy \} \cup \{ \neg rich, \neg on \}$	$\checkmark$

Now we have another two states in the transition system: we have also  $\{ \neg rich, \neg happy, on \}$  and  $\{ \neg rich, \neg happy, \neg on \}$ . Notice that when  $happy$  is statically determined, the pair of static laws  $happy$  if  $rich$  and **default**  $\neg happy$  is equivalent to the pair  $happy$  if  $rich$  and  $\neg happy$  if  $\neg rich$ . Both say that  $happy$  is true in a state iff  $rich$  is true. (When  $happy$  is not statically determined, the two pairs are not equivalent, as the figure above demonstrates: we have **default**  $\neg happy$  in the diagram on the left and  $\neg happy$  if  $\neg rich$  in the diagram on the right. In particular, **default**  $\neg happy$  admits states that are eliminated by the other formulation.)

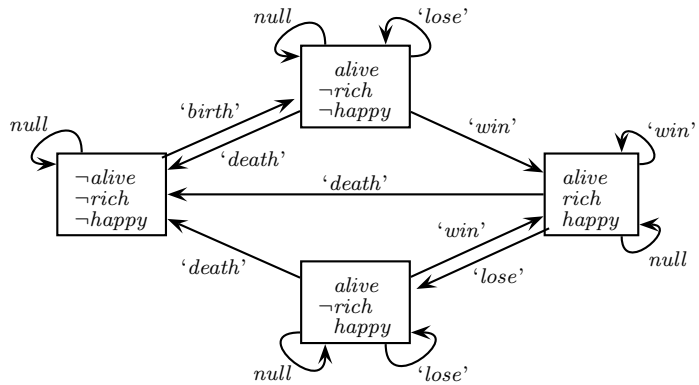
Finally, we cannot add inertial  $happy$  when  $happy$  is statically determined, since inertial declarations are fluent dynamic laws and cannot have statically determined fluent constants in their head.

#### 4.7 Example (Winning the lottery)

Winning the lottery causes one to become (or remain) rich. Losing one's wallet causes one to become (or remain) not rich. A person who is rich is happy.

Signature: simple Boolean fluent constants  $alive, rich, happy$ ; exogenous Boolean action constants  $birth, death, win, lose$ .

inertial  $alive$   
inertial  $rich$   
inertial  $happy$   
 $birth$  causes  $alive$   
nonexecutable  $birth$  if  $alive$   
 $death$  causes  $\neg alive$   
nonexecutable  $death$  if  $\neg alive$   
 $win$  causes  $rich$   
nonexecutable  $win$  if  $\neg alive$   
 $lose$  causes  $\neg rich$   
nonexecutable  $lose$  if  $\neg alive$   
 $happy$  if  $rich$   
 $\neg rich$  if  $\neg alive$   
 $\neg happy$  if  $\neg alive$   
nonexecutable  $birth \wedge death$   
nonexecutable  $birth \wedge win$   
nonexecutable  $birth \wedge lose$   
nonexecutable  $win \wedge lose$



The example is constructed partly to show how  $\mathcal{C}+$  deals with indirect effects of actions (ramifications). Here we have two fluents, *happy* and *rich*, whose values are constrained by the presence of static laws. The first is not affected directly by any action, whereas the second one is.

Because of the static laws, there are only four states in the transition system and not  $2^3 = 8$ . Transition labels ‘*birth*’, ‘*death*’, ‘*win*’, ‘*lose*’ in the diagram are shorthand for the events  $\{birth, \neg death, \neg win, \neg lose\}$ ,  $\{\neg birth, death, \neg win, \neg lose\}$ ,  $\{\neg birth, \neg death, win, \neg lose\}$ ,  $\{\neg birth, \neg death, \neg win, lose\}$ , respectively. The label *null* is shorthand for  $\{\neg birth, \neg death, \neg win, \neg lose\}$ .

Notice that as formulated here, the example allows for reincarnation: a person can be born, die, and be born again. The possibility of reincarnation can be eliminated easily enough, for example by adding another fluent constant *dead* together with a static law  $\perp$  if *alive*  $\wedge$  *dead*; the simpler version with reincarnation is perfectly adequate for present purposes.

The diagram does not show transitions of type *death*  $\wedge$  *lose* (i.e., transitions with label  $\{\neg birth, death, \neg win, lose\}$ ). Their effects in this example are exactly the same as those of ‘*death*’ transitions. There are no transitions of type *win*  $\wedge$  *death*: they would lead to states with *rich*  $\wedge$   $\neg$ *alive*, and there can be no such states because of the static law  $\perp$  if *rich*  $\wedge$   $\neg$ *alive*. However, it seems unreasonable to insist that transitions of type *win*  $\wedge$  *death* are non-executable. We can admit the possibility of *win*  $\wedge$  *death* transitions by replacing the *causes* law for *win* by the following:

$$win \text{ causes } rich \text{ if } \neg death$$

or equivalently *win*  $\wedge$   $\neg$ *death* *causes* *rich*. (These statements are equivalent because they are shorthand for *rich* if  $\top$  after *win*  $\wedge$   $\neg$ *death* and *rich* if  $\top$  after  $(win \wedge \neg death) \wedge \top$  respectively.) The effects of the ‘*win*’ transitions are unchanged, but the transition system now contains transitions of type *win*  $\wedge$  *death*: their effects are exactly the same as those of ‘*death*’ transitions and transitions of type *death*  $\wedge$  *lose* transitions. More generally, we might want to say that winning the lottery makes a person rich *by default*, and is overridden in exceptional circumstances, such as dying as the lottery is being won.

Notice that *happy* is declared inertial, and so still persists even if one becomes not rich. That is why the ‘*lose*’ transition from state  $\{alive, rich, happy\}$  results in the state  $\{alive, \neg rich, happy\}$ . We could of course modify the action description so that *happy* is no longer inertial but defined to be true if and only if *rich* is true. Or we might prefer to make *happy* non-inertial and let the ‘*lose*’ transition be non-deterministic. The interactions between these various adjustments are rather subtle, however, and not always immediately obvious. We will discuss some options below.

For illustration of the semantics, let us consider state  $\{\neg alive, \neg rich, \neg happy\}$  and event ‘*birth*’ (i.e.,  $\{birth, \neg death, \neg win, \neg lose\}$ ). There are no action dynamic laws to consider; the fluent dynamic laws that need to be considered are the law *birth* *causes* *alive*, the three inertial statements, and the nonexecutable statements.

- Is  $(\{\neg alive, \neg rich, \neg happy\}, \text{‘birth’}, \{alive, \neg rich, \neg happy\})$  a transition defined by the action description? It is, because  $E(\{\neg alive, \neg rich, \neg happy\}, \text{‘birth’}, \{alive, \neg rich, \neg happy\}) = \{alive, \neg rich, \neg happy\}$ .

- Is  $(\{\neg\text{alive}, \neg\text{rich}, \neg\text{happy}\}, \text{'birth'}, \{\text{alive}, \neg\text{rich}, \text{happy}\})$  a transition defined by this action description? No, because  $E(\{\neg\text{alive}, \neg\text{rich}, \neg\text{happy}\}, \text{'birth'}, \{\text{alive}, \neg\text{rich}, \text{happy}\}) = \{\text{alive}, \neg\text{rich}\} \neq \{\text{alive}, \neg\text{rich}, \text{happy}\}$ .
- Is there a transition of the form  $(\{\text{alive}, \neg\text{rich}, \neg\text{happy}\}, \text{'birth'}, s')$  defined by this action description? No: the law **nonexecutable** *birth* if *alive* means that  $E(\{\text{alive}, \neg\text{rich}, \neg\text{happy}\}, \text{'birth'}, s')$  must contain the element  $\perp$ , and no state  $s'$  can contain  $\perp$ .

Simple though it is, the example still raises a number of interesting questions. For example:

- Given the stated results of *birth* and *death*, the law **nonexecutable** *birth*  $\wedge$  *death* is implied by the remaining parts of the action description, in the sense that it could be deleted without affecting the transition system that is defined. This is because we have the implied laws  $(\text{birth} \wedge \text{death}) \text{ causes } (\text{alive} \wedge \neg\text{alive})$ , or equivalently  $(\text{birth} \wedge \text{death}) \text{ causes } \perp$ , which is just **nonexecutable** *birth*  $\wedge$  *death*. Similarly, in the original formulation, we have the implied law  $(\text{win} \wedge \text{lose}) \text{ causes } (\text{rich} \wedge \neg\text{rich})$ , which gives **nonexecutable** *win*  $\wedge$  *lose*. But note that after the adjustment to take into account the possible overriding effects of *death*, we have only the implied law  $(\text{win} \wedge \text{lose}) \text{ causes } (\text{rich} \wedge \neg\text{rich})$  if  $\neg\text{death}$ , or equivalently,  $(\text{win} \wedge \neg\text{death} \wedge \text{lose}) \text{ causes } (\text{rich} \wedge \neg\text{rich})$ , that is, **nonexecutable** *win*  $\wedge$   $\neg\text{death} \wedge \text{lose}$ , or equivalently **nonexecutable** *win*  $\wedge$  *lose* if  $\neg\text{death}$ . This clearly does *not* imply **nonexecutable** *win*  $\wedge$  *lose*.
- In the original formulation we have the implied law **nonexecutable** *win*  $\wedge$  *death* because: we have the implied law  $(\text{win} \wedge \text{death}) \text{ causes } (\text{rich} \wedge \neg\text{alive})$ , and there is no state in which  $\text{rich} \wedge \neg\text{alive}$  holds because of the static law  $\perp$  if  $\text{rich} \wedge \neg\text{alive}$ ; so  $(\text{win} \wedge \text{death}) \text{ causes } \perp$ .
- The laws **nonexecutable** *birth*  $\wedge$  *win* and **nonexecutable** *birth*  $\wedge$  *lose* are also implied. We have *birth* **causes**  $\perp$  if *alive* and *win* **causes**  $\perp$  if  $\neg\text{alive}$ . So we can derive  $(\text{birth} \wedge \text{win}) \text{ causes } \perp$  if *alive* and  $(\text{birth} \wedge \text{win}) \text{ causes } \perp$  if  $\neg\text{alive}$ . Now we derive  $(\text{birth} \wedge \text{win}) \text{ causes } \perp$  if  $(\text{alive} \vee \neg\text{alive})$ . And similarly for **nonexecutable** *birth*  $\wedge$  *lose*.

Such questions are addressed in a separate paper (Sergot and Craven, 2005) on the logical properties of ‘causal theories’ and of the language  $\mathcal{C}+$ .

Instead of the static laws  $\neg\text{rich}$  if  $\neg\text{alive}$  and  $\neg\text{happy}$  if  $\neg\text{alive}$ , suppose we have  $\perp$  if  $\text{rich} \wedge \neg\text{alive}$  and  $\perp$  if  $\text{happy} \wedge \neg\text{alive}$ . What difference does it make, in this example, and in general? What about *alive* if *rich* and *alive* if *happy*?

First, notice that all three formulations are equivalent from the point of view of *states*. A state  $s$  satisfies the static law  $F$  if  $G$  iff it satisfies  $\neg G$  if  $\neg F$  iff it satisfies  $\perp$  if  $G \wedge \neg F$  iff  $s \models G \rightarrow F$ . But they are not equivalent from the point of view of transitions. If we replace  $\neg\text{happy}$  if  $\neg\text{alive}$  by either of *alive* if *happy* or  $\perp$  if  $\text{happy} \wedge \neg\text{alive}$ , the only way that  $\neg\text{happy}$  can be ‘caused’ is by inertia. Consequently, we eliminate all the *death* transitions (‘*death*’, and transitions of type *death*  $\wedge$  *lose* and *death*  $\wedge$  *win*) from states in which *rich* holds. (Also, look at the causality reading of these rules.) If we replace  $\neg\text{rich}$  if  $\neg\text{alive}$  by either of *alive* if *rich* or  $\perp$  if  $\text{rich} \wedge \neg\text{alive}$ , the argument is similar but a little more complicated. The only way that  $\neg\text{rich}$  can be ‘caused’ is by a *lose* transition or by inertia. Consequently, the transitions ‘*death*’ and transitions of type *death*  $\wedge$  *win* become non-executable in the state  $\{\text{alive}, \text{rich}, \text{happy}\}$ .

There is one way in which we can use constraints  $\perp$  if  $rich \wedge \neg alive$  and  $\perp$  if  $happy \wedge \neg alive$  (or  $alive$  if  $rich$  and  $alive$  if  $happy$ ) without losing transitions. That is by adding a pair of extra fluent dynamic laws: either

$$death \text{ causes } \neg rich \quad \text{and} \quad death \text{ causes } \neg happy$$

or the weaker pair

$$death \text{ may cause } \neg rich \quad \text{and} \quad death \text{ may cause } \neg happy$$

The second pair seems preferable to the first, but neither is entirely satisfactory since they require all ramifications of  $death$  to be identified in advance and then modelled explicitly using causal laws.

### Statically determined fluents

To illustrate the difference between ‘simple’ and ‘statically determined’ fluents, suppose now that fluent constant  $happy$  is statically determined not simple. We have to remove the inertial  $happy$  declaration since this is not well formed when  $happy$  is statically determined. Suppose we have the static laws  $happy$  if  $rich$ ,  $\perp$  if  $rich \wedge \neg alive$ , and  $\perp$  if  $happy \wedge \neg alive$ . There is only one state in the transition system defined by this action description, namely  $\{alive, rich, happy\}$ .

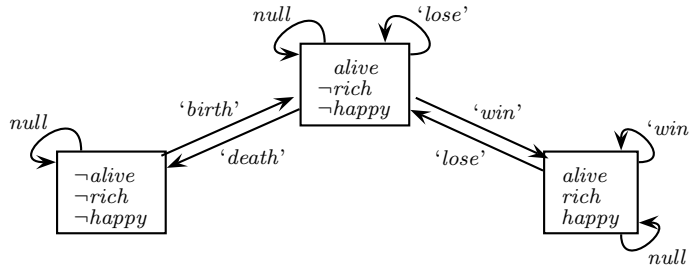
$$\begin{array}{rcl} s & = & T_{static}(s) \cup Simple(s) \\ \hline \{\neg alive, \neg rich, \neg happy\} & \neq & \{\} \cup \{\neg alive, \neg rich\} \quad \times \\ \{alive, \neg rich, \neg happy\} & \neq & \{\} \cup \{alive, \neg rich\} \quad \times \\ \{alive, \neg rich, happy\} & \neq & \{\} \cup \{alive, \neg rich\} \quad \times \\ \{alive, rich, happy\} & = & \{happy\} \cup \{alive, rich\} \quad \checkmark \end{array}$$

(We omit the four interpretations satisfying  $\neg alive \wedge rich$ ,  $\neg alive \wedge happy$ , or  $rich \wedge \neg happy$  since they are obviously not states.)

If however we complete the definition of  $happy$  by adding another static law  $\neg happy$  if  $\neg happy$  (i.e., default  $\neg happy$ ) then the only state eliminated from the original transition system is  $\{alive, \neg rich, happy\}$ , and we get three states.

$$\begin{array}{rcl} s & = & T_{static}(s) \cup Simple(s) \\ \hline \{\neg alive, \neg rich, \neg happy\} & = & \{\neg happy\} \cup \{\neg alive, \neg rich\} \quad \checkmark \\ \{alive, \neg rich, \neg happy\} & = & \{\neg happy\} \cup \{alive, \neg rich\} \quad \checkmark \\ \{alive, \neg rich, happy\} & \neq & \{\} \cup \{alive, \neg rich\} \quad \times \\ \{alive, rich, happy\} & = & \{happy\} \cup \{alive, rich\} \quad \checkmark \end{array}$$

We have the following transition system:



Notice that with the static laws as given above, there is no ‘*death*’ transition from  $\{alive, rich, happy\}$  to  $\{\neg alive, \neg rich, \neg happy\}$ : there is nothing to ‘cause’ the transition from *rich* to  $\neg rich$ . We can obtain the additional ‘*death*’ transition by adding, for example, the static law  $\neg rich$  if  $\neg alive$  (which implies, but is not implied by,  $\perp$  if  $rich \wedge \neg alive$ ).

#### 4.8 Example (Going to work 1)

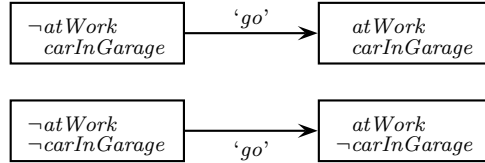
This is a modified version of an example used in (Giunchiglia and Lifschitz, 1998) to illustrate non-deterministic actions in  $\mathcal{C}+$  and its precursor  $\mathcal{C}$ . There is a slightly more complicated version in (Giunchiglia et al., 2004). The simpler version is sufficient for the points we want to make here.

Let the (exogenous) Boolean action constant *go* represent ‘Jack goes to work’. Jack can go to work by walking or, if his car is in his garage, he can drive. For simplicity, to simplify the diagrams, we ignore the possibility that Jack goes in the opposite direction.

The following action description

*inertial atWork*  
*inertial carInGarage*  
*go causes atWork*  
*nonexecutable go if atWork*

makes ‘*go*’ deterministic in all states, as shown in the following diagram

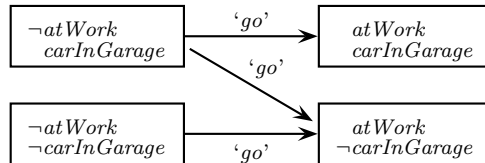


(Reflexive edges corresponding to the event  $\{\neg go\}$  are not shown.)

However, what we expect (or want) is that ‘*go*’ is non-deterministic in those states where *carInGarage* is true, because here Jack can either walk to work or drive and thereby move his car. To obtain this effect one adds another statement to the action description:

*go may cause ¬carInGarage if carInGarage*

*go may cause ¬carInGarage if carInGarage* is an abbreviation for the fluent dynamic law  $\neg carInGarage$  if  $\neg carInGarage$  after  $carInGarage \wedge go$ . With the additional statement we obtain the following transition diagram ( $\{\neg go\}$  events omitted):

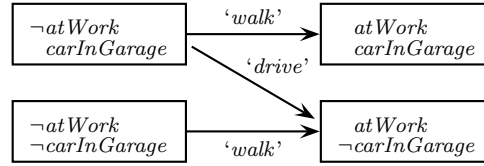


## 4.9 Example (Going to work 2)

An alternative formulation of the preceding example would distinguish between walking to work and driving to work. Let us have two (exogenous) Boolean action constants *walk* and *drive* to represent walking and driving to work respectively. The action description

*inertial atWork*  
*inertial carInGarage*  
*walk causes atWork*  
*drive causes atWork*  
*drive causes ¬carInGarage if carInGarage*  
 nonexecutable *walk* if *atWork*  
 nonexecutable *drive* if *atWork*  
 nonexecutable *drive* if  $\neg carInGarage$   
 nonexecutable *walk*  $\wedge$  *drive*

defines the following transition system ( $\{\neg walk, \neg drive\}$  omitted):



The first two *causes* laws could be replaced by the (equivalent) law:

$(walk \vee drive)$  *causes atWork*

We could also represent that *walk* and *drive* are both kinds of *go* by means of the action dynamic laws:

*go* if *walk*  
*go* if *drive*  
 $\neg go$  if  $\neg go$

The third law ( $\neg go$  if  $\neg go$ ) can be dropped when *go* is an exogenous action constant. When *go* is exogenous, we can also get the same effects by writing the pair of fluent dynamic laws:

nonexecutable *walk*  $\wedge$   $\neg go$   
 nonexecutable *drive*  $\wedge$   $\neg go$

We might also wish to add (in the absence of another kind of *go*, such as cycling):

nonexecutable *go*  $\wedge$   $\neg walk$   $\wedge$   $\neg drive$

This would not change the form of the transition system shown above except to replace transition labels '*walk*' and '*drive*' by  $\{go, walk\}$  and  $\{go, drive\}$  respectively.

Notice that the transition label  $\{go, walk\}$ , and indeed the action dynamic laws, cannot distinguish between two concurrent but unrelated actions *go* and *walk* and one action 'go by walking'. We return to this point when we discuss 'counts as' and act generation in Section 8 below.



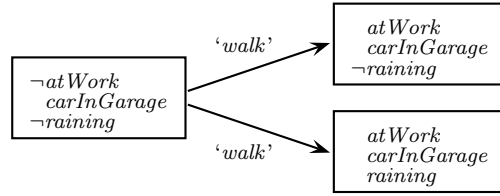
#### 4.10 Example (Going to work 3)

There is another source of non-determinism: fluents which vary from state to state but are not ‘caused’ by any kind of action. An example is ‘raining’. These are fluents which might properly be termed ‘exogenous’.

For example: add to the signature a simple Boolean constant *raining*, and to the action description the pair of static laws:

$$\begin{aligned} & \textit{raining} \text{ if } \textit{raining} \\ \neg \textit{raining} \text{ if } \neg \textit{raining} \end{aligned}$$

We obtain a transition system containing, for example, the fragment:



The pair of static laws for *raining* above may also be written more concisely in  $\mathcal{C}+$  as:

$$\text{exogenous } \textit{raining}$$

Recall that in general, for a fluent constant  $f$ , the abbreviation **exogenous**  $f$  stands for the set of static laws  $f = v \text{ if } f = v$ , for every  $v \in \text{dom}(f)$ .

## 5 Nonmonotonic causal theories

The language  $\mathcal{C}+$  can be regarded as a higher-level notation for particular classes of theories in the formalism of ‘nonmonotonic causal theories’, and indeed this is how it is presented in (Giunchiglia et al., 2004). The formalism is also referred to as ‘causal theories’, and as ‘the logic of causal explanation’ in (Lifschitz, 1997). Turner (1999) has a more general formalism which he calls the ‘logic of universal causation’.

**Syntax** A causal theory of signature  $\sigma$  is a set of expressions (‘causal rules’) of the form

$$F \Leftarrow G$$

where  $F$  and  $G$  are formulas of signature  $\sigma$ . A rule of this form is to be read as saying that  $F$  is ‘caused’ if  $G$  is true (which is not the same as saying that  $G$  ‘causes’  $F$ ).

**Semantics** Let  $\Gamma$  be a causal theory of signature  $\sigma$  and let  $X$  be an interpretation of its signature. The *reduct*  $\Gamma^X$  is the set of heads of rules of  $\Gamma$  whose bodies are satisfied by the interpretation  $X$ :

$$\Gamma^X =_{\text{def}} \{F \mid F \Leftarrow G \text{ is a rule in } \Gamma \text{ and } X \models G\}$$

$X$  is a *model* of  $\Gamma$  iff  $X$  is the unique model of  $\Gamma^X$ , i.e., the unique interpretation of  $\sigma$  that satisfies all formulas in  $\Gamma^X$ .

The rationale is this.  $\Gamma^X$  is the set of formulas that are ‘caused’ or ‘causally explained’ according to the rules of  $\Gamma$ , under interpretation  $X$ . If  $\Gamma^X$  has no models, or more than one model, or a unique model different from  $X$ , then  $X$  is not considered to be a model of  $\Gamma$ .  $\Gamma$  is *consistent* or *satisfiable* if it has a model.

**Definite clausal theories** A causal theory  $\Gamma$  is *definite* if

- the head of every rule of  $\Gamma$  is an atom or  $\perp$ , and
- no atom is the head of infinitely many rules of  $\Gamma$ .

(Compare with the definition of definite action descriptions given earlier.) Recall that as defined earlier, when  $p$  is a Boolean constant,  $\neg p$  is shorthand for the atom  $p = \text{f}$  and so covered by the definition.

## 5.1 Translation into (classical) propositional logic

*Definite causal theories* can be translated via the process of ‘literal completion’ into expressions of (classical) propositional logic. The process is analogous to the Clark (predicate) completion that provides the original semantics for negation by failure in logic programs.

The completion of a definite causal theory  $\Gamma$ ,  $\text{comp}(\Gamma)$ , is defined as follows.

Consider a definite causal theory  $\Gamma$  of signature  $\sigma$ . For each non-trivial atom  $A$  the *completion formula* for  $A$  is the formula

$$A \leftrightarrow G_1 \vee \dots \vee G_n$$

where  $G_1, \dots, G_n$  ( $n \geq 0$ ) are the bodies of the rules of  $\Gamma$  which have head  $A$ . The *completion* of  $\Gamma$ ,  $\text{comp}(\Gamma)$ , is the set of formulas obtained by taking the completion formulas for all non-trivial atoms of  $\sigma$  along with the formula  $\neg F$  for each rule of the form  $\perp \Leftarrow F$  in  $\Gamma$ .

Notice:

1. The completion formula for any non-trivial atom  $A$  that is the head of no rule in  $\Gamma$  (the case  $n = 0$ ) is  $A \leftrightarrow \perp$ .
2. For a *trivial* atom  $A$  nothing is included in the completion. Alternatively, we could add  $A$  to the completion. It would make no difference to the models of the completion but might perhaps make it clearer.
3. We are treating Boolean constants as a special case of two-valued constants, with  $\neg p$  standing for  $p = \text{f}$  when  $p$  is a Boolean constant.

**Proposition** (Proposition 6 in (Giunchiglia et al., 2004)) *The models of a definite causal theory are precisely the models of its completion.*

Careful: the causal theory  $\{p \Leftarrow q\}$  has no models. The completion of this theory is *not*  $\{p \leftrightarrow q\}$ , which has models, but  $\{p \leftrightarrow q, \neg p \leftrightarrow \perp, q \leftrightarrow \perp, \neg q \leftrightarrow \perp\}$ , which has no models.

As already observed, the task of finding a model of the completion  $comp(\Gamma)$  of a definite causal theory  $\Gamma$ —and generally the task of finding a model of any set of formulas of a multi-valued signature—can be reduced to the task of finding a model of a set of classical propositional formulas.

## 6 Translation of $\mathcal{C}+$ to causal theories

For any action description  $D$  in  $\mathcal{C}+$ , and any non-negative integer  $m$ , it is possible to construct a causal theory  $\Gamma_m^D$  such that the models of  $\Gamma_m^D$  correspond to the paths of length  $m$  of the transition system defined by  $D$ . The language  $\mathcal{C}+$  can thus be regarded as a higher-level notation for defining causal theories of a particular kind, and indeed this is exactly as it is presented in (Giunchiglia et al., 2004).

When the action signature of  $D$  is  $(\sigma^f, \sigma^a)$  the signature of  $\Gamma_m^D$  is obtained by time-stamping every fluent and action constant of  $D$  with non-negative integers between 0 and  $m$ : in the notation of Section 3.3 this is the signature  $\sigma_m = \sigma_m^f \cup \sigma_m^a$ , where  $\sigma_m^f = \sigma^f[0] \cup \dots \cup \sigma^f[m]$  and  $\sigma_m^a = \sigma^a[0] \cup \dots \cup \sigma^a[m-1]$ . Note in particular that the signature  $\sigma_0$  of  $\Gamma_0^D$  is  $\sigma^f[0]$  and the signature  $\sigma_1$  of  $\Gamma_1^D$  is  $\sigma^f[0] \cup \sigma^a[0] \cup \sigma^f[1]$ .

Given an action description  $D$ , and a non-negative integer  $m$ , the causal theory  $\Gamma_m^D$  is constructed as follows. For every static law  $F$  if  $G$  in  $D$ , include a causal rule of the form

$$F[i] \Leftarrow G[i] \quad \text{for every } i \in 0..m.$$

For every fluent dynamic law  $F$  if  $G$  after  $\psi$  in  $D$ , include a causal rule of the form

$$F[i+1] \Leftarrow G[i+1] \wedge \psi[i] \quad \text{for every } i \in 0..m-1.$$

For every action dynamic law  $\alpha$  if  $\psi$  in  $D$ , include a causal rule of the form

$$\alpha[i] \Leftarrow \psi[i] \quad \text{for every } i \in 0..m-1.$$

We also require the following ‘exogeneity laws’. For every *simple* fluent constant  $f$  and every  $v \in dom(f)$ , include a causal rule:

$$f[0] = v \Leftarrow f[0] = v$$

And for every *exogenous* action constant  $a$  and every  $v \in dom(a)$ , include a causal rule:

$$a[i] = v \Leftarrow a[i] = v \quad \text{for every } i \in 0..m-1.$$

As observed earlier, the presentation of  $\mathcal{C}+$  in (Giunchiglia et al., 2004) does not treat the specification of exogenous action constants as part of the signature but instead requires statements of the form

exogenous  $a$

to be included explicitly in the action description for every action constant  $a$  intended to be exogenous. *exogenous  $a$*  is shorthand for the fluent dynamic laws

$$a = v \text{ if } a = v, \quad \text{for all } v \in dom(a)$$

One can see that the translation to the causal theory  $\Gamma_m^D$  is the same whichever treatment of exogenous action constants is taken.

For convenience, we list here the translated forms of the most commonly used abbreviations of  $\mathcal{C}+$ :

default $F$	$F[i] \Leftarrow F[i]$
default $F$ if $G$	$F[i] \Leftarrow F[i] \wedge G[i]$
inertial $f$	$f[i+1] = v \Leftarrow f[i+1] = v \wedge f[i] = v$ , for all $v \in \text{dom}(f)$
inertial $f$ if $\psi$	$f[i+1] = v \Leftarrow f[i+1] = v \wedge f[i] = v \wedge \psi[i]$ , for all $v \in \text{dom}(f)$
$\alpha$ causes $F$	$F[i+1] \Leftarrow \alpha[i]$
$\alpha$ causes $F$ if $\psi$	$F[i+1] \Leftarrow \alpha[i] \wedge \psi[i]$
nonexecutable $\alpha$	$\perp \Leftarrow \top \wedge \alpha[i]$
nonexecutable $\alpha$ if $\psi$	$\perp \Leftarrow \top \wedge \alpha[i] \wedge \psi[i]$
$\alpha$ may cause $F$	$F[i+1] \Leftarrow F[i+1] \wedge \alpha[i]$
$\alpha$ may cause $F$ if $\psi$	$F[i+1] \Leftarrow F[i+1] \wedge \alpha[i] \wedge \psi[i]$

Clearly any interpretation  $X$  of the signature  $\sigma_m$  of  $\Gamma_m^D$  can be written in the form

$$s_0[0] \cup \varepsilon_0[0] \cup s_1[1] \cup \dots \cup s_{m-1}[m-1] \cup \varepsilon_{m-1}[m-1] \cup s_m[m]$$

where each  $s_i[i]$  and  $\varepsilon_i[i]$  is an interpretation of  $\sigma^f[i]$  and  $\sigma^a[i]$ , respectively.

Now we relate models of causal theories  $\Gamma_m^D$  to the definitions of states and transitions of an action description  $D$  as given in Section 4.2. The second part of Theorem 9 below corresponds to Proposition 8 of (Giunchiglia et al., 2004). The statements and proofs are different, however, because in (Giunchiglia et al., 2004) states and transitions are defined in terms of models of  $\Gamma_0^D$  and  $\Gamma_1^D$  respectively, and not explicitly in terms of transition systems as in Section 4.2. The following result establishes that the definitions of Section 4.2 are equivalent to those of (Giunchiglia et al., 2004). The proofs are not difficult but they are rather long to present in detail and so we defer them to a separate section so as not to disrupt the flow of the presentation.

**Theorem 9** *Let  $D$  be an action description with signature  $(\sigma^f, \sigma^a)$ . An interpretation  $s[0]$  of  $\sigma^f[0]$  is a model of  $\Gamma_0^D$  iff  $s$  is a state of  $D$ .*

*An interpretation  $X \in \text{I}(\sigma_m)$  of the form  $s_0[0] \cup \varepsilon_0[0] \cup \dots \cup \varepsilon_{m-1}[m-1] \cup s_m[m]$  is a model of  $\Gamma_m^D$  iff  $(s_0 \varepsilon_0 \dots \varepsilon_{m-1} s_m)$  is a path of length  $m \geq 0$  of the transition system defined by  $D$ .*

In particular,  $\Gamma_1^D$  represents paths of length 1 of  $D$ , i.e., the transitions defined by the action description  $D$ .  $\Gamma_0^D$  represents the states of the transition system defined by  $D$ : an interpretation  $s$  of the fluent constants  $\sigma^f$  is a state of  $D$  precisely when  $s[0]$  is a model of  $\Gamma_0^D$ .

Since action descriptions in  $\mathcal{C}+$  can be regarded as abbreviations for causal theories of a certain kind (above), literal completions for laws of  $\mathcal{C}+$  can be defined also. Action descriptions in  $\mathcal{C}+$  are thereby translated into formulas of (classical) propositional logic, which in turn can be submitted to a standard

propositional satisfiability solver. Although this method is limited to action descriptions that are *definite*, this is not an important restriction in practical applications.

## 6.1 Proof of Theorem 9

First, consider the reduct  $(\Gamma_0^D)^{s[0]}$ . It consists of:

- formulas  $F[0]$  for every static law  $F$  if  $G$  in  $D$  such that  $s[0] \models_{\sigma^f[0]} G[0]$ , i.e.,  $F[0]$  for every  $F \in T_{static}(s)$ ; we write  $T_{static}(s)[0]$  for this set of formulas;
- atoms  $f[0] = v$  for every simple fluent constant  $f$  such that  $s[0] \models_{\sigma^f[0]} f[0] = v$ ; this is the set  $Simple(s[0])$ .

Let  $s'$  be an interpretation of  $\sigma^f[0]$ . It is easy to see that  $s'[0] \models_{\sigma^f[0]} (\Gamma_0^D)^{s[0]}$  iff  $s'[0] \models_{\sigma^f[0]} T_{static}(s)[0] \cup Simple(s[0])$  iff  $s' \models T_{static}(s) \cup Simple(s)$ .

And  $s'[0]$  is the only model of  $(\Gamma_0^D)^{s[0]}$  iff  $s$  is the only model of  $T_{static}(s) \cup Simple(s)$  iff  $s$  is a state of  $D$ .  $\square$

Now let  $X$  be an interpretation of  $\sigma_m$  of the form

$$s_0[0] \cup \varepsilon_0[0] \cup \dots \cup \varepsilon_{m-1}[m-1] \cup s_m[m]$$

First, consider the reduct  $(\Gamma_m^D)^X$ . It consists of:

- formulas  $F[i]$  for every static law  $F$  if  $G$  in  $D$  such that  $s_i[i] \models_{\sigma^f[i]} G[i]$  for  $0 \leq i \leq m$ , i.e.,  $F[i]$  for every  $F \in T_{static}(s_i)$  for  $0 \leq i \leq m$ ; this is the set  $T_{static}(s_0)[0] \cup \dots \cup T_{static}(s_1)[0] \cup \dots \cup T_{static}(s_m)[m]$ ;
- atoms  $f[0] = v$  for every simple fluent constant  $f$  such that  $s_0[0] \models_{\sigma^f[0]} f[0] = v$ , i.e.,  $Simple(s_0[0])$ ;
- formulas  $F[i+1]$  for every fluent dynamic law  $F$  if  $G$  after  $\psi$  in  $D$  such that  $s_{i+1}[i+1] \models_{\sigma^f[i+1]} G[i+1]$  and  $s_i[i] \cup \varepsilon_i[i] \models_{\sigma[\psi]} \psi[i]$  for  $0 \leq i < m$ , i.e.,  $F[i+1]$  for every fluent dynamic law  $F$  if  $G$  after  $\psi$  in  $D$  such that  $s_{i+1} \models G$  and  $s_i \cup \varepsilon_i \models \psi$  for  $0 \leq i < m$ ; following the previous notation, this is the union of sets  $E(s_i, \varepsilon_i, s_{i+1})[i+1]$  for  $0 \leq i < m$ ;
- formulas  $\alpha[i]$  for every action dynamic law  $\alpha$  if  $\psi$  in  $D$  such that  $s_i[i] \cup \varepsilon_i[i] \models_{\sigma[\psi]} \psi[i]$  for  $0 \leq i < m$ , i.e.,  $\alpha[i]$  for every action dynamic law  $\alpha$  if  $\psi$  in  $D$  such that  $s_i \cup \varepsilon_i \models \psi$  for  $0 \leq i < m$ ; this is the union of sets  $A(\varepsilon_i, s_i)[i]$  for  $0 \leq i < m$ ;
- atoms  $\alpha[i] = v$  for every exogenous action atom in  $\sigma^a$  such that  $\varepsilon_i[i] \models_{\sigma^a[i]} \alpha[i] = v$  for  $0 \leq i < m$ , i.e., the union of sets  $Exog(\varepsilon_i[i])$  for  $0 \leq i < m$ .

Let  $X'$  be the interpretation  $s'[0] \cup e'[0] \cup \dots \cup e'[m-1] \cup s'[m]$  of  $\sigma_m$ .

$X' \models_{\sigma_m} (\Gamma_m^D)^X$  iff

$$\begin{aligned}
s'[0] &\models_{\sigma^f[0]} T_{static}(s_0)[0] \cup Simple(s_0[0]) \\
e'[0] &\models_{\sigma^a[0]} A(\varepsilon_0, s_0)[0] \cup Exog(\varepsilon_0[0]) \\
s'[1] &\models_{\sigma^f[1]} T_{static}(s_1)[1] \cup E(s_0, \varepsilon_0, s_1)[1] \\
&\vdots \\
s'[i] &\models_{\sigma^f[i]} T_{static}(s_i)[i] \cup E(s_{i-1}, \varepsilon_{i-1}, s_i)[i] \\
e'[i] &\models_{\sigma^a[i]} A(\varepsilon_i, s_i)[i] \cup Exog(\varepsilon_i[i]) \\
s'[i+1] &\models_{\sigma^f[i+1]} T_{static}(s_{i+1})[i+1] \cup E(s_i, \varepsilon_i, s_{i+1})[i+1] \\
&\vdots \\
e'[m-1] &\models_{\sigma^a[m-1]} A(\varepsilon_{m-1}, s_{m-1})[m-1] \cup Exog(\varepsilon_{m-1}[m-1]) \\
s'[m] &\models_{\sigma^f[m]} T_{static}(s_m)[m] \cup E(s_{m-1}, \varepsilon_{m-1}, s_m)[m]
\end{aligned}$$

It follows that  $X' \models_{\sigma_m} (\Gamma_m^D)^X$  iff

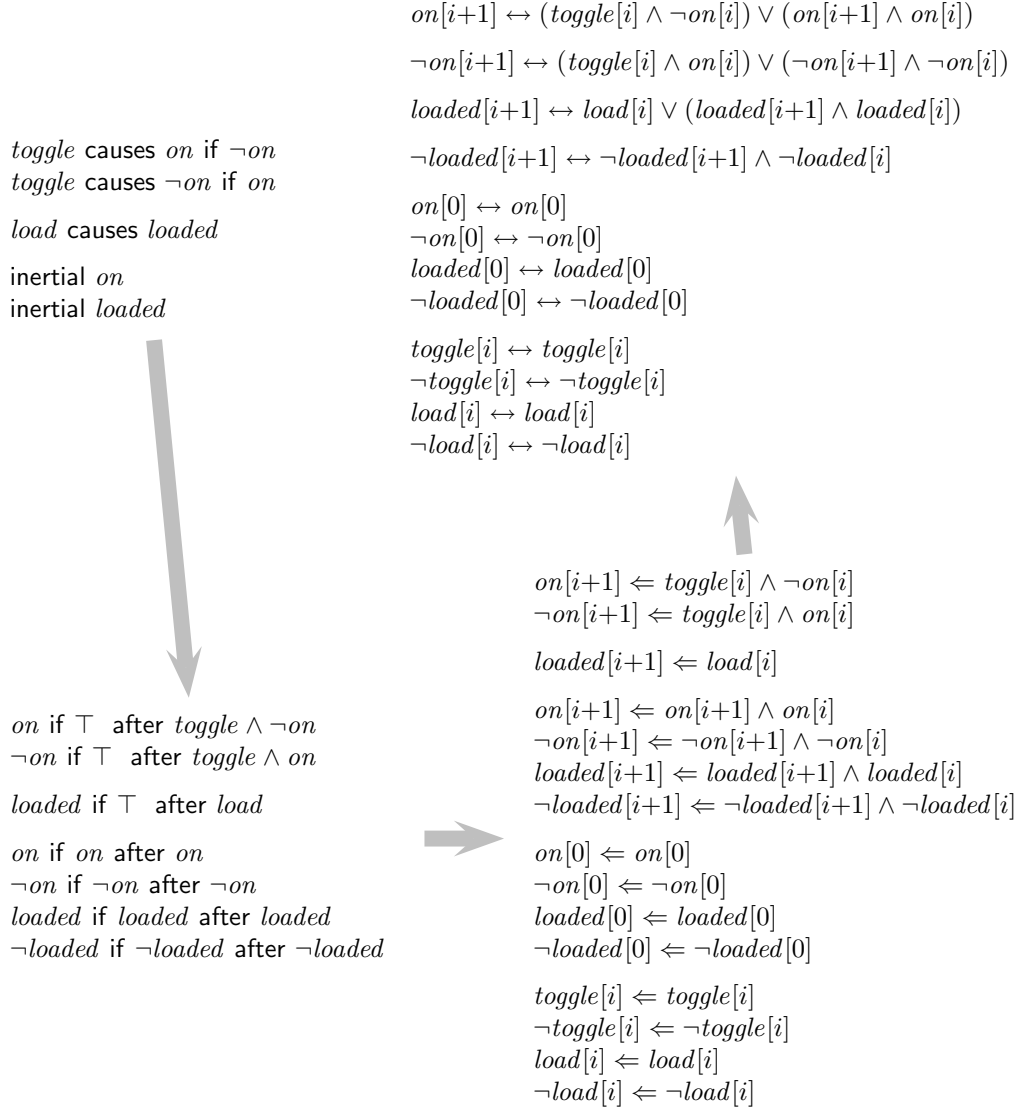
$$\begin{aligned}
s'_0 &\models T_{static}(s_0) \cup Simple(s_0) \\
\varepsilon_0 &\models A(\varepsilon_0, s_0) \cup Exog(\varepsilon_0) \\
s'_1 &\models T_{static}(s_1) \cup E(s_0, \varepsilon_0, s_1) \\
&\vdots \\
s'_i &\models T_{static}(s_i) \cup E(s_{i-1}, \varepsilon_{i-1}, s_i) \\
\varepsilon_i &\models A(\varepsilon_i, s_i) \cup Exog(\varepsilon_i) \\
s'_{i+1} &\models T_{static}(s_{i+1}) \cup E(s_i, \varepsilon_i, s_{i+1}) \\
&\vdots \\
\varepsilon_{m-1} &\models A(\varepsilon_{m-1}, s_{m-1}) \cup Exog(\varepsilon_{m-1}) \\
s'_m &\models T_{static}(s_m) \cup E(s_{m-1}, \varepsilon_{m-1}, s_m)
\end{aligned}$$

And  $X$  is the only model of  $(\Gamma_m^D)^X$  iff

$$\begin{aligned}
s'_0 &\not\models T_{static}(s_0) \cup Simple(s_0) \\
\varepsilon_0 &\not\models A(\varepsilon_0, s_0) \cup Exog(\varepsilon_0) \\
s'_1 &\not\models T_{static}(s_1) \cup E(s_0, \varepsilon_0, s_1) \\
&\vdots \\
s'_i &\not\models T_{static}(s_i) \cup E(s_{i-1}, \varepsilon_{i-1}, s_i) \\
\varepsilon_i &\not\models A(\varepsilon_i, s_i) \cup Exog(\varepsilon_i) \\
s'_{i+1} &\not\models T_{static}(s_{i+1}) \cup E(s_i, \varepsilon_i, s_{i+1}) \\
&\vdots \\
\varepsilon_{m-1} &\not\models A(\varepsilon_{m-1}, s_{m-1}) \cup Exog(\varepsilon_{m-1}) \\
s'_m &\not\models T_{static}(s_m) \cup E(s_{m-1}, \varepsilon_{m-1}, s_m)
\end{aligned}$$

But by Proposition 5, these are precisely the conditions for  $s_0 \varepsilon_0 \cdots \varepsilon_{m-1} s_m$  to be a path of length  $m$  of  $D$ .  $\square$

## 6.2 Example



Note the introduction of the ‘exogeneity laws’ in the second step.

## 7 Computational tasks

$\mathcal{C}+$  is a language for defining labelled transition systems (of a certain kind). A wide variety of languages—‘query languages’—can be interpreted on the transition systems so defined. These include simple propositional languages, as well as temporal logics such as CTL and LTL, and potentially many others. The query language supported by the Causal Calculator CCALC is the time-stamped language discussed in Section 3.3 for expressing properties of paths/runs of some

specified length  $m$  of a transition system. For an action description of signature  $(\sigma^f, \sigma^a)$  and non-negative integer  $m$ , this is the propositional language of signature  $\sigma_m^f \cup \sigma_{m-1}^a$ , which is also the signature of the causal theory  $\Gamma_m^D$ .

For example, for the ‘lottery example’ of Section 4.7, the query

$$\neg \text{alive}[0] \wedge \text{alive}[m] \wedge \text{happy}[m]$$

asks whether there is a path of length  $m$  in the transition system such that *alive* is false in the initial state while *alive* and *happy* are both true in state  $m$ .

For a definite action description  $D$ , a formula  $\psi$  of the query language can be evaluated by finding (classical) models of  $\text{comp}(\Gamma_m^D)$  (and therefore paths of length  $m$  of the transition system defined by  $D$ ) which are also models of  $\psi$ , i.e., finding models of  $\text{comp}(\Gamma_m^D) \cup \{\psi\}$ .

This then is the basic operation of the Causal Calculator CCALC. For a given action description  $D$  of signature  $(\sigma^f, \sigma^a)$ , non-negative integer  $m$ , and query  $\psi$  of the time-stamped signature  $\sigma_m$ , CCALC

- performs the translation of  $D$  to  $\Gamma_m^D$ ,
- constructs  $\text{comp}(\Gamma_m^D)$ ,
- invokes a standard propositional sat-solver to find (classical) models of  $\text{comp}(\Gamma_m^D) \cup \psi$ , and then
- post-processes the sat-solver output to show the models obtained.

In practice, the first two steps may be combined into one, possibly with some additional optimisations to simplify the set of formulas passed to the sat-solver.

There is also a sub-language for specifying the action signature—sorts, variables, shorthand for various common forms of fluent, implementations of integer arithmetic, and so on—and a language for expressing particular common forms of queries. Details of these components are omitted here. They are not always well documented and vary slightly from one version of CCALC to another. See <http://www.cs.utexas.edu/users/tag/cc> for details of the current version of CCALC.

This general method covers a wide range of computational tasks. The account below follows a standard classification (prediction, ‘postdiction’, planning) current in the AI literature. CCALC has been applied successfully to a number of non-trivial benchmark examples in the temporal reasoning literature (see e.g. (Akman et al., 2004) and the CCALC website). In our own work we have used it to construct executable specifications of agent societies (see e.g. (Artikis et al., 2003a,b)).

**Prediction** Given an action description  $D$  of signature  $(\sigma^f, \sigma^a)$ :

- Initially  $F$  holds.
- Partially specified events of type  $\alpha_0, \alpha_1, \dots, \alpha_k$  happen.
- Does  $G$  hold in state  $k + 1$ ?

In other words is there a path/run/history  $s_0 \varepsilon_0 s_1 \cdots s_k \varepsilon_k s_{k+1}$  such that  $s_0 \models F$ ,  $\varepsilon_i \models \alpha_i$  for each  $i \in 0..k$ , and  $s_{k+1} \models G$ ?  $F$  and  $G$  are formulas of  $\sigma^f$  and  $\alpha_i$  are formulas of  $\sigma^a$ .



We want to know whether

$$\text{comp}(\Gamma_{k+1}^D) \models (F[0] \wedge \alpha_0[0] \wedge \alpha_1[1] \wedge \cdots \wedge \alpha_k[k] \rightarrow G[k+1])$$

We use a sat-solver to check whether

$$\text{comp}(\Gamma_{k+1}^D) \cup \{F[0] \wedge \alpha_0[0] \wedge \alpha_1[1] \wedge \cdots \wedge \alpha_k[k] \wedge \neg G[k+1]\}$$

is satisfiable.

A variant of the problem:

- Initially  $F$  holds.
- Partially specified events of type  $\alpha_0, \alpha_1, \dots, \alpha_k$  happen.
- Is it *possible* that  $G$  holds in state  $k+1$ ? In other words, is there a possible run/path/history through the transition system such that  $G$  holds at its final state?

We check whether

$$\text{comp}(\Gamma_{k+1}^D) \cup \{F[0] \wedge \alpha_0[0] \wedge \alpha_1[1] \wedge \cdots \wedge \alpha_k[k] \wedge \neg G[k+1]\}$$

is satisfiable. If satisfiable, a propositional sat-solver will return all models.

**‘Postdiction’** (stupid term)

- Partially specified events of type  $\alpha_0, \alpha_1, \dots, \alpha_k$  happen.
- $G$  holds now.
- Does it follow that initially  $F$ ?

We want to know whether

$$\text{comp}(\Gamma_{k+1}^D) \models (\alpha_0[0] \wedge \alpha_1[1] \wedge \cdots \wedge \alpha_k[k] \wedge G[k+1] \rightarrow F[0])$$

We check whether

$$\text{comp}(\Gamma_{k+1}^D) \cup \{\alpha_0[0] \wedge \alpha_1[1] \wedge \cdots \wedge \alpha_k[k] \wedge G[k+1] \wedge \neg F[0]\}$$

is satisfiable. And as before, checking whether

$$\text{comp}(\Gamma_{k+1}^D) \cup \{\alpha_0[0] \wedge \alpha_1[1] \wedge \cdots \wedge \alpha_k[k] \wedge G[k+1] \wedge F[0]\}$$

is satisfiable deals with the variant of the problem in which we want to know whether it is *possible* that initially  $F$ .

**Temporal interpolation** Prediction and ‘postdiction’ are both special cases of the general problem in which:

- Partially specified events of type  $\alpha_0, \alpha_1, \dots, \alpha_k$  happen.
- Certain combinations of fluents (partially specified states) hold at given times.

We want to determine what holds in each state, or what *possibly* holds in each state.

## Planning

- Initially  $F$ .
- Goal:  $G$ .

Find the shortest sequence of fully specified actions (i.e., events, or transition labels)  $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{k-1}$  such that there is a path/run/history  $s_0 \varepsilon_0 s_1 \dots s_{k-1} \varepsilon_{k-1} s_k$  in which  $s_0 \models F$  and  $s_k \models G$ .

We try *consecutively* for  $k = 0, 1, \dots$  up to some specified maximum value  $m$ :

$$\text{comp}(\Gamma_k^D) \cup \{F[0] \wedge G[k]\} \quad \text{satisfiable?}$$

The sat-solver returns all models, and these contain a representation of the plan:  $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{k-1}$ .

(The reference here is to events  $\varepsilon_i$  because we want fully specified actions in the plan. In the other problems, a formula  $\alpha_i$  represents a partially specified event/transition.)

Notice that this formulation of planning is rather simplistic, in several respects. If any of the events in the plan  $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{k-1}$  are non-deterministic, execution of the sequence of events  $\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{k-1}$  in the initial state may *possibly* result in the goal state but there is no guarantee that it will. Moreover, there is no distinction here between events that are effectively realisable and those that or not. For example, a possible ‘plan’ for becoming rich in the lottery example is to win the lottery. Possible refinements to deal with practical planning problems are outside the scope of this document.

## Other possible problems

- Given a sequence of (partially specified) events  $\alpha_0 \alpha_1 \dots \alpha_k$  (no gaps), is this consistent with a given transition system (action description)  $D$ ? This can be combined with partial information about these actions, and about some or all of the states. This is an instance of the temporal interpolation problem above.
- Given a sequence of (partially specified) events  $\alpha_0 \alpha_1 \dots \alpha_k$ , but with possible gaps, is *this* consistent with a given transition system (action description)  $D$ ? What are the complete (no gap) sequences of events?

In principle this problem could also be solved, by iterating various sat-solver tests as in the planning problem to cover in some systematic fashion single gaps of one missing event, single gaps of  $j$  missing events, multiple gaps of single and  $j$  missing events, and so on, and so on, in all combinations. Clearly this method is not feasible without some further restrictions on possible gaps because of the (infinitely many) combinations to be considered.

## Model checkers for temporal logic

A main attraction of the  $\mathcal{C}+$  formalism compared to other action languages in the AI literature is that it has an explicit semantics in terms of transition systems, providing a bridge between AI formalisms and methods in other areas of computer science. One way of exploiting this link is by applying model checkers

for temporal logics to verify system properties of transition systems defined using the language  $\mathcal{C}+$  (and its extensions  $(\mathcal{C}+)^{++}$  to be presented later).

A small example of the use of CTL with an action description of  $\mathcal{C}+$  (actually  $(\mathcal{C}+)^{++}$  to be introduced in later sections) is sketched in (Sergot, 2005). Further development, both of  $\mathcal{C}+$  as an input language to standard temporal logic model checkers, and of the use of a modified form of CCALC as a bounded model checker, is a topic of current investigation and will be covered in a separate paper.

## Translation into (extended) logic programs

An alternative implementation route for  $\mathcal{C}+$  is provided via a translation of (a subclass of) causal theories into (extended) logic programs (Giunchiglia et al., 2004, Section 7.2). We will not pursue this line of development here. The following summary is included for completeness.

Consider a causal theory whose rules are all of the form

$$L_0 \Leftarrow L_1 \wedge \cdots \wedge L_n \quad (4)$$

where  $L_0, \dots, L_n (n \geq 0)$  are literals (i.e., of the form  $A$  or  $\neg A$  for  $A$  an atom of the signature of the theory).

Each such causal rule is translated to a clause of an (extended) logic program:

$$L_0 \leftarrow \text{not } \bar{L}_1, \dots, \text{not } \bar{L}_n$$

where *not* is negation by failure and  $\bar{L}_i$  stands for the literal complementary to  $L_i$ .

For  $\mathcal{C}+$ , a static law  $F$  if  $G$  thus translates to the logic program

$$F[i] \leftarrow \text{not } \neg G[i],$$

a law inertial  $p$  for Boolean constant  $p$  translates to the clauses

$$p[i+1] \leftarrow \text{not } \neg p[i+1], \text{not } \neg p[i] \quad (5)$$

$$\neg p[i+1] \leftarrow \text{not } p[i+1], \text{not } p[i], \quad (6)$$

and a law  $\alpha$  causes  $F$  if  $\psi$  translates to the clause

$$F[i+1] \leftarrow \text{not } \neg \alpha[i], \text{not } \neg \psi[i].$$

Every causal theory of the form (4) can be translated to an equivalent (extended) logic program, in the following sense.

**Proposition** (Proposition 11 in (Giunchiglia et al., 2004)) *Let  $\Gamma$  be a causal theory whose rules have the form*

$$L_0 \Leftarrow L_1 \wedge \cdots \wedge L_n \quad (n \geq 0)$$

where  $L_0, \dots, L_n$  are literals. Let  $lp(\Gamma)$  denote the (extended) logic program obtained by writing each such causal rule as a clause of the form:

$$L_0 \leftarrow \text{not } \bar{L}_1, \dots, \text{not } \bar{L}_n$$

An interpretation  $X$  of the signature of  $\Gamma$  is a model of the causal theory  $\Gamma$  iff  $X$  is an answer set for the logic program  $lp(\Gamma)$ .

Lifschitz and Turner (1999) present other, equivalent but more direct, translations of action descriptions into (extended) logic programs. For certain (common) special forms of action description, the translations of  $\mathcal{C}+$  shown above can be simplified. A static law  $F$  if  $G$  becomes

$$F[i] \leftarrow G[i]$$

a law inertial  $p$  for Boolean constant  $p$  translates to the clauses

$$p[i+1] \leftarrow \text{not } \neg p[i+1], p[i] \tag{7}$$

$$\neg p[i+1] \leftarrow \text{not } p[i+1], \neg p[i], \tag{8}$$

and a law  $\alpha$  causes  $F$  if  $\psi$  translates to the clause

$$F[i+1] \leftarrow \alpha[i], \psi[i].$$

The translation into logic programs provides another implementation route, employing answer set solvers such as SMODELs (Niemelä, 1999) and DLV (Eiter et al., 2001) in place of the satisfiability solvers used in CCALC. We shall not pursue this line of development.

We have a separate development, not presented here, that is attempting to develop an alternative translation to a different form of extended logic program with the same answer sets but with a different, ‘event calculus’ style, of computation (Craven and Sergot, 2003).

## 8 The language $(\mathcal{C}+)^+$ : Action generation and ‘counts as’

We now extend the language  $\mathcal{C}+$  by adding expressions of the form  $\alpha$  counts<sub>as</sub>  $\beta$  for representing (a version of) the ‘counts as’ relation of (Jones and Sergot, 1996) or, essentially the same idea, what Goldman (1970) referred to as ‘conventional generation’. We will call this extended language  $(\mathcal{C}+)^+$ . The counts<sub>as</sub> atoms of  $(\mathcal{C}+)^+$  can be seen as a distinguished type of fluent atom which has a special, fixed interpretation in an augmented kind of transition system. We want to be able to say that every transition of type  $\alpha$  counts also (in specified circumstances) as a transition of type  $\beta$ . (Or more precisely: every transition of the type denoted by the action formula  $\alpha$  counts (in specified circumstances) as a transition of the type denoted by action formula  $\beta$ .) The effects of such a transition are those that it has by virtue of being a transition of type  $\alpha$  and those it has by virtue of being a transition of type  $\beta$ . This in turn provides us with a treatment of institutionalised power, as outlined in Section 2.

One simple possibility, which we shall reject, is to leave the notion of transition system itself unchanged, and simply require that the transition labels defined by a  $(\mathcal{C}+)^+$  action description are closed under the generation/counts as rules. This is unsatisfactory, however. First, we would not be able to distinguish between  $\alpha$  counts<sub>as</sub>  $\beta$  and nonexecutable  $\alpha \wedge \neg\beta$  ( $\perp$  after  $\alpha \wedge \neg\beta$ ). Second,

we want ‘counts as’ to be a semantic feature of a transition system and not just a syntactic feature of the  $(\mathcal{C}+)^+$  language. And third, we want a *dynamic* notion of ‘counts as’ that can change under the effects of actions, so that we can say in the  $(\mathcal{C}+)^+$  language, for example, that  $\alpha'$  **causes**  $(\alpha \text{ counts\_as } \beta)$ , or that  $\alpha'$  **causes**  $\text{pow}(X, \alpha)$  where  $\text{pow}(X, \alpha)$  is an abbreviation for special forms of **counts\_as** expressions. So we are left with the task of augmenting transition systems in one way or another.

We will present the development in two stages: first, a simple form of augmented transition system with a *static* notion of ‘counts as’ to introduce the basic idea, and then a more complicated structure to support the dynamic notion. We will then discuss how action descriptions of  $(\mathcal{C}+)^+$  can be translated to  $\mathcal{C}+$  and causal theories.

A more general version of the whole account, allowing arbitrary formulas of signature  $\sigma^f \cup \sigma^a$  in **counts\_as** expressions and not just action formulas, could also be constructed along similar lines but will not be explored here.

## 8.1 Preliminaries

We want to be able to say that, under specified conditions, every transition  $(s, \varepsilon, s')$  of type  $X$  counts, in institution  $\mathcal{I}$ , as a transition of type  $Y$ . First we must decide on what to take as the transition types. One obvious candidate is the subsets  $\wp(R)$  of the set  $R$  of transitions. We would then define a counts as relation  $C_{\mathcal{I}} \subseteq \wp(R) \times \wp(R)$ . However, this is much more general than we need. It would allow, for instance,  $(s_1, \varepsilon_1, s'_1)$  counts as  $(s_2, \varepsilon_2, s'_2)$  even when  $s_1 \neq s'_1$  and/or  $s_2 \neq s'_2$ . This is meaningful, but is much more general than we need.

Instead we take as types of transition the subsets  $\wp(\mathbf{A})$  of the set  $\mathbf{A}$  of transition labels. The transition  $(s, \varepsilon, s')$  is of type  $X$  ( $X \subseteq \mathbf{A}$ ) when  $\varepsilon \in X$ . Since we want to use action formulas to denote types of transitions, we again consider the special case where transition labels are interpretations of  $\sigma^a$  (as opposed to specifying a separate valuation function for action constants of  $\sigma^a$ ). For action formula  $\alpha$  and a given transition system, let  $\|\alpha\|$  denote the set of transition labels/events that satisfy  $\alpha$ :

$$\|\alpha\| =_{\text{def}} \{\varepsilon \in \mathbf{I}(\sigma^a) \mid \varepsilon \models \alpha\}$$

As in previous sections, we will say that a transition is of type  $\alpha$  when we mean that the transition is of type  $\|\alpha\|$ . Notice that if action formulas  $\alpha$  and  $\beta$  are equivalent in the signature  $\sigma^a$  ( $\alpha \equiv_{\sigma^a} \beta$ ) then  $\|\alpha\| = \|\beta\|$ .

Now we add to the labelled transition system an extra component

$$C_{\mathcal{I}} \subseteq S \times \wp(\mathbf{A}) \times \wp(\mathbf{A})$$

$C_{\mathcal{I}}(s, X, Y)$  represents that any transition  $(s, e, s')$  of type  $X$  counts, in institution  $\mathcal{I}$ , as a transition of type  $Y$ . In principle there could be several such  $C_{\mathcal{I}}$  components, one for each institution  $\mathcal{I}$  that we wish to model. However, there are some further points of detail to be settled concerning the structure of states when there are multiple institutions, and to avoid those complications we stick to the case of a single institution for now (i.e., a single ‘counts as’ relation  $C_{\mathcal{I}}$ ).

One might also consider a four-argument ‘counts as’ relation

$$C_{\mathcal{I}} \subseteq S \times \wp(\mathbf{A}) \times \wp(\mathbf{A}) \times S$$

but this is also more general than we need. It would only make a difference in the case of non-deterministic transitions  $\varepsilon$  where  $X$  counts as  $Y$  depends on the final state as well as the initial state of the transition. This is so rare it is not worth the extra complications.

Now the basic semantic structure, which we will call a ‘statically augmented transition system’, or ‘static ATS’ for short, is a structure of the form

$$\langle S, \mathbf{A}, R, C_{\mathcal{I}} \rangle$$

where  $S$ ,  $\mathbf{A}$ , and  $R$  are the states, transition labels/events, and transitions of a labelled transition system, as usual, and  $C_{\mathcal{I}}$  is the ‘counts as’ relation,  $C_{\mathcal{I}} \subseteq S \times \wp(\mathbf{A}) \times \wp(\mathbf{A})$ .

We further require that  $C_{\mathcal{I}}$  satisfies the following condition, for all states  $s \in S$  and all transition types  $X, Y \subseteq \mathbf{A}$ :

$$\text{if } C_{\mathcal{I}}(s, X, Y) \text{ then for any } (s, \varepsilon, s') \in R, \text{ if } \varepsilon \in X \text{ then } \varepsilon \in Y \quad (C'_\varepsilon)$$

Condition  $(C'_\varepsilon)$  captures the key idea that any transition of type  $X$  is also a transition of type  $Y$ . It can also be expressed in the form

$$\text{if } C_{\mathcal{I}}(s, X, Y) \text{ then } \tau_R(s, X) \subseteq \tau_R(s, Y)$$

where  $\tau_R(s, X) =_{\text{def}} \{(s, \varepsilon, s') \in R \mid \varepsilon \in X\}$  denotes the set of transitions of type  $X$  that are executable in state  $s$ . Note that we do *not* want the converse of  $(C'_\varepsilon)$ : even if  $\varepsilon \in X$  implies  $\varepsilon \in Y$  for all  $(s, \varepsilon, s') \in R$  we do not necessarily want to conclude that  $C_{\mathcal{I}}(s, X, Y)$ . There may be other rules operative in institution  $\mathcal{I}$  (such as mere coincidence, and other kinds of causal and logical connections) which mean that all transitions of type  $X$  are also transitions of type  $Y$ , and we do not want to include these other connections as instances of the ‘counts as’ variety (Cf. (Jones and Sergot, 1996)). We will make some further comments about condition  $(C'_\varepsilon)$  in Section 8.10.

Note that transition types and the counts as relation  $C_{\mathcal{I}}$  are defined here in terms of states and transition labels/events, and are independent of the actual transition relation  $R$ .

The intention is this. We have some query language  $\mathcal{L}_Q$  of signature  $\sigma^f \cup \sigma^a$ , or possibly its time-stamped version, of signature  $\sigma^f[0] \cup \sigma^a[0] \cup \dots \cup \sigma^a[m-1] \cup \sigma^f[m]$  for some non-negative integer  $m$ . This language will have an additional binary connective, say `counts_as $\mathcal{I}$` , allowing expressions of the form  $\alpha \text{ counts\_as}_{\mathcal{I}} \beta$  for  $\alpha$  and  $\beta$  belonging to some designated set of action formulas, or  $\alpha \text{ counts\_as}_{\mathcal{I}}[i] \beta$  in the time-stamped language. This `counts_as $\mathcal{I}$`  connective is intended to be the analogue of the ‘counts as’ conditional of (Jones and Sergot, 1996) but restricted here to expressing ‘counts as’ relations between (designated pairs of) action formulas rather than formulas in general. Atoms of the language  $\mathcal{L}_Q$  will be interpreted on states and transition labels as usual. `counts_as $\mathcal{I}$`  expressions will be evaluated on states as follows:

$$s \models \alpha \text{ counts\_as}_{\mathcal{I}} \beta \quad \text{iff} \quad (s, \|\alpha\|, \|\beta\|) \in C_{\mathcal{I}}$$

or if  $\mathcal{L}_Q$  is a time-stamped language evaluated on paths, then

$$\pi \models \alpha \text{ counts\_as}_{\mathcal{I}}[i] \beta \quad \text{iff} \quad (\pi_s[i], \|\alpha\|, \|\beta\|) \in C_{\mathcal{I}}$$

where  $\pi_s[i]$  denotes the  $i$ th state of the path  $\pi$ , i.e.,  $\pi_s[i] = s_i$  when  $\pi = s_0 \varepsilon_0 s_1 \varepsilon_1 \dots \varepsilon_{i-1} s_i \dots$ .

(We are using different symbols  $\text{counts\_as}_{\mathcal{I}}$  and  $\text{counts\_as}$  so as not to confuse expressions of the query language with expressions of the  $(\mathcal{C}+)^+$  action descriptions used to define an ATS. Informally, they are intended to represent the same concept. It will be convenient to keep the distinction when we discuss implementation options and translations of  $(\mathcal{C}+)^+$  to causal theories in Section 8.9 below.)

**Augmented action signatures** We will again consider the special case in which states are interpretations of some set  $\sigma^f$  of fluent constants and transition labels/events are the interpretations of a set  $\sigma^a$  of action constants. We add a (disjoint) set  $\sigma^\gamma$  of Boolean constants of the form  $\alpha \text{ counts\_as} \beta$  for  $\alpha$  and  $\beta$  belonging to some designated set of action formulas of  $\sigma^a$ . We do not necessarily want to allow arbitrary formulas of  $\alpha$  and  $\beta$  of  $\sigma^a$  in  $\alpha \text{ counts\_as} \beta$  expressions.  $\sigma^\gamma$  allows these restrictions to be specified if desired.  $(\sigma^f, \sigma^\gamma, \sigma^a)$  is an *augmented action signature*. We will say that constants  $\alpha \text{ counts\_as} \beta$  and  $\alpha' \text{ counts\_as} \beta'$  of  $\sigma^\gamma$  are  $\sigma^a$ -equivalent when  $\alpha \equiv_{\sigma^a} \alpha'$  and  $\beta \equiv_{\sigma^a} \beta'$ . Notice that since  $\top$  and  $\perp$  are not action formulas (they are 0-ary connectives) there are no constants of the form  $\alpha \text{ counts\_as} \perp$  or  $\top \text{ counts\_as} \beta$ . As for the query language  $\mathcal{L}_Q$ ,  $\alpha \text{ counts\_as}_{\mathcal{I}} \beta$  is a well-formed expression of  $\mathcal{L}_Q$  when  $\alpha \text{ counts\_as} \beta$  is a constant of  $\sigma^\gamma$  (and likewise for the time-stamped query language).

**Definition 10 (Static ATS)** A statically augmented transition system (static ATS) of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  is a structure of the form

$$\langle \sigma^f, \sigma^\gamma, S, I(\sigma^a), R, C_{\mathcal{I}} \rangle$$

where

- $S \subseteq I(\sigma^f)$  is a set of states,
- $I(\sigma^a)$  is the set of transition labels/events,
- $R$  is the set of labelled transitions,  $R \subseteq S \times I(\sigma^a) \times S$ ,
- $C_{\mathcal{I}}$  is the ‘counts as’ relation,  $C_{\mathcal{I}} \subseteq S \times \wp(I(\sigma^a)) \times \wp(I(\sigma^a))$ ,

and, for all states  $s \in S$  and all transition types  $X, Y \subseteq I(\sigma^a)$ :

$$\text{if } C_{\mathcal{I}}(s, X, Y) \text{ then for any } (s, \varepsilon, s') \in R, \text{ if } \varepsilon \in X \text{ then } \varepsilon \in Y \quad (C'_{\varepsilon})$$

**Remark (unsupported institutional actions)** When  $\alpha \text{ counts\_as} \beta$ ,  $\beta$  represents an institutional action (as opposed to a ‘brute’ or ‘basic’ action).  $\alpha$  may represent an institutional action or a brute/basic action. In the present version, we have chosen not to make any explicit distinction between brute/basic and institutional actions in the semantics. Further, we have chosen not to eliminate, in the semantics, the possibility of *unsupported* institutional actions.  $\beta$  is an unsupported institutional action in a transition  $(s, \varepsilon, s')$  when  $\alpha \text{ counts\_as} \beta$  is a

constant in  $\sigma^\gamma$  (so  $\beta$  represents an institutional action),  $\varepsilon \models \beta$ , but  $\varepsilon \not\models \alpha'$  for any constant  $\alpha'$  `counts_as`  $\beta$  in  $\sigma^\gamma$ . For example, if `wave_flag counts_as start_race` is a constant in  $\sigma^\gamma$ , and there is no other means of performing `start_race` besides `wave_flag` specified in  $\sigma^\gamma$ , then `start_race` is an unsupported institutional action in any transition of type `start_race`  $\wedge$   $\neg$ `wave_flag`. It is a moot point whether unsupported institutional actions are meaningful given the intended informal reading. We have chosen not to eliminate them in the semantics (though this could be done easily) because unsupported institutional actions can be eliminated straightforwardly as desired when the action description is formulated, as will be illustrated by an example later. Not building this into the semantics gives more flexibility. We return to this point later in Section 8.8.

The treatment of ‘counts as’ presented here is guided by the account of ‘counts as’ conditionals and institutionalised power in (Jones and Sergot, 1996). A direct comparison is not possible because the semantic frameworks of the two accounts are very different, and because the dynamic setting adopted here raises a number of further considerations outside the scope of (Jones and Sergot, 1996). We limit ourselves to some brief remarks and some comments on the related concept of conventional generation of Goldman (1970).

The semantic counterpart of  $C_{\mathcal{I}}$  in (Jones and Sergot, 1996) is transitive, at least in the absence of convincing counter-examples to the contrary. This suggests that we might strengthen the  $C_{\mathcal{I}}$  relation, for example by imposing the following additional constraint, for all transition types  $X, Y, Z \subseteq \mathbf{A}$  and all states  $s$  of the transition system:

$$\text{if } C_{\mathcal{I}}(s, X, Y) \text{ and } C_{\mathcal{I}}(s, Y, Z) \text{ then } C_{\mathcal{I}}(s, X, Z) \quad (C_{\text{trans}})$$

This would mean that all states  $s$  have the property that

$$s \models (\alpha \text{ counts\_as}_{\mathcal{I}} \beta \wedge \beta \text{ counts\_as}_{\mathcal{I}} \gamma) \rightarrow \alpha \text{ counts\_as}_{\mathcal{I}} \gamma$$

though only when  $\alpha \text{ counts\_as} \beta$ ,  $\beta \text{ counts\_as} \gamma$ , and  $\alpha \text{ counts\_as} \gamma$  are all constants of  $\sigma^\gamma$ .

The possibility that ‘counts as’ could also be reflexive—that is, in the present context, that the following constraint could also be imposed

$$C_{\mathcal{I}}(s, X, X)$$

for all states  $s$  and all transition types  $X \subseteq \mathbf{A}$ —is explicitly rejected in (Jones and Sergot, 1996). The reason is simply that ‘counts as’ is not a mere form of classification, but is intended to convey the idea that certain kinds of action, or states of affairs more generally, have a special conventional significance within an institution. The action of opening a window, say, might have no special conventional significance within an institution: we would then not wish to say that opening a window ‘counts as’ opening a window. A reflexivity constraint could be imposed, but only in a restricted form, for special categories of transition types  $X$ , those corresponding to actions of conventional significance in institution  $\mathcal{I}$ . How are these special categories of transition types are to be characterised? This is not a question we will attempt to address in this paper. For present purposes we can say that  $X$  is a transition type in this special category when  $X = \|\alpha\|$  or  $X = \|\beta\|$  for some constant  $\alpha \text{ counts\_as} \beta$  in  $\sigma^\gamma$ . For



ease of reference, let us write  $\mathbf{A}_{\mathcal{I}}$  for the transition types of special conventional significance in institution  $\mathcal{I}$ , as represented by the signature  $\sigma^\gamma$ :  $\mathbf{A}_{\mathcal{I}} =_{\text{def}} \{X \subseteq \mathbf{A} \mid X = \|\alpha\| \text{ or } X = \|\beta\| \text{ for some constant } \alpha \text{ counts\_as } \beta \text{ in } \sigma^\gamma\}$ . A suitably qualified reflexivity constraint would take the form

$$C_{\mathcal{I}}(s, X, X) \text{ if } X \subseteq \mathbf{A}_{\mathcal{I}} \quad (\text{C}_{\text{refl}})$$

Note that relying on the signature  $\sigma^\gamma$  in this way finesses the practical consequences of deciding what form of reflexivity constraint to adopt. The constraint  $(\text{C}_{\text{refl}})$  yields the property that all states  $s$  satisfy

$$s \models \alpha \text{ counts\_as}_{\mathcal{I}} \alpha$$

but only when  $\alpha \text{ counts\_as } \alpha$  is specified as a constant of  $\sigma^\gamma$ .

Jones and Sergot (1996) comment briefly on the similarities and some of the differences between their ‘counts as’ conditional and Goldman’s notion of conventional generation. Let  $G_{\mathcal{I}}(s, X, Y)$  represent that a transition of type  $X$  (under conditions  $s$ ) conventionally generates, in the norm/rule system (institution)  $\mathcal{I}$ , a transition of type  $Y$ . This move already deserves some comment<sup>2</sup>. Goldman’s act generation is a relationship between act *tokens* not act *types*. So, for example, a particular instance of Jim’s hammering a nail into a wall might ‘causally generate’ in Goldman’s terminology a particular instance of Jim’s waking Mary, but that would not be a relationship between all instances of Jim’s hammering a nail into a wall and all instances of Jim’s waking Mary. Conventional generation is different because norms/conventions tend to be general rather than specific, and so it is often meaningful to speak about conventional generation as a relationship between act *types*, as we do here.

Goldman argues that conventional generation, like other forms of act generation, is irreflexive and transitive. So let us impose the following constraints, for all states  $s$  in  $S$  and all subsets  $X, Y$  of  $\mathbf{A}$ :

- $(s, X, X) \notin G_{\mathcal{I}}$
- if  $G_{\mathcal{I}}(s, X, Y)$  and  $G_{\mathcal{I}}(s, Y, Z)$  then  $G_{\mathcal{I}}(s, X, Z)$

We see that  $G_{\mathcal{I}}(s, \cdot, \cdot)$  (or more precisely, the projection  $G_{\mathcal{I}}(s, \cdot, \cdot)$  for fixed  $s$ ) is a strict ordering on transition types. The natural interpretation is to see the relationship between ‘counts as’  $C_{\mathcal{I}}$  and Goldman’s ‘conventional generation’  $G_{\mathcal{I}}$  as analogous to that between a partial order and its corresponding strict partial order, i.e., to expect that the following should hold:  $C_{\mathcal{I}}(s, X, Y)$  iff either  $G_{\mathcal{I}}(s, X, Y)$  or  $X = Y$ , for all states  $s$  and all transition types  $X, Y \subseteq \mathbf{A}_{\mathcal{I}}$ .  $C_{\mathcal{I}}$  then comes out to be anti-symmetric, which suggests imposing the following further constraint on  $C_{\mathcal{I}}$ , for all states  $s$  in  $S$  and all subsets  $X, Y$  of  $\mathbf{A}_{\mathcal{I}}$ :

$$\text{if } C_{\mathcal{I}}(s, X, Y) \text{ and } C_{\mathcal{I}}(s, Y, X) \text{ then } X = Y \quad (\text{C}_{\text{a-symm}})$$

The ‘counts as’ relation  $C_{\mathcal{I}}$  (or rather, each projection  $C_{\mathcal{I}}(s, \cdot, \cdot)$  for fixed  $s$ ) is a partial ordering, not on the set of all transition types, but on the set  $\mathbf{A}_{\mathcal{I}}$  of transition types with special conventional significance in institution  $\mathcal{I}$ . The non-minimal elements of this ordering correspond to institutional actions, and

<sup>2</sup>The author is grateful to Andreas Herzig for a helpful reminder of this point.

the minimal ones to what we might call ‘basic’ actions. (We prefer the term ‘basic’ actions to Searle’s ‘brute’ actions because many examples, such as the auction protocol presented in Section 8.7 below, can be usefully formalised at the level of basic actions such as ‘signal a bid’ without necessarily specifying in detail the ‘brute’ actions (physical actions, the transmission of specific forms of electronic messages, and so on) by means of which the making of a bid is actually signalled in the auction.)

Although this seems quite plausible, we shall not explore these possibilities in detail here. Since ‘counts as’ and conventional generation appear inter-definable (at least in principle, even if not exactly as suggested above), we shall fix on the ‘counts as’ reading  $C_{\mathcal{I}}$  from now on. Moreover, we will not adopt the constraints  $(C_{\text{refl}})$ ,  $(C_{\text{trans}})$ ,  $(C_{\text{a-symm}})$  as fixed features of the semantics. Although it is tempting to do so, the practical significance of the reflexivity constraint  $(C_{\text{refl}})$  is negligible, while the transitivity constraint is problematic to express in the formalism of causal theories. There are moreover alternative formulations, particularly of transitivity, that could also be considered. We defer further discussion until Section 8.9 below.

Finally, given the intended reading of `counts_as` constants, it would be natural to impose the additional restriction that the ‘counts as’ signature  $\sigma^\gamma$  must be cycle-free, in the sense that it must contain no constants  $\alpha_0 \text{ counts\_as } \alpha_1, \alpha_1 \text{ counts\_as } \alpha_2, \dots, \alpha_{n-1} \text{ counts\_as } \alpha_n$  ( $n \geq 1$ ) where  $\alpha_n \equiv_{\sigma^a} \alpha_0$ , and  $\alpha_i \equiv_{\sigma^a} \alpha_{i+1}$  for  $0 \leq i < n$ . Again, this seems reasonable but there are some further points of detail to resolve, and we make no such assumption about  $\sigma^\gamma$  in the version presented in this paper. All the examples that will be discussed, however, will have  $\sigma^\gamma$  cycle-free in the sense described above.

## 8.2 Augmented transition systems

The ‘static ATS’ structure defined in the previous section is too simple for what we want. It provides only a *static* ‘counts as’ relation between transition types, whereas we want to support a dynamic version, that is, one in which transitions can change the  $C_{\mathcal{I}}$  relation whether or not they also change the values of fluents. We need a more complicated structure. One possibility is to extend the static ATS structure with another transition relation, this time between ‘counts as’ relations: the tuple  $(C_{\mathcal{I}}, \varepsilon, C'_{\mathcal{I}})$  would represent that a transition with label  $\varepsilon \in \mathbf{A}$  transforms the ‘counts as’ relation from  $C_{\mathcal{I}}$  to  $C'_{\mathcal{I}}$ . We would then need to formulate suitable constraints to ensure that transitions between states and transitions between  $C_{\mathcal{I}}$  relations are coherent in some appropriate sense. Instead of that, we can get the same effect by defining a structure with a *single* transition relation, but one defined between *augmented states*  $(s_f, \gamma_{\mathcal{I}})$ , each of which is characterised by the values  $s_f$  of the fluent constants and by a set  $\gamma_{\mathcal{I}}$  of pairs  $(X, Y)$  specifying that in this augmented state a transition of type  $X$  counts as a transition of type  $Y$ . The condition corresponding to  $(C'_\varepsilon)$  is then:

$$\text{if } ((s_f, \gamma_{\mathcal{I}}), \varepsilon, (s'_f, \gamma'_{\mathcal{I}})) \in R, \text{ then for any } (X, Y) \in \gamma_{\mathcal{I}}, \text{ if } \varepsilon \in X \text{ then } \varepsilon \in Y$$

As usual, we will concentrate on the special case where states (now augmented states) are interpretations of the action signature. In an augmented action signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ ,  $\sigma^f$  and  $\sigma^\gamma$  are disjoint, and every augmented state  $s$  can

be represented by an interpretation of  $\sigma^f \cup \sigma^\gamma$ , i.e., in the form  $s_f \cup \gamma$  where  $s_f$  is an interpretation of the fluent constants  $\sigma^f$  and  $\gamma$  is an interpretation of the ‘counts as’ constants  $\sigma^\gamma$ . When interpretations are represented by the set of atoms that they satisfy,  $s_f$  is a set of fluent atoms and each  $\gamma$  will be a set of `counts_as` atoms. We need one further restriction: since  $\|\alpha\| = \|\beta\|$  when  $\alpha \equiv_{\sigma^a} \beta$ , the interpretation  $s$  must agree on the value of all  $\sigma^a$ -equivalent pairs of `counts_as` constants in  $\sigma^\gamma$ ;  $s$  (and therefore  $\gamma$ ) must satisfy the condition  $s \models (\alpha \text{ counts\_as } \beta \leftrightarrow \alpha' \text{ counts\_as } \beta')$  for all (distinct) pairs of constants  $\alpha \text{ counts\_as } \beta$ ,  $\alpha' \text{ counts\_as } \beta'$  in  $\sigma^\gamma$  such that  $\alpha \equiv_{\sigma^a} \alpha'$  and  $\beta \equiv_{\sigma^a} \beta'$ .

As a small point of detail, since an augmented state is characterised by the values of the  $\sigma^f$  and  $\sigma^\gamma$  constants, we can allow the augmented action signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  to have either, but not both, of  $\sigma^f$  and  $\sigma^\gamma$  empty. (We will use this in small illustrative examples to follow.)

**Definition 11** *Let  $(\sigma^f, \sigma^\gamma, \sigma^a)$  be an augmented action signature.  $s$  is an interpretation of the augmented action signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  when  $s$  is an interpretation of  $\sigma^f \cup \sigma^\gamma$  that agrees on the values of all  $\sigma^a$ -equivalent pairs of `counts_as` atoms in  $\sigma^\gamma$ . We will write  $I(\sigma^f \cup \sigma^\gamma; \sigma^a)$  for the set of all such interpretations. For convenience later, we write*

$$\text{Equiv}(\sigma^\gamma, \sigma^a) =_{\text{def}} \{ \alpha \text{ counts\_as } \beta \leftrightarrow \alpha' \text{ counts\_as } \beta' \mid \\ \alpha \text{ counts\_as } \beta \text{ and } \alpha' \text{ counts\_as } \beta' \text{ are distinct} \\ \text{constants of } \sigma^\gamma \text{ such that } \alpha \equiv_{\sigma^a} \alpha' \text{ and } \beta \equiv_{\sigma^a} \beta' \}$$

$s$  is then an interpretation of  $(\sigma^f, \sigma^\gamma, \sigma^a)$  when  $s$  is an interpretation of  $\sigma^f \cup \sigma^\gamma$  such that  $s \models \text{Equiv}(\sigma^\gamma, \sigma^a)$ .

Here is the special case of interest.

**Definition 12 (Augmented transition system (ATS))** *An augmented transition system (ATS) of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  is a structure of the form*

$$\langle \sigma^f, \sigma^\gamma, S, I(\sigma^a), R, C_{\mathcal{I}}^+ \rangle$$

where

- $S$  is the set of augmented states,  $S \subseteq I(\sigma^f \cup \sigma^\gamma; \sigma^a)$ , i.e., every  $s$  in  $S$  is an interpretation of  $\sigma^f \cup \sigma^\gamma$  that agrees on the values of all  $\sigma^a$ -equivalent pairs of `counts_as` atoms in  $\sigma^\gamma$ ,  $s \models \text{Equiv}(\sigma^\gamma, \sigma^a)$ ,
- $I(\sigma^a)$  is the set of transition labels/events,
- $R$  is the set of labelled transitions,  $R \subseteq S \times I(\sigma^a) \times S$ ,
- $C_{\mathcal{I}}^+$  is the dynamic ‘counts as’ relation,  $C_{\mathcal{I}}^+ \subseteq S \times \wp(I(\sigma^a)) \times \wp(I(\sigma^a))$ , defined as

$$C_{\mathcal{I}}^+ =_{\text{def}} \{ (s, \|\alpha\|, \|\beta\|) \mid s \in S, s \models \alpha \text{ counts\_as } \beta \}$$

- $R$  satisfies the following condition, for all states  $s \in S$  and all transition types  $X, Y \subseteq I(\sigma^a)$ :

$$\text{if } C_{\mathcal{I}}^+(s, X, Y) \text{ then for any } (s, \varepsilon, s') \in R, \text{ if } \varepsilon \in X \text{ then } \varepsilon \in Y \quad (C_{\varepsilon}^+)$$

There is clearly some redundancy in this structure, because the  $C_{\mathcal{I}}^+$  component is already determined by the augmented states, but keeping  $C_{\mathcal{I}}^+$  explicit is convenient for expressing constraints on an ATS, and for ease of comparison with the static ATS structure of the previous section. The condition  $(C_{\varepsilon}^+)$  is equivalently stated (see below) as:

if  $(s, \varepsilon, s') \in R$  then, for any  $s \models \alpha \text{ counts\_as } \beta$ , if  $\varepsilon \in \|\alpha\|$  then  $\varepsilon \in \|\beta\|$

or equivalently again

if  $(s, \varepsilon, s') \in R$  then, for any  $s \models \alpha \text{ counts\_as } \beta$ , we have  $\varepsilon \models \alpha \rightarrow \beta$

Now, by construction, we have  $(s, \|\alpha\|, \|\beta\|) \in C_{\mathcal{I}}$  iff  $s \models \alpha \text{ counts\_as } \beta$ , and so evaluation of  $\alpha \text{ counts\_as}_{\mathcal{I}} \beta$  expressions in a query language evaluated on the augmented transition system can be reduced to evaluation of  $\alpha \text{ counts\_as } \beta$  atoms on *states*. This is a very desirable feature.

Note again that we have chosen not to eliminate the possibility of unsupported institutional actions in the ATS but leave them to be eliminated as desired when action descriptions are formulated. We return to this point in Section 8.8.

**Remark** One can see that a static ATS  $\langle \sigma^f, \sigma^\gamma, S_f, \mathbf{A}, R_f, C_{\mathcal{I}} \rangle$  is the special case of an ATS  $\langle \sigma^f, \sigma^\gamma, S, \mathbf{A}, R, C_{\mathcal{I}}^+ \rangle$  in which there are no states  $s_f \cup \gamma$  and  $s_f \cup \gamma'$  in  $S$  with  $\gamma \neq \gamma'$ , and where  $S_f$  is the set of states  $\{s_f \mid s_f \cup \gamma \in S\}$ ,  $R_f = \{(s_f, \varepsilon, s'_f) \mid (s_f \cup \gamma, \varepsilon, s'_f \cup \gamma') \in R\}$ , and  $C_{\mathcal{I}} = \{(s_f, X, Y) \mid (s_f \cup \gamma, X, Y) \in C_{\mathcal{I}}^+\}$ . We note that a dynamic ATS can be approximated by a static ATS, with no real loss of expressive power. For instead of allowing actions to change the values of `counts_as` constants, we can always define `counts_as` in terms of ordinary fluent constants, possibly with the introduction of new fluent constants. We then specify how actions affect these (new) fluent constants rather than the `counts_as` constants directly. Since the values of `counts_as` constants are then determined completely by the values of ordinary fluent constants in each state, a static ‘counts as’ relation is sufficient. It may seem therefore that there is no real benefit in adopting the more complicated dynamically augmented structure. However, there is also no real benefit in restricting attention to the simpler static version. As an implementation issue, the causal theories into which  $(\mathcal{C}+)^+$  action descriptions are translated turn out to be essentially the same whether we choose the static or the dynamic version. Moreover, even in the static version, the ‘counts as’ relation  $C_{\mathcal{I}}$  has to be represented somehow. If we choose to represent it by a set of `counts_as` atoms, which seems natural, then we end up with augmented states again, exactly as in the dynamic version but now without the flexibility to manipulate `counts_as` constants directly. We save nothing, at the price of having to impose a restriction on the  $(\mathcal{C}+)^+$  language to prohibit formulas containing `counts_as` constants from appearing in the heads of fluent dynamic laws. We might as well support the dynamic notion, at no extra cost, and with the additional flexibility.

### 8.3 The language $(\mathcal{C}+)^+$

We want to be able to include ‘counts as’ expressions in both the heads and bodies of static and dynamic laws, and to write, for example, laws of the form  $\alpha'$  *causes*  $(\alpha \text{ counts\_as } \beta)$  if  $G$ . Syntactically, expressions of the form  $\alpha \text{ counts\_as } \beta$  will behave just like (Boolean) fluent constants. In any given application the set of action formulas  $\alpha$  and  $\beta$  that can appear in a *counts\_as* constant may be further restricted when the ‘counts as’ signature  $\sigma^\gamma$  is specified. In particular it is often sufficient to restrict  $\alpha \text{ counts\_as } \beta$  constants to the form where both  $\alpha$  and  $\beta$  are action *atoms*, or more generally conjunctions of action atoms, from some designated set. The general development does not depend on any such restriction however, so we will not adopt it yet, but defer further discussion until we consider the translation of  $(\mathcal{C}+)^+$  action descriptions to causal theories. Moreover, to provide flexibility in formulating action descriptions, it is convenient to retain the distinction between *simple* and *statically determined* constants in the ‘counts as’ signature  $\sigma^\gamma$ .

**Syntax** The syntax of static and dynamic laws of  $(\mathcal{C}+)^+$  is exactly as in the language  $\mathcal{C}+$ , including the restrictions that apply to statically determined constants, and including the abbreviations *inertial*, *causes*, etc., described earlier. Syntactically, an action description of  $(\mathcal{C}+)^+$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  is also an action description of  $\mathcal{C}+$  of signature  $(\sigma^f \cup \sigma^\gamma, \sigma^a)$ . Semantically, *counts\_as* constants are given a special status in  $(\mathcal{C}+)^+$  that they do not have in  $\mathcal{C}+$  where they are just regular fluent constants.

To avoid (unintended) unsupported institutional actions, any action constant in the formula  $\beta$  of an  $\alpha \text{ counts\_as } \beta$  atom is not (automatically) exogenous, though it may be deliberately declared as exogenous when the action description is formulated. (An example follows presently.)

**Semantics** An action description of  $(\mathcal{C}+)^+$  defines an augmented transition system: the *counts\_as* constants are given special treatment in that they define its ‘counts as’ relation, as follows.

**Definition 13** Let  $D$  be an action description of  $(\mathcal{C}+)^+$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ . The ATS defined by  $D$  is

$$\langle \sigma^f, \sigma^\gamma, S, I(\sigma^a), R, C_D^+ \rangle$$

where:

- $s \in S$  is a state of the ATS defined by  $D$  (in short: a  $(\mathcal{C}+)^+$ -defined state of  $D$ ) iff  $s$  is the only interpretation of  $\sigma^f \cup \sigma^\gamma$  that agrees on the values of all  $\sigma^a$ -equivalent pairs of *counts\_as* atoms in  $\sigma^\gamma$  and satisfies  $T_{\text{static}}(s) \cup \text{Simple}(s)$ , i.e., the only interpretation of  $\sigma^f \cup \sigma^\gamma$  such that<sup>3</sup>

$$s \models T_{\text{static}}(s) \cup \text{Simple}(s) \cup \text{Equiv}(\sigma^\gamma, \sigma^a)$$

---

<sup>3</sup>It is easy to check that  $s$  is the only interpretation of  $\sigma^f \cup \sigma^\gamma$  that agrees on the values of all  $\sigma^a$ -equivalent pairs of *counts\_as* atoms in  $\sigma^\gamma$  and satisfies a set of formulas  $X$  iff  $s \models X \cup \text{Equiv}(\sigma^\gamma, \sigma^a)$ .

- $(s, \varepsilon, s') \in R$  is a transition of the ATS defined by  $D$  (in short: a  $(\mathcal{C}+)^+$ -defined transition of  $D$ ) iff  $s$  and  $s'$  are interpretations of  $\sigma^f \cup \sigma^\gamma$  and  $\varepsilon$  is an interpretation of  $\sigma^a$  such that
  - $s \models T_{static}(s) \cup Simple(s) \cup Equiv(\sigma^\gamma, \sigma^a)$
  - $s' \models T_{static}(s') \cup E(s, \varepsilon, s') \cup Equiv(\sigma^\gamma, \sigma^a)$
  - $\varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon) \cup \{\alpha \rightarrow \beta \mid s \models \alpha \text{ counts\_as } \beta\}$
- The ‘counts as’ relation  $C_D^+$  of the ATS defined by  $D$  (in short: the ‘counts as’ relation defined by  $D$ ) is

$$C_D^+ = \{(s, \|\alpha\|, \|\beta\|) \mid s \models \alpha \text{ counts\_as } \beta\}$$

Compare in particular the conditions  $\varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon) \cup \{\alpha \rightarrow \beta \mid s \models \alpha \text{ counts\_as } \beta\}$  required for  $(s, \varepsilon, s')$  to be a transition of the  $(\mathcal{C}+)^+$  action description  $D$  with the weaker conditions  $\varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon)$  required for  $(s, \varepsilon, s')$  to be a transition of the  $\mathcal{C}+$  action description  $D$ .

It remains to check that the transition system defined by a  $(\mathcal{C}+)^+$  action description  $D$  is indeed an ATS of the correct form. We need to check that the transitions  $(\mathcal{C}+)^+$ -defined by  $D$  satisfy the condition  $(C_\varepsilon^+)$ , and that when  $(s, \varepsilon, s')$  is a transition  $(\mathcal{C}+)^+$ -defined by  $D$  then  $s'$  is an augmented state of  $D$ .

**Proposition 14** *Let  $D$  be an action description of  $(\mathcal{C}+)^+$ . The transition system  $\langle \sigma^f, \sigma^\gamma, S, I(\sigma^a), R, C_D^+ \rangle$  defined by  $D$  (i) satisfies the condition  $(C_\varepsilon^+)$  for an ATS, and (ii) has the property that when  $(s, \varepsilon, s') \in R$  then  $s' \in S$ .*

*Proof.* For (i): suppose  $(s, X, Y) \in C_D^+$ . Then  $X = \|\alpha\|$  and  $Y = \|\beta\|$  for some action formulas  $\alpha$  and  $\beta$  such that  $s \models \alpha \text{ counts\_as } \beta$ . Suppose  $(s, \varepsilon, s') \in R$ . Then by Definition 13, we have  $\varepsilon \models (\alpha \rightarrow \beta)$ , i.e.,  $\varepsilon \in \|\alpha\|$  implies  $\varepsilon \in \|\beta\|$ , and so  $\varepsilon \in X$  implies  $\varepsilon \in Y$ , as required.

For (ii), we need to show that if  $s'$  is the only model of  $T_{static}(s') \cup E(s, \varepsilon, s') \cup Equiv(\sigma^\gamma, \sigma^a)$  then  $s'$  is the only model of  $T_{static}(s') \cup Simple(s') \cup Equiv(\sigma^\gamma, \sigma^a)$ . This follows from the observation that  $E(s, \varepsilon, s')$  contains no statically determined constants, by exactly the same line of reasoning as in the proof of Proposition 4 for labelled transition systems.  $\square$

**Example (waving a flag)** Consider the  $(\mathcal{C}+)^+$  action description  $D$  with Boolean action constants *wave* and *start*, where *wave* (representing the waving of a flag, say) is exogenous and *start* (representing the institutional action of starting a race, say) is not.  $\sigma^\gamma$  contains the simple (Boolean) constant *wave counts\_as start*. Suppose  $D$  contains just the static laws

$$\begin{aligned} & \textit{wave counts\_as start} \text{ if } \textit{wave counts\_as start} \\ & \neg(\textit{wave counts\_as start}) \text{ if } \neg(\textit{wave counts\_as start}) \end{aligned}$$

(i.e., exogenous *wave counts\_as start*) and the action dynamic law

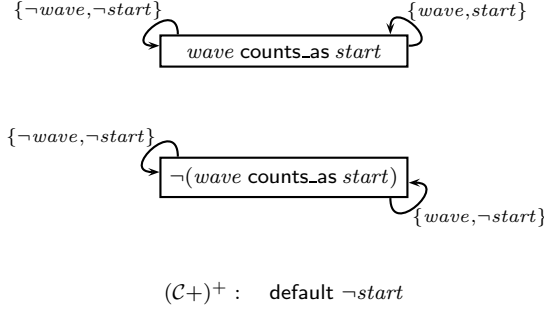
$$\neg\textit{start} \text{ if } \neg\textit{start} \quad (\text{i.e., default } \neg\textit{start})$$

There are no simple fluent constants in  $\sigma^f$ , no trivial constants in either  $\sigma^f$  or  $\sigma^\gamma$ , no pairs of distinct  $\sigma^a$ -equivalent *counts\_as* constants in  $\sigma^\gamma$ , and there are

only two static laws, so it is easy to confirm that there are precisely two states in the ATS defined by  $D$ :

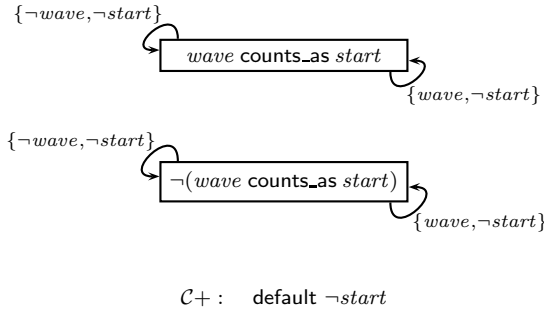
$$\{wave\ counts\_as\ start\} \quad \text{and} \quad \{\neg(wave\ counts\_as\ start)\}.$$

There are no trivial constants in  $\sigma^a$  and only one action dynamic law, and so the transitions in the ATS defined by  $D$  can be easily determined, as shown in the diagram below.

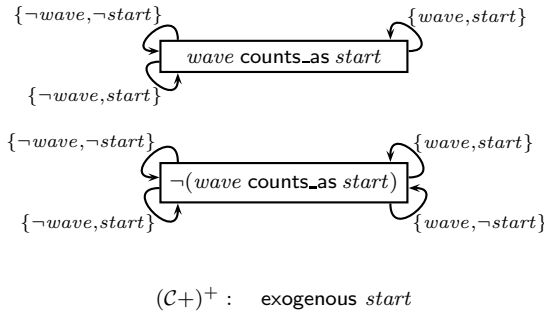


When  $wave\ counts\_as\ start$  holds, all transitions of type  $wave$  are also transitions of type  $start$ . Notice that  $\text{default } \neg start$  is strong enough (in this example) to eliminate unsupported institutional actions  $\{\neg wave, start\}$  even when  $wave\ counts\_as\ start$  does not hold.

For comparison, here is the transition system obtained when the same action description is viewed as an action description of  $C+$ .



We can admit unsupported institutional actions in the  $(C+)^+$ -defined ATS, if we want them, by making  $start$  exogenous, effectively adding the action dynamic law  $start$  if  $start$  to  $\text{default } \neg start$ . We then get



This is the transition system that would be obtained from the  $\mathcal{C}+$  action description, except that the  $\mathcal{C}+$ -defined system also contains a  $\{wave, \neg start\}$  transition from  $\{wave \text{ counts\_as } start\}$  to itself.

Note that the elimination of unsupported institutional actions in this example does not depend on the fact that there is only one means of starting a race. Suppose we extend the example with another action constant  $shoot$  (for shooting a gun, say) and another constant  $shoot \text{ counts\_as } start$  in the action signature. We add the static laws  $exogenous \ shoot \text{ counts\_as } start$ . There are four states in the ATS so defined. The state  $\{wave \text{ counts\_as } start, shoot \text{ counts\_as } start\}$  has (reflexive) transitions  $\{wave, shoot, start\}$ ,  $\{wave, \neg shoot, start\}$ ,  $\{\neg wave, shoot, start\}$ , and  $\{\neg wave, \neg shoot, \neg start\}$  but no transition  $\{\neg wave, \neg shoot, start\}$  where there is an unsupported institutional action. The state  $\{wave \text{ counts\_as } start, \neg(shoot \text{ counts\_as } start)\}$  has (reflexive) transitions  $\{wave, shoot, start\}$ ,  $\{wave, \neg shoot, start\}$ , and  $\{\neg wave, \neg shoot, \neg start\}$ , but no transition  $\{\neg wave, shoot, start\}$  and no transition  $\{\neg wave, \neg shoot, start\}$  with an unsupported institutional action. Similarly for  $\{\neg(wave \text{ counts\_as } start), shoot \text{ counts\_as } start\}$ . The state  $\{\neg(wave \text{ counts\_as } start), \neg(shoot \text{ counts\_as } start)\}$  has (reflexive) transitions  $\{wave, shoot, \neg start\}$ ,  $\{wave, \neg shoot, \neg start\}$ ,  $\{\neg wave, shoot, \neg start\}$ , and  $\{\neg wave, \neg shoot, \neg start\}$  but no transitions of type  $start$ , and in particular no transition  $\{\neg wave, \neg shoot, start\}$  with an unsupported institutional action.

We have eliminated unsupported institutional actions in this example by means of the action dynamic law **default**  $\neg start$ . This was chosen to emphasise the difference between  $(\mathcal{C}+)^+$  and  $\mathcal{C}+$  action descriptions. It is simpler and clearer to eliminate institutional actions by means of a law of the form

$$\text{nonexecutable } start \wedge \neg wave \wedge \neg shoot$$

This is the style we will adopt in subsequent examples.

It is now possible to investigate the relationships that exist in general between the states and transitions of a  $(\mathcal{C}+)^+$  action description of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  and those of the  $\mathcal{C}+$  action description of signature  $(\sigma^f \cup \sigma^\gamma, \sigma^a)$ . These relationships are somewhat involved however and we omit the details here.

## 8.4 Translation to $\mathcal{C}+$ and causal theories

Now: given an action description  $D$  of  $(\mathcal{C}+)^+$  we construct a causal theory  $\Gamma_m^+{}^D$  such that models of  $\Gamma_m^+{}^D$  correspond one-to-one with the paths of length  $m$  of the ATS defined by  $D$ . We do not need to worry about representing the ‘counts as’ relation  $C_D^+$  in models of  $\Gamma_m^+{}^D$  because **counts\_as** atoms in states do that for us already.

In fact, we can define a translation from a  $(\mathcal{C}+)^+$  action description of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  to a  $\mathcal{C}+$  action description  $D \cup \delta$  of signature  $(\sigma^f \cup \sigma^\gamma, \sigma^a)$  where the additional  $\mathcal{C}+$  laws  $\delta$  capture the special treatment given to **counts\_as** constants in  $(\mathcal{C}+)^+$ . Models of the causal theory  $\Gamma_m^{D \cup \delta}$  obtained from the  $\mathcal{C}+$  action description  $D \cup \delta$  correspond one-to-one with the paths of length  $m$  of the ATS defined by  $D$ ; in other words,  $\Gamma_m^+{}^D =_{\text{def}} \Gamma_m^{D \cup \delta}$ . Moreover (subject to some restrictions identified below) when  $D$  is a definite action description of



$(\mathcal{C}+)^+$ , the additional laws  $\delta$  are all definite laws of  $\mathcal{C}+$ , and so we obtain the very desirable property that definite action descriptions of  $(\mathcal{C}+)^+$  translate to definite action descriptions of  $\mathcal{C}+$ . This means that we can employ the method of literal completion, and the Causal Calculator CCALC, to perform computations on definite  $(\mathcal{C}+)^+$  action descriptions.

The additional causal laws  $\delta$  required to translate a  $(\mathcal{C}+)^+$  action description  $D$  to the equivalent  $\mathcal{C}+$  action description consist of two (disjoint) sets,  $\delta_D$  and  $\delta_{\sigma^\gamma}$ . The basic idea is that  $\delta_D$  contains a  $\mathcal{C}+$  action dynamic law

$$\beta \text{ if } \alpha \wedge \alpha \text{ counts\_as } \beta$$

for every constant  $\alpha \text{ counts\_as } \beta$  in  $\sigma^\gamma$ . However, in the case where  $\beta$  is a conjunction of (not necessarily atomic) action formulas it is much more convenient to include a separate action dynamic law for each conjunct of  $\beta$ , dealing separately with the special cases where the conjunct is  $\perp$  or  $\top$ . More precisely: every formula  $\beta$  in a constant  $\alpha \text{ counts\_as } \beta$  of  $\sigma^\gamma$  is a conjunction  $\beta_1 \wedge \dots \wedge \beta_m$  ( $m \geq 1$ ) of (not necessarily atomic) formulas  $\beta_i$ .  $\delta_D$  is the set of  $\mathcal{C}+$  laws obtained as follows. For every constant  $\alpha \text{ counts\_as } \beta$  in  $\sigma^\gamma$  such that  $\perp$  is not a conjunct of  $\beta$ ,  $\delta_D$  includes a  $\mathcal{C}+$  action dynamic law of the form

$$\beta_i \text{ if } \alpha \wedge \alpha \text{ counts\_as } \beta \tag{9}$$

for every conjunct  $\beta_i \neq \top$  of  $\beta$ ; and for every constant  $\alpha \text{ counts\_as } \beta$  in  $\sigma^\gamma$  such that  $\perp$  is a conjunct of  $\beta$ ,  $\delta_D$  includes a  $\mathcal{C}+$  fluent dynamic law of the form

$$\perp \text{ if } \top \text{ after } \alpha \wedge \alpha \text{ counts\_as } \beta \tag{10}$$

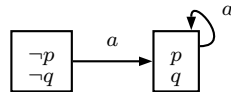
(i.e., nonexecutable  $\alpha$  if  $\alpha \text{ counts\_as } \beta$  in abbreviated form). The separate form (10) is required in the latter case because the expression  $\perp \text{ if } \alpha \wedge \alpha \text{ counts\_as } \beta$  is not a well-formed action dynamic law of  $\mathcal{C}+$  ( $\perp$  is a 0-ary connective and not an action formula of  $\sigma^a$ ). And note that although there is no constant of the form  $\alpha \text{ counts\_as } \top$  in  $\sigma^\gamma$  ( $\top$  is a 0-ary connective) there could be a conjunct  $\top$  in the formula  $\beta$  of a constant  $\alpha \text{ counts\_as } \beta$ . No  $\mathcal{C}+$  law is included in  $\delta_D$  in this case. ( $\top \text{ if } \top \text{ after } \alpha \text{ counts\_as } \beta$  is well-formed but it has no effect on the models of a  $\mathcal{C}+$  action description and so is omitted from  $\delta_D$ .)

The other component  $\delta_{\sigma^\gamma}$  of  $\delta$  is the set of  $\mathcal{C}+$  static laws of the form

$$(\alpha \text{ counts\_as } \beta \leftrightarrow \alpha' \text{ counts\_as } \beta') \text{ if } \top \tag{11}$$

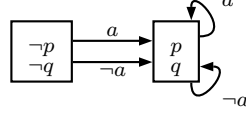
for every pair of distinct constants  $\alpha \text{ counts\_as } \beta$  and  $\alpha' \text{ counts\_as } \beta'$  in  $\sigma^\gamma$  such that  $\alpha \equiv_{\sigma^a} \alpha'$  and  $\beta \equiv_{\sigma^a} \beta'$ . These are required to ensure that interpretations/states agree on the values of  $\sigma^a$ -equivalent  $\text{counts\_as}$  constants.

It may seem that the causal laws  $\delta_{\sigma^\gamma}$  are unnecessarily cumbersome to get the desired effect. But consider: suppose we want simple Boolean constants  $p$  and  $q$  to have the same values in all states. (We choose  $p$  and  $q$  here to make the point rather than  $\text{counts\_as}$  atoms to reduce writing.) Suppose, for illustration, that we have the single fluent dynamic law  $a \text{ causes } p$ . Adding the static law  $(p \leftrightarrow q) \text{ if } \top$  gives the following set of states and transitions

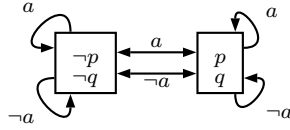


which is also what is required by the semantics.

If in place of  $(p \leftrightarrow q)$  if  $\top$  we try the pair of static laws  $\perp$  if  $p \wedge \neg q$  and  $\perp$  if  $q \wedge \neg p$  we get the right states but we get no transitions. (There is nothing to ‘cause’  $q$  in the state after  $a$ .) If we try instead the pair of static laws  $p$  if  $q$  and  $q$  if  $p$  we again get the right states but now *too many* transitions. We have the implied laws  $p$  if  $p$  and  $q$  if  $q$ , and we get the following states and transitions:



If we add two further laws  $\neg p$  if  $\neg q$  and  $\neg q$  if  $\neg p$  then  $p$  and  $q$  both become (implied) exogenous, and we get even more transitions:



The form of causal laws  $\delta_D$  in (9) and (10) is chosen so that when  $D$  is a definite action description of  $(\mathcal{C}+)^+$  (defined below), and  $Equiv(\sigma^\gamma, \sigma^a)$  and  $\delta_{\sigma^\gamma}$  are empty, then  $D \cup \delta_D = D \cup \delta_D \cup \delta_{\sigma^\gamma}$  is a definite action description of  $\mathcal{C}+$ . There will be no causal laws of form (10) in  $\delta_D$  in this case. However, the soundness of the translation from  $(\mathcal{C}+)^+$  to  $\mathcal{C}+$  for the general case does not depend on any assumptions about the form of  $\beta$  in  $\alpha$  counts\_as  $\beta$  constants.

**Theorem 15** *Let  $D$  be a  $(\mathcal{C}+)^+$  action description of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ . Let  $\delta_D$  be the set of  $\mathcal{C}+$  action dynamic laws*

$$\beta_i \text{ if } \alpha \wedge \alpha \text{ counts\_as } \beta_1 \wedge \dots \wedge \beta_i \wedge \dots \wedge \beta_m$$

*for every constant  $\alpha$  counts\_as  $\beta_1 \wedge \dots \wedge \beta_i \wedge \dots \wedge \beta_m$  ( $m \geq 1$ ) in  $\sigma^\gamma$  for which  $\beta_k \neq \perp$  for any  $1 \leq k \leq m$  and every  $\beta_i$  ( $1 \leq i \leq m$ ) such that  $\beta_i \neq \top$ , together with  $\mathcal{C}+$  fluent dynamic laws*

$$\perp \text{ if } \top \text{ after } \alpha \wedge \alpha \text{ counts\_as } \beta$$

*for every constant  $\alpha$  counts\_as  $\beta$  in  $\sigma^\gamma$  for which  $\perp$  is a conjunct of  $\beta$ . Let  $\delta_{\sigma^\gamma}$  be the set of  $\mathcal{C}+$  static laws*

$$(\alpha \text{ counts\_as } \beta \leftrightarrow \alpha' \text{ counts\_as } \beta') \text{ if } \top$$

*for every pair of distinct constants  $\alpha$  counts\_as  $\beta$  and  $\alpha'$  counts\_as  $\beta'$  in  $\sigma^\gamma$  such that  $\alpha \equiv_{\sigma^a} \alpha'$  and  $\beta \equiv_{\sigma^a} \beta'$ .*

*$s$  is a state of the ATS defined by  $D$  iff  $s$  is a state of the labelled transition system defined by the  $\mathcal{C}+$  action description  $D \cup \delta_D \cup \delta_{\sigma^\gamma}$  of signature  $(\sigma^f \cup \sigma^\gamma, \sigma^a)$ .  $(s, \varepsilon, s')$  is a transition of the ATS defined by  $D$  iff  $(s, \varepsilon, s')$  is a transition of the labelled transition system defined by the  $\mathcal{C}+$  action description  $D \cup \delta_D \cup \delta_{\sigma^\gamma}$  of signature  $(\sigma^f \cup \sigma^\gamma, \sigma^a)$ .*

*Proof.* The expressions  $T_{static}(s)$ ,  $E(s, \varepsilon, s')$ ,  $A(\varepsilon, s)$  refer here to the causal laws of the  $(\mathcal{C}+)^+$  action description  $D$ .

*States:* The conditions for  $s \in S$  to be a  $(\mathcal{C}+)^+$ -defined augmented state of  $D$  are:

$$s \models T_{static}(s) \cup Simple(s) \cup Equiv(\sigma^\gamma, \sigma^a)$$

Compare the conditions for  $s$  to be a state of the  $\mathcal{C}+$  action description  $D \cup \delta_D \cup \delta_{\sigma^\gamma}$ :

$$s \models (T_{static}(s) \cup Equiv(\sigma^\gamma, \sigma^a)) \cup Simple(s)$$

These are clearly equivalent.

*Transitions:* To prove the result for transitions we show that the conditions

$$\begin{aligned} s' \models T_{static}(s') \cup E(s, \varepsilon, s') \cup Equiv(\sigma^\gamma, \sigma^a) \\ \varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon) \cup \{\alpha \rightarrow \beta \mid s \models \alpha \text{ counts\_as } \beta\} \end{aligned}$$

are equivalent to the conditions

$$\begin{aligned} s' \models (T_{static}(s') \cup Equiv(\sigma^\gamma, \sigma^a)) \cup (E(s, \varepsilon, s') \cup E_{\delta_D}(s, \varepsilon, s')) \\ \varepsilon \models (A(\varepsilon, s) \cup A_{\delta_D}(\varepsilon, s)) \cup Exog(\varepsilon) \end{aligned}$$

where  $E_{\delta_D}(s, \varepsilon, s')$  and  $A_{\delta_D}(\varepsilon, s')$  denote respectively the heads of the fluent dynamic and action dynamic laws in  $\delta_D$  whose bodies are satisfied by the transition  $(s, \varepsilon, s')$ , i.e.,

$$\begin{aligned} E_{\delta_D}(s, \varepsilon, s') &= \{\perp \mid \alpha \text{ counts\_as } \beta \text{ in } \sigma^\gamma, \perp \text{ is a conjunct of } \beta, \\ &\quad s' \models \top, s \cup \varepsilon \models \alpha \wedge \alpha \text{ counts\_as } \beta\} \\ &= \{\perp \mid s \models \alpha \text{ counts\_as } \beta, \perp \text{ is a conjunct of } \beta, \varepsilon \models \alpha\} \end{aligned}$$

and

$$\begin{aligned} A_{\delta_D}(\varepsilon, s) &= \{\beta_i \mid \alpha \text{ counts\_as } \beta \text{ in } \sigma^\gamma, \perp \text{ is not a conjunct of } \beta, \\ &\quad \beta_i \text{ is a conjunct of } \beta, \beta_i \neq \top, s \cup \varepsilon \models \alpha \wedge \alpha \text{ counts\_as } \beta\} \\ &= \{\beta_i \mid s \models \alpha \text{ counts\_as } \beta, \perp \text{ is not a conjunct of } \beta, \\ &\quad \beta_i \text{ is a conjunct of } \beta, \beta_i \neq \top, \varepsilon \models \alpha\} \end{aligned}$$

The conditions on  $s'$  and  $\varepsilon$  above are equivalently expressed (because  $s'$  and  $\varepsilon$  are disjoint) as

$$\begin{aligned} s' \cup \varepsilon \models T_{static}(s') \cup Equiv(\sigma^\gamma, \sigma^a) \cup E(s, \varepsilon, s') \cup E_{\delta_D}(s, \varepsilon, s') \cup \\ A(\varepsilon, s) \cup A_{\delta_D}(\varepsilon, s) \cup Exog(\varepsilon) \end{aligned}$$

which in turn can be expressed (because  $s'$  and  $\varepsilon$  are disjoint) as

$$\begin{aligned} s' \models T_{static}(s') \cup E(s, \varepsilon, s') \cup Equiv(\sigma^\gamma, \sigma^a) \\ \varepsilon \models A(\varepsilon, s) \cup Exog(\varepsilon) \cup A_{\delta_D}(\varepsilon, s) \cup E_{\delta_D}(s, \varepsilon, s') \end{aligned}$$

Now

$$\begin{aligned} \varepsilon \models A_{\delta_D}(\varepsilon, s) \text{ iff } \varepsilon \models \{\beta_i \mid s \models \alpha \text{ counts\_as } \beta, \perp \text{ is not a conjunct of } \beta, \\ \beta_i \text{ is a conjunct of } \beta, \beta_i \neq \top, \varepsilon \models \alpha\} \\ \text{iff } \varepsilon \models \{\beta_i \mid s \models \alpha \text{ counts\_as } \beta, \perp \text{ is not a conjunct of } \beta, \\ \beta_i \text{ is a conjunct of } \beta, \varepsilon \models \alpha\} \end{aligned}$$

and

$$\begin{aligned} \varepsilon \models E_{\delta_D}(\varepsilon, s) \text{ iff } \varepsilon \models \{\perp \mid s \models \alpha \text{ counts\_as } \beta, \perp \text{ is a conjunct of } \beta, \varepsilon \models \alpha\} \\ \text{iff } \varepsilon \models \{\beta_i \mid s \models \alpha \text{ counts\_as } \beta, \perp \text{ is a conjunct of } \beta, \\ \beta_i \text{ is a conjunct of } \beta, \varepsilon \models \alpha\} \end{aligned}$$

And so

$$\begin{aligned} \varepsilon \models A_{\delta_D}(\varepsilon, s) \cup E_{\delta_D}(\varepsilon, s) \\ \text{iff } \varepsilon \models \{\beta_i \mid s \models \alpha \text{ counts\_as } \beta, \beta_i \text{ is a conjunct of } \beta, \varepsilon \models \alpha\} \\ \text{iff } \varepsilon \models \{\beta \mid s \models \alpha \text{ counts\_as } \beta, \varepsilon \models \alpha\} \\ \text{iff } \varepsilon \models \{\alpha \rightarrow \beta \mid s \models \alpha \text{ counts\_as } \beta\} \end{aligned}$$

which completes the proof.  $\square$

**Corollary 16** *Paths of length  $m$  in the ATS defined by a  $(\mathcal{C}+)^+$  action description  $D$  are in one-to-one correspondence with the models of the causal theory  $\Gamma_m^{D \cup \delta_D \cup \delta_{\sigma^\gamma}}$ .*

## 8.5 Reduced action descriptions

The causal laws  $\delta_{\sigma^\gamma}$  required for translation to  $\mathcal{C}+$  are a nuisance—for one thing, they mean that in general a  $(\mathcal{C}+)^+$  action description cannot translate to a definite action description of  $\mathcal{C}+$ . They are included primarily for the sake of completeness. There are a number of straightforward methods by which their effect can be obtained in other ways, depending on the level of generality required.

One obvious method is as follows. We want to restrict attention to interpretations that agree on the values of constants  $c_1, c_2, \dots, c_n$ . In the present context  $c_1, c_2, \dots, c_n$  are  $\sigma^a$ -equivalent `counts_as` constants in  $\sigma^\gamma$ . Instead of representing such interpretations by the set of all atoms they satisfy, it is obviously sufficient to represent an interpretation by the value of a single representative member of  $c_1, c_2, \dots, c_n$ . So: we partition the `counts_as` constants in  $\sigma^\gamma$  into  $\sigma^a$ -equivalent classes and pick one representative member from each equivalence class. This representative member can be identified when the signature  $\sigma^\gamma$  is specified or it can be chosen arbitrarily. Now we replace every occurrence of a `counts_as` constant in the action description by the representative member of the  $\sigma^a$ -equivalence class to which it belongs. This is straightforward to implement, and has the additional benefit that the representation of interpretations is more concise. It only remains to ensure that the evaluation of query languages on the ATS similarly translates formulas before they are evaluated, which is also easy to arrange.

We will call a  $(\mathcal{C}+)^+$  action description  $D$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  *reduced* if there are no distinct constants  $\alpha \text{ counts\_as } \beta$  and  $\alpha' \text{ counts\_as } \beta'$  in  $\sigma^\gamma$  such that  $\alpha \equiv_{\sigma^a} \alpha'$  and  $\beta \equiv_{\sigma^a} \beta'$ . Clearly every action description of  $(\mathcal{C}+)^+$  can be translated to a reduced action description, for example by the process outlined above. For a reduced action description  $D$  the set  $\text{Equiv}(\sigma^\gamma, \sigma^a)$  is empty and so the conditions for  $s$  to be a state of the ATS  $(\mathcal{C}+)^+$ -defined by  $D$  are simply

$s \models T_{static}(s) \cup Simple(s)$  (as for a  $\mathcal{C}+$  action description), and the translation of the  $(\mathcal{C}+)^+$  action description to  $\mathcal{C}+$  can dispense with the static laws  $\delta_{\sigma^\gamma}$  since these are also empty.

In what follows, we will not employ this general method but instead restrict attention to a special form. We will say that an action description of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ , is *atom-conjunctive* when every  $\alpha \text{ counts\_as } \beta$  constant of  $\sigma^\gamma$  is such that the formulas  $\alpha$  and  $\beta$  are both (non-empty) conjunctions of action *atoms*, written in some standard order and without repetitions. As an implementation issue, it is easy to ensure that the conjuncts in `counts_as` constants are re-arranged into the standard order if necessary in a simple pre-processing stage. Clearly all atom-conjunctive action descriptions are reduced. (And as a minor consideration, the translation of  $(\mathcal{C}+)^+$  to  $\mathcal{C}+$  does not introduce any fluent dynamic laws of the ‘nonexecutable’ form (10).)

Why this particular form? Although in practice it will often be the case that both  $\alpha$  and  $\beta$  in  $\alpha \text{ counts\_as } \beta$  constants are *atoms* of  $\sigma^a$ , we want to be able to support the more general case where they are conjunctions of atoms, for the following main reason. An action of type  $a$  with attributes  $x_1, x_2, \dots, x_n$  can be represented in  $\mathcal{C}+$  (and  $(\mathcal{C}+)^+$ ) either as an action atom  $a(x_1, x_2, \dots, x_n) = v$  with  $n$  arguments or as a conjunction of  $n + 1$  action atoms  $a = v \wedge a_1 = x_1 \wedge a_2 = x_2 \wedge \dots \wedge a_n = x_n$  where  $a_1, \dots, a_n$  are (usually non-exogenous) action constants which pick out the values of the attributes. This second style of representation has significant advantages, both in terms of flexibility or ‘elaboration tolerance’ (Giunchiglia et al., 2004, Sect. 5.6), and economy, in that the number of action constants required in the first style is exponential in the number of attributes  $n$  whereas in the conjunction-of-attributes style it is linear.

The atom-conjunctive form accommodates all the examples that we can expect to encounter in practice. Moreover, it is likely that the properties of the ‘counts as’ relation  $C_T^+$  of an ATS can be strengthened quite naturally, allowing all  $\alpha \text{ counts\_as } \beta$  atoms, for arbitrary action formulas  $\alpha$ , to be written equivalently in a normal form as a truth-functional compound of  $\alpha' \text{ counts\_as } \beta$  atoms in which  $\alpha'$  is a conjunction of action atoms. This requires only that  $(\alpha_1 \vee \alpha_2) \text{ counts\_as } \beta$  comes out to be equivalent to the conjunction of  $\alpha_1 \text{ counts\_as } \beta$  and  $\alpha_2 \text{ counts\_as } \beta$ , which seems non-controversial. However, this remains to be checked, and we leave the details to a future version.

## 8.6 Definite action descriptions

As in the case of  $\mathcal{C}+$ , the action descriptions of primary interest are those that are definite.

**Definition 17** *A  $(\mathcal{C}+)^+$  action description  $D$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$  is definite iff the  $\mathcal{C}+$  action description  $D$  of signature  $(\sigma^f \cup \sigma^\gamma, \sigma^a)$  is definite, and for every constant  $\alpha \text{ counts\_as } \beta$  in  $\sigma^\gamma$ ,  $\beta$  is a (non-empty) conjunction of atoms of  $\sigma^a$ , written in some standard order and without repetitions.*

Note that this definition does not depend on the form of action formulas  $\alpha$  in  $\alpha \text{ counts\_as } \beta$  constants, so an action description can be definite without being atom-conjunctive (or even ‘reduced’). However, as already explained, we are

primarily interested in the special category of definite action descriptions which are also atom-conjunctive.

For definite, atom-conjunctive action descriptions, the characterisation of states and transitions of the ATS they define can be simplified as usual. For reference, the relevant conditions are as follows.

**Proposition 18** *Let  $D$  be a definite, atom-conjunctive action description of  $(\mathcal{C}+)^+$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ .  $s$  is a state of  $D$  and  $(s, \varepsilon, s')$  is a transition of  $D$  iff  $s$  and  $s'$  are interpretations of  $\sigma^f \cup \sigma^\gamma$  and  $\varepsilon$  is an interpretation of  $\sigma^a$  such that:*

- $s = T_{static}(s) \cup Simple(s) \cup Trivial(\sigma^f \cup \sigma^\gamma)$
- $s' = T_{static}(s') \cup E(s, \varepsilon, s') \cup Trivial(\sigma^f \cup \sigma^\gamma)$
- $\varepsilon = A(\varepsilon, s) \cup Exog(\varepsilon) \cup Trivial(\sigma^a) \cup \{ \beta_i \mid s \models \alpha \text{ counts\_as } \beta_1 \wedge \dots \wedge \beta_i \wedge \dots \wedge \beta_n, \varepsilon \models \alpha \}$

Further,  $D \cup \delta_D$  is a definite action description of  $\mathcal{C}+$  when  $D$  is a definite action description of  $(\mathcal{C}+)^+$ , and so we have the following.

**Proposition 19** *Let  $D$  be a definite, atom-conjunctive action description of  $(\mathcal{C}+)^+$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ . Paths of length  $m$  of the ATS defined by  $D$  are in one-to-one correspondence with the (classical) models of  $comp(\Gamma_m^{D \cup \delta_D})$ . In particular, states of the ATS are encoded by the models of  $comp(\Gamma_0^{D \cup \delta_D})$  and transitions by the models of  $comp(\Gamma_1^{D \cup \delta_D})$ .*

**Example** Consider the following example. It is rather artificial but it is difficult to construct a convincing example without making it too large and complicated. A more realistic example is presented in Section 8.7 below.

Signature: exogenous action constants *wave\_flag* and *wave\_hand* (Boolean) and *switch* with values *on* and *off*; non-exogenous action constants *start\_race*, *signalX* (Boolean); simple fluent constants *statusX* (Boolean) and *light* with values *on* and *off*; simple (Boolean) constants *wave\_flag counts\_as start\_race*, *switch = on counts\_as signalX*, *wave\_hand counts\_as signalX*.

The following is a definite, atom-conjunctive action description of  $(\mathcal{C}+)^+$ :

```

inertial light
inertial wave_flag counts_as start_race
inertial switch = on counts_as signalX
inertial wave_hand counts_as signalX
inertial statusX

switch = on causes light = on
switch = off causes light = off

wave_flag counts_as start_race if  $\top$  % fixed, for the sake of the example
start_race causes (switch = on counts_as signalX)
start_race causes (wave_hand counts_as signalX)
signalX causes statusX

nonexecutable start_race  $\wedge$   $\neg$ wave_flag
nonexecutable signalX  $\wedge$   $\neg$ switch = on  $\wedge$   $\neg$ wave_hand

```

(We allow the possibility of waving the flag, waving a hand, and switching the light at the same time.)

The translation to  $\mathcal{C}+$  adds the following action dynamic laws:

$$\begin{aligned} & \textit{start\_race} \text{ if } \textit{wave\_flag} \wedge (\textit{wave\_flag} \text{ counts\_as } \textit{start\_race}) \\ & \textit{signalX} \text{ if } \textit{switch} = \textit{on} \wedge (\textit{switch} = \textit{on} \text{ counts\_as } \textit{signalX}) \\ & \textit{signalX} \text{ if } \textit{wave\_hand} \wedge (\textit{wave\_hand} \text{ counts\_as } \textit{signalX}) \end{aligned}$$

In the ATS so defined, there are no transitions of type  $\textit{wave\_flag} \wedge \neg \textit{start\_race}$ : if a  $\textit{wave\_flag}$  action occurs then the transition must also be of type  $\textit{start\_race}$ . And since  $\textit{start\_race}$  ‘causes’  $\textit{switch} = \textit{on} \text{ counts\_as } \textit{signalX}$ , after a  $\textit{wave\_flag}$  event, there are no transitions of type  $\textit{switch} = \textit{on} \wedge \neg \textit{signalX}$ : after a  $\textit{wave\_flag}$  event, any transition of type  $\textit{switch} = \textit{on}$  is also of type  $\textit{signalX}$ ; its effects are those it has by virtue of being a transition of type  $\textit{switch} = \textit{on}$  ( $\textit{light} = \textit{on}$ ) and those it has by virtue of being a transition of type  $\textit{signalX}$  ( $\textit{statusX}$  is true). Likewise for  $\textit{wave\_hand}$  transitions following a  $\textit{wave\_flag}$ . The two nonexecutable laws of the action description are included to eliminate unsupported institutional actions.

## 8.7 Example: an auction protocol

Consider a simple kind of auction in which agents make bids as follows. Any agent can *open* the bidding (when there is no current bid) at an amount  $N$  up to some designated maximum. Any agent can *raise* the current bid by  $N$  units, where again  $N$  must not exceed the designated maximum. Bids can be made by agents in any order, except that no agent may bid twice in succession. An agent may *revoke* its last bid, but only until the next (valid) bid is made. An agent may *withdraw* from the auction at any time, unless it has made the last (valid) bid. Once withdrawn, an agent may not resume bidding. (The ‘winner’ is the last agent left in the auction when all others have withdrawn, though for simplicity we will ignore this here.)

The terms ‘may’ and ‘can’ are to be understood not as permission, but as part of the definition of what counts as a valid bid. The representation should be able to express both attempted but possibly invalid bids, and valid bids.

We have trivial (and so ‘rigid’) fluent constants  $\textit{player}(X)$  for  $X$  ranging over the names of participants in the auction. There is also a trivial fluent atom  $\textit{max\_raise} = \textit{Max}$  for some positive integer  $\textit{Max}$ . The action constants are Boolean constants ( $X \text{ signals } \alpha$ ) and ( $X:\alpha$ ) for  $\alpha$  ranging over  $\textit{open}(N)$ ,  $\textit{raise}(N)$ ,  $\textit{withdraw}$ ,  $\textit{revoke}$ , and  $N$  ranging over all integers between 1 and  $\textit{Max}$ . The action constants ( $X \text{ signals } \alpha$ ) are exogenous, and ( $X:\alpha$ ) are not. We could also have made ( $X \text{ signals } \textit{open}$ ), ( $X \text{ signals } \textit{raise}$ ), ( $X:\textit{open}$ ) and ( $X:\textit{raise}$ ) multi-valued constants with values in the range  $1..M$  for some integer  $M$  but there is little to be gained from doing so in this example.

To represent the state of the auction we have (inertial) multi-valued fluent constants  $\textit{current\_bidder}$  and  $\textit{current\_bid}$ , and an inertial (Boolean) fluent constant  $\textit{withdrawn}(X)$  for each player  $X$ . The possible values of  $\textit{current\_bid}$  are all integers in the range  $1..M_{\text{max}}$  for some suitably chosen integer  $M_{\text{max}}$ .

We also employ the abbreviation

$$\textit{pow}(X, \alpha) =_{\text{def}} (X \text{ signals } \alpha) \text{ counts\_as } (X:\alpha)$$

The ‘counts as’ signature  $\sigma^\gamma$  contains all instances of the **pow** expressions. And instead of making them simple and inertial, in this example we will make them statically determined and false by default.

**Version 1** Although it is a simple idea, revocation complicates the formulation substantially. So first we show a version without revocation, so as not to distract attention from the main points. We will discuss revocation separately after that. Preliminary:

$$\begin{aligned} & \text{inertial } \textit{current\_bidder}, \textit{current\_bid}, \textit{withdrawn}(X) \\ & \text{nonexecutable } (X \textit{ signals } \alpha) \wedge (X \textit{ signals } \alpha') \end{aligned}$$

for all  $\alpha \neq \alpha'$  ranging over *open* ( $N$ ), *raise* ( $N$ ), *withdraw*, *revoke*. The above allows for concurrent actions by different players. These can be eliminated if desired by replacing the **nonexecutable** laws above by adding

$$\begin{aligned} & \text{inertial } \textit{current\_bidder}, \textit{current\_bid}, \textit{withdrawn}(X) \\ & \text{nonexecutable } (X \textit{ signals } \alpha) \wedge (X' \textit{ signals } \alpha') \end{aligned}$$

for all  $X \neq X'$  ranging over the names of the players.

Now we specify the powers. The **nonexecutable** laws below eliminate unsupported institutional actions.

$$\begin{aligned} & \text{nonexecutable } (X:\alpha) \wedge \neg(X \textit{ signals } \alpha) \\ & \text{default } \neg\text{pow}(X, \alpha) \quad \% \text{ statically determined, and not inertial} \\ & \text{pow}(X, \textit{open}(N)) \text{ if} \\ & \quad \textit{player}(X) \wedge \neg\textit{withdrawn}(X) \wedge \\ & \quad \textit{max\_raise} = \textit{Max} \wedge 0 < N \leq \textit{Max} \wedge \\ & \quad \textit{current\_bid} = \textit{none} \\ & \text{pow}(X, \textit{raise}(N)) \text{ if} \\ & \quad \textit{player}(X) \wedge \neg\textit{withdrawn}(X) \wedge \\ & \quad \textit{max\_raise} = \textit{Max} \wedge 0 < N \leq \textit{Max} \wedge \\ & \quad \textit{current\_bidder} \neq X \\ & \text{pow}(X, \textit{withdraw}) \text{ if} \\ & \quad \textit{player}(X) \wedge \neg\textit{withdrawn}(X) \wedge \\ & \quad \textit{current\_bidder} \neq X \end{aligned}$$

Now we specify the effects of (successful, valid) actions:

$$\begin{aligned} & X:\textit{open}(N) \text{ causes } \textit{current\_bidder} = X \\ & X:\textit{open}(N) \text{ causes } \textit{current\_bid} = N \\ & X:\textit{raise}(N) \text{ causes } \textit{current\_bidder} = X \\ & X:\textit{raise}(N) \text{ causes } \textit{current\_bid} = \textit{Current} + N \text{ if} \\ & \quad \textit{current\_bid} = \textit{Current} \\ & X:\textit{withdraw} \text{ causes } \textit{withdrawn}(X) \end{aligned}$$

The above uses some simple arithmetic, easily added. See e.g. the implementation provided as part of the Causal Calculator CCALC<sup>4</sup>.

<sup>4</sup><http://www.cs.utexas.edu/users/tag/cc>



**Version 2** The complication with revocation is that we have to remember at least the last previous bid made so we know what to re-instate when revocation takes place. We do not want to implement a *stack* of bids—not in this example anyway. Notice that the extra complication is part of the temporal structure of the bidding process and not a result of the need to specify `counts_as` or `pow`. It is a limitation of  $\mathcal{C}+$ , and of transition systems in general, that a state is independent of the history of transitions by which it was reached.

So introduce another pair of inertial (multi-valued) fluent constants, *previous\_bid* and *previous\_bidder*. They are inertial because some actions (such as withdrawals) do not affect the current bid. A *raise* action now changes the value of *current\_bid* and *current\_bidder* as before, but first sets *previous\_bid* and *previous\_bidder*. In addition to the laws of version 1 we have:

$$\begin{aligned} X:\text{raise}(N) \text{ causes } \text{previous\_bidder} = Y \text{ if } \text{current\_bidder} = Y \\ X:\text{raise}(N) \text{ causes } \text{previous\_bid} = \text{Tot} \text{ if } \text{current\_bid} = \text{Tot} \end{aligned}$$

where  $Y$  ranges over the names of the players and  $\text{Tot}$  ranges over integers in the range  $1..M_{\max}$ .

Revocation when it happens re-instates the *previous\_bid* and *previous\_bidder* as the *current\_bid* and *current\_bidder*:

$$\begin{aligned} X:\text{revoke} \text{ causes } \text{current\_bidder} = Y \text{ if } \text{previous\_bidder} = Y \\ X:\text{revoke} \text{ causes } \text{current\_bid} = \text{Tot} \text{ if } \text{previous\_bid} = \text{Tot} \end{aligned}$$

Notice that this formulation also deals with revocations of the opening bid. In that case only the second of the clauses above has any effect: the first and third have no effect because there is no *previous\_bid* following an opening bid.

It remains to specify when agent  $X$  is empowered to revoke. One way is to attempt to specify it by formulating appropriate combinations of *previous\_bid* and *current\_bid*. Clearer, and more natural, is to do it by specifying further effects of *open*, *raise*, and *revoke*, as follows:

$$\begin{aligned} X:\text{open}(N) \text{ causes } \text{pow}(X, \text{revoke}) \\ X:\text{raise}(N) \text{ causes } \text{pow}(X, \text{revoke}) \\ X:\text{raise}(N) \text{ causes } \neg\text{pow}(Y, \text{revoke}) \text{ if } \text{pow}(Y, \text{revoke}) \\ X:\text{revoke} \text{ causes } \neg\text{pow}(X, \text{revoke}) \end{aligned}$$

This completes the specification of the example.

To illustrate that instances of `counts_as` can sometimes be chained together, suppose that agent Bob is to act as a representative in the auction of another agent Charles. This arrangement might be expressed as follows:

$$\text{Bob signals raise}(N) \text{ counts\_as } \text{Charles signals raise}(N) \text{ if } \top$$

There are other forms of representation, also expressible in the language  $(\mathcal{C}+)^+$ . We do not attempt any systematic account of the possibilities here. Some suggestions along these lines, using a different formalism, are discussed for example in (Gelati et al., 2004).

Note also that actions can have institutional effects which are not of the ‘power’ kind. Suppose, for example, that Charles promises or is instructed that every

time he makes a bid, he must inform his superior, Clare. Here, a bid by Charles creates an obligation on him to perform another action. There is some ambiguity here about whether Charles is obliged to inform Clare of every attempted bid ('signals') or only of successful bids. The distinction is evidently expressible in the  $(\mathcal{C}+)^+$  language.

**Further structuring** It would be advantageous to introduce some further structure in the specification of `pow`, for instance by separating out the time-dependent from the time-independent parts. More generally, we might separate the components that make a bid 'valid' (better perhaps: 'well formed') from the conditions defining `pow` and `counts_as`. The difficulty here is not in devising such representations but in fixing on terminology. Terms such as 'valid', 'well formed', 'inadmissible', 'improper', 'out of order', 'void', 'voidable' have specific meanings in certain contexts, but these contexts are relatively few, and the same meaning is not always given in each. There is unfortunately much scope for confusion here. It would be better to identify and fix on some neutral terms, but these are difficult to find. Some preliminary suggestions are provided in (Prakken, 1998; Prakken and Gordon, 1999).

Examples of what can be done with such formalisations are provided in (Artikis et al., 2002, 2003a,b). There, similar protocols are formalised using event calculus and/or the language  $\mathcal{C}+$  as the action formalism and then subjected to various kinds of analysis. These examples can be reformulated using the additional resources of the language  $(\mathcal{C}+)^+$  instead. (We do not present these reformulations here.) The advantage of using  $(\mathcal{C}+)^+$  is that it provides a built-in treatment of `counts_as` and `pow`, and their properties, which otherwise have to be anticipated and coded up explicitly as  $\mathcal{C}+$  laws.

## 8.8 Unsupported institutional actions

We have chosen to leave the elimination of unsupported institutional actions to the formulation of action descriptions. A case can be made that this is unreasonable, that the elimination of unsupported institutional actions is fundamental to the notion of 'counts as' being modelled, and that this should therefore not be a feature of action descriptions but rather a fixed semantic feature of ATS structures. It is easy to formulate the required conditions. Where  $\alpha_1 \text{ counts\_as } \beta, \dots, \alpha_n \text{ counts\_as } \beta$  are all the constants of the form  $\alpha \text{ counts\_as } \beta$  in the signature  $\sigma^\gamma$ , we require that every transition  $(s, \varepsilon, s')$  of the ATS has the property that  $\varepsilon \not\models \beta \wedge \neg \alpha_1 \wedge \dots \wedge \neg \alpha_n$ , i.e., that  $\varepsilon \models \beta \rightarrow (\alpha_1 \vee \dots \vee \alpha_n)$ . (It may seem that this is equivalent to the condition  $\varepsilon \models \beta \leftrightarrow (\alpha_1 \vee \dots \vee \alpha_n)$  but this is not so, since we require  $\varepsilon \models \alpha_i \rightarrow \beta$  only when  $\alpha_i \text{ counts\_as } \beta$  holds in state  $s$  of the transition  $(s, \varepsilon, s')$ , whereas the condition required to eliminate unsupported institutional actions depends only on the signature  $\sigma^\gamma$  and not on the `counts_as` atoms that actually hold in state  $s$ .)

However, this treatment would also pre-suppose an assumption of completeness in the specification of the  $\sigma^\gamma$  signature. It would assume that there are no other means of effecting the institutional action  $\beta$  besides those actions  $\alpha_i$  that appear in  $\alpha_i \text{ counts\_as } \beta$  constants in  $\sigma^\gamma$ . It is at least conceivable that this assumption would not always be desired, and in these circumstances transition labels/events

with unsupported institutional actions would be meaningful even if in general they are not. It is, in any case, very easy to ensure that an action description eliminates (undesired) unsupported institutional actions. To eliminate an unsupported institutional action  $\beta$  one includes in the action description a fluent dynamic law

$$\text{nonexecutable } \beta \wedge \neg\alpha_1 \wedge \dots \wedge \neg\alpha_n$$

where  $\alpha_1 \text{ counts\_as } \beta, \dots, \alpha_n \text{ counts\_as } \beta$  are all the constants of the form  $\alpha \text{ counts\_as } \beta$  in the signature  $\sigma^\gamma$ . This guarantees that all transitions  $(s, \varepsilon, s')$  of the ATS defined satisfy the property  $\varepsilon \models \beta \rightarrow (\alpha_1 \vee \dots \vee \alpha_n)$ , without introducing any new transitions. For these two reasons together (no assumptions of completeness in  $\sigma^\gamma$ , and the ease of adding the required **nonexecutable** laws in the action description where desired) we do not include the elimination of unsupported institutional actions as a fixed feature of the semantics.

Some care has to be taken when deciding on the ‘counts as’ signature  $\sigma^\gamma$ . Compare a signature containing the single constant  $\alpha \text{ counts\_as } \beta_1 \wedge \beta_2$ , on the one hand, with a signature containing the pair of constants  $\alpha \text{ counts\_as } \beta_1$  and  $\alpha \text{ counts\_as } \beta_2$ , on the other. These have different effects. In both cases, a transition of type  $\beta_1 \wedge \beta_2 \wedge \neg\alpha$  represents (assuming no other relevant **counts\\_as** constants) an unsupported institutional action:  $\beta_1 \wedge \beta_2$  is unsupported in the first case, and both  $\beta_1$  and  $\beta_2$  are unsupported in the second. However, in the first case a transition of type  $\beta_1 \wedge \neg\beta_2 \wedge \neg\alpha$  does not represent an unsupported institutional action, whereas in the second case it does ( $\beta_1$  is unsupported).

Implicit in this treatment of ‘counts as’ is that the conjunction  $\beta_1 \wedge \beta_2$  in a constant  $\alpha \text{ counts\_as } \beta_1 \wedge \beta_2$  represents a single action whose attributes are specified by  $\beta_1 \wedge \beta_2$ . In contrast, when  $\beta_1 \wedge \beta_2$  is intended to represent the concurrent execution of two separate actions  $\beta_1$  and  $\beta_2$ , two separate constants  $\alpha \text{ counts\_as } \beta_1$  and  $\alpha \text{ counts\_as } \beta_2$  should be employed.

As a further example, consider the difference in terms of unsupported institutional actions between a signature  $S1$  containing two constants

$$\begin{array}{l} \alpha_1 \text{ counts\_as } \beta \wedge \beta_1 \\ \alpha_2 \text{ counts\_as } \beta \wedge \beta_2 \end{array} \quad (S1)$$

and a signature  $S2$  containing four separate constants

$$\begin{array}{ll} \alpha_1 \text{ counts\_as } \beta & \alpha_1 \text{ counts\_as } \beta_1 \\ \alpha_2 \text{ counts\_as } \beta & \alpha_1 \text{ counts\_as } \beta_2 \end{array} \quad (S2)$$

In the case  $S1$ ,  $\beta$  would typically be a conjunction of some common set of core attribute-values and  $\beta_1$  and  $\beta_2$  conjunctions of additional attribute-values to represent the two different types of  $\beta$  actions. In the case  $S2$ ,  $\beta$ ,  $\beta_1$ , and  $\beta_2$  would each represent a particular kind of action, though of course there is nothing in the syntax of  $\mathcal{C}+$  (or of  $(\mathcal{C}+)^+$ ) that indicates the difference explicitly.

In the case  $S1$ , a transition has an unsupported institutional action when it satisfies  $\beta \wedge \beta_1 \wedge \neg\alpha_1$  or  $\beta \wedge \beta_2 \wedge \neg\alpha_2$ ; in the case  $S2$  when it satisfies  $\beta \wedge \neg\alpha_1 \wedge \neg\alpha_2$  or  $\beta_1 \wedge \neg\alpha_1$  or  $\beta_2 \wedge \neg\alpha_2$ . One can see that if a transition has an unsupported institutional action with signature  $S1$  then it also has an unsupported institutional action with signature  $S2$ , though not *vice-versa*. However, if every transition label/event also satisfies  $\beta \leftrightarrow (\beta_1 \vee \beta_2)$ , then the two are equivalent. So: suppose

that in  $S1$  the conjunction  $\beta$  does represent some common set of core attribute values, and that there are exactly two kinds of  $\beta$  actions, one represented by the full set of attribute values  $\beta \wedge \beta_1$ , and the other represented by  $\beta \wedge \beta_2$ . Now suppose we make this explicit in the action description, e.g., by including the three fluent dynamic laws **nonexecutable**  $\beta \wedge \neg\beta_1 \wedge \neg\beta_2$ , **nonexecutable**  $\beta_1 \wedge \neg\beta$ , and **nonexecutable**  $\beta_2 \wedge \neg\beta$ . Then every transition label/event in the ATS defined satisfies  $\beta \leftrightarrow (\beta_1 \vee \beta_2)$ , and the signature  $S1$  and  $S2$  agree on the conditions for unsupported institutional actions.

There is clearly more to this. It may be worthwhile extending the syntax of  $(\mathcal{C}+)^+$  (and  $\mathcal{C}+$ ) to make explicit when a conjunction of action atoms is intended to represent the attributes of a single action and when it is intended to represent instead the concurrent execution of several distinct actions. One can see how this might be done, by having a separate category of action-attribute constants in the signature  $\sigma^a$ , for instance. We leave the details for future work.

## 8.9 Discussion: Possible strengthenings

We turn now to the possibility of strengthening the ‘counts as’ relation  $C_{\mathcal{I}}^+$  to a transitive relation or to an ordering on (the special category of ‘institutionally significant’) transition types, as outlined in the earlier discussion in Section 8.1 on comparisons with ‘counts as’ conditionals in (Jones and Sergot, 1996) and ‘conventional generation’ in (Goldman, 1970). Since reflexivity of ‘counts as’ is of negligible practical interest, there are two main issues: the difficulty of implementing the transitivity conditions in the translation from  $(\mathcal{C}+)^+$  to causal theories, and the choice of the most appropriate formulation of transitivity. We limit ourselves to some remarks about possibilities and the issues arising, and leave detailed investigation to a future version.

We note first that the constraint  $(C_{\varepsilon}^+)$  in every ATS already gives a kind of ‘pseudo-transitivity’ of  $C_{\mathcal{I}}^+$ . Expressed in terms of  $\tau_R(s, X)$ , the set of transitions of type  $X$  executable in state  $s$ ,  $(C_{\varepsilon}^+)$  takes the form (for all transition types  $X$ ,  $Y$  and all states  $s$ ):

$$\text{if } C_{\mathcal{I}}^+(s, X, Y) \text{ then } \tau_R(s, X) \subseteq \tau_R(s, Y)$$

from which follows the ‘pseudo-transitivity’ property (for all transition types  $X$ ,  $Y$ ,  $Z$  and all states  $s$ ):

$$\text{if } C_{\mathcal{I}}^+(s, X, Y) \text{ and } C_{\mathcal{I}}^+(s, X, Z) \text{ then } \tau_R(s, X) \subseteq \tau_R(s, Z)$$

as well as a (very weak) kind of ‘pseudo-anti-symmetry’

$$\text{if } C_{\mathcal{I}}^+(s, X, Y) \text{ and } C_{\mathcal{I}}^+(s, Y, X) \text{ then } \tau_R(s, X) = \tau_R(s, Y)$$

and a (trivial) ‘pseudo-reflexivity’

$$\text{if } C_{\mathcal{I}}^+(s, X, X) \text{ then } \tau_R(s, X) = \tau_R(s, X)$$

But suppose we wished to strengthen these properties as suggested in Section 8.1, to

$$\begin{array}{ll} C_{\mathcal{I}}^+(s, X, X) & (C_{\text{refl}}^+) \\ \text{if } C_{\mathcal{I}}^+(s, X, Y) \text{ and } C_{\mathcal{I}}^+(s, Y, Z) \text{ then } C_{\mathcal{I}}^+(s, X, Z) & (C_{\text{trans}}^+) \\ \text{if } C_{\mathcal{I}}^+(s, X, Y) \text{ and } C_{\mathcal{I}}^+(s, Y, X) \text{ then } X = Y & (C_{\text{a-symm}}^+) \end{array}$$

for all states  $s$  and for all transition types  $X, Y, Z \subseteq \mathbf{A}_{\mathcal{I}}$ , where, recall,  $\mathbf{A}_{\mathcal{I}}$  denotes the transition types of special conventional significance in institution  $\mathcal{I}$ , defined in this paper as  $\mathbf{A}_{\mathcal{I}} =_{\text{def}} \{X \subseteq \mathbf{A} \mid X = \|\alpha\| \text{ or } X = \|\beta\| \text{ for some constant } \alpha \text{ counts\_as } \beta \text{ in } \sigma^\gamma\}$ .

The first step is easy. It is easy to adjust the definition of the ATS defined by a  $(\mathcal{C}+)^+$  action description (Definition 13) so that the ‘counts as’ relation has the desired additional properties. Let  $D$  be an action description of  $(\mathcal{C}+)^+$  of signature  $(\sigma^f, \sigma^\gamma, \sigma^a)$ . The *super-augmented* transition system defined by  $D$  is the structure

$$\langle \sigma^f, \sigma^\gamma, S, \mathbf{A}, R, C_D^{+++} \rangle$$

where  $\langle \sigma^f, \sigma^\gamma, S, \mathbf{A}, R, C_D^+ \rangle$  is the ATS defined by  $D$  according to Definition 13 and where the new ‘counts as’ relation  $C_D^{+++}$  is defined to be the ‘reflexive transitive closure’ of  $C_D^+$ , or more precisely, the smallest relation containing  $C_D^+$  that satisfies the conditions  $(C_{\text{refl}}^+)$  and  $(C_{\text{trans}}^+)$ . If we further specify that  $C_D^+$  must be cycle-free, then  $C_D^{+++}$  also satisfies  $(C_{\text{a-symm}}^+)$ , and so is an ordering on the ‘institutionally significant’ transition types  $\mathbf{A}_{\mathcal{I}}$  in the sense discussed in Section 8.1. It is easy to check that the condition  $(C_\varepsilon^+)$  is satisfied by  $C_D^{+++}$  if it is satisfied by  $C_D^+$ , and so the structure defined in this way is an ATS.

So far so good. The problem comes in trying to find a causal theory, or a translation to a  $\mathcal{C}+$  action description, that will encode the paths of the ‘super-augmented’ transition system, and specifically additional causal laws that will construct  $C_D^{+++}$  as the ‘reflexive transitive’ closure of  $C_D^+$ .

It might appear that  $(C_{\text{refl}}^+)$  and  $(C_{\text{trans}}^+)$  can be captured straightforwardly by including further laws of the form:

$$\alpha \text{ counts\_as } \alpha \text{ if } \top \tag{12}$$

for every constant  $\alpha$  counts\_as  $\alpha$  in  $\sigma^\gamma$ ,

$$\alpha \text{ counts\_as } \beta \text{ if } \alpha \text{ counts\_as } \beta' \wedge \beta' \text{ counts\_as } \beta \tag{13}$$

for all constants  $\alpha$  counts\_as  $\beta$ ,  $\alpha$  counts\_as  $\beta'$ , and  $\beta'$  counts\_as  $\beta$  in  $\sigma^\gamma$ , and

$$\text{default } \neg(\alpha \text{ counts\_as } \beta) \tag{14}$$

for every constant  $\alpha$  counts\_as  $\beta$  in  $\sigma^\gamma$  to capture the closure condition.

But this does not work. As pointed out in (Giunchiglia et al., 2004, sect.7.2) expressing the reflexive transitive closure of a relation in the language of causal theories is not straightforward. For example, at first sight the causal laws

$$q(x, x) \Leftarrow \top, \quad q(x, y) \Leftarrow p(x, z) \wedge q(z, y), \quad \neg q(x, y) \Leftarrow \neg q(x, y)$$

may appear to express that  $q$  is the reflexive transitive closure of  $p$ . But they do not: there *is* a model of the causal theory in which  $q$  is the reflexive transitive closure of  $p$  but there are also other models in which  $q$ , though reflexive and transitive, contains more than the reflexive transitive closure of  $p$ . Consider  $\{p(1, 1), p(2, 2)\}$ . The reflexive transitive closure is  $\{q(1, 1), q(2, 2)\}$ , represented by the model  $\{q(1, 1), q(2, 2), \neg q(1, 2), \neg q(2, 1)\}$ , but the causal theory also has three other models:  $\{q(1, 1), q(2, 2), q(1, 2), \neg q(2, 1)\}$ ,  $\{q(1, 1), q(2, 2), \neg q(1, 2)\}$ ,  $\{q(1, 1), q(2, 2), \neg q(2, 1)\}$ ,

$q(2, 1)\}$ , and  $\{q(1, 1), q(2, 2), q(1, 2), q(2, 1)\}$ . It is perhaps easier to see why this is so by comparison with symmetric closure. The causal laws

$$q(x, y) \Leftarrow p(x, y), \quad q(x, y) \Leftarrow q(y, x), \quad \neg q(x, y) \Leftarrow \neg q(x, y)$$

similarly fail to express that  $q$  is the symmetric closure of  $p$ . In the example  $\{p(1, 1), p(2, 2)\}$ , the causal theory will contain the instances  $q(1, 2) \Leftarrow q(2, 1)$ ,  $q(2, 1) \Leftarrow q(1, 2)$ ,  $\neg q(1, 2) \Leftarrow \neg q(2, 1)$ , and  $\neg q(2, 1) \Leftarrow \neg q(1, 2)$  and it is easy to see that we get two models:  $\{q(1, 1), q(2, 2), \neg q(1, 2), \neg q(2, 1)\}$  as intended, but also  $\{q(1, 1), q(2, 2), q(1, 2), q(2, 1)\}$ .

It is not clear how this problem can be resolved. Since we view translations to causal theories as essentially an implementation issue, one possibility is to place the burden elsewhere in the implementation, that is to say, not in the component that constructs causal theories from  $(\mathcal{C}+)^+$  action descriptions but in the component that evaluates queries on the causal theories so constructed. To evaluate formulas of some query language on the ‘super-augmented’ ATS defined by a  $(\mathcal{C}+)^+$  action description  $D$ , we construct the causal theory  $\Gamma_m^{+D} = \Gamma_m^{D \cup \delta_D \cup \delta_{\sigma\gamma}}$  representing the (regular) ATS defined by  $D$ , which encodes the ‘counts as’ relation  $C_D^+$ , but arrange that the associated query evaluator evaluates formulas containing `counts_as $\mathcal{I}$`  expressions not on  $C_D^+$  but on the relation  $C_D^{++}$ . The process of constructing (the relevant instances of)  $C_D^{++}$  from (the relevant instances of)  $C_D^+$  is the job of the query evaluator, not part of the translation from  $D$  to  $\Gamma_m^{+D}$ .

This is not difficult to implement, but it is not satisfactory either. It would mean that `counts_as` constants in the  $(\mathcal{C}+)^+$  language do not behave in the same way, do not have the same properties as, the corresponding `counts_as $\mathcal{I}$`  expressions in the query language. One could say that this does not matter, that `counts_as` constants in  $(\mathcal{C}+)^+$  are used only to *define* the relations  $C_D^+$  and (indirectly)  $C_D^{++}$ , and that the semantics of  $(\mathcal{C}+)^+$  is unambiguous in this respect: one should not ascribe more properties to `counts_as` constants than are stated by the semantics. But still: it is only natural to assume that `counts_as` constants in the  $(\mathcal{C}+)^+$  language will behave in the same way as the relations they define, and it is at best confusing if it is otherwise. For instance, suppose we write  $a$  `counts_as`  $b$  if  $\top$  in a  $(\mathcal{C}+)^+$  action description, and also a fluent dynamic law  $\alpha$  `causes` ( $b$  `counts_as`  $c$ ) (say). It would be entirely reasonable to suppose, given the intended (stronger, transitive) reading of `counts_as`, that in any state after an  $\alpha$  transition,  $a$  `counts_as`  $c$  will hold. But it does not—this kind of inference is not supported by the language  $(\mathcal{C}+)^+$  itself. In the existing version, only the weaker ‘pseudo-transitivity’ is present: after an  $\alpha$  transition, every transition of type  $a$  will be a transition of type  $c$ , but this does not imply that the constant  $a$  `counts_as`  $c$  will be true. In the query language, in contrast,  $a$  `counts_as $\mathcal{I}$`   $c$  *would* be true, and this sort of mis-match is clearly to be avoided.

Implementation issues aside, there are in any case some other possible formulations of transitivity that can be considered as alternatives to  $(\mathcal{C}_{\text{trans}}^+)$ , and which remain to be investigated. Suppose for example that  $\alpha$  counts as  $\beta_1 \wedge \dots \wedge \beta_j \wedge \beta_{j+1} \wedge \dots \wedge \beta_n$ , and that  $\beta_1 \wedge \dots \wedge \beta_j$  counts as  $\alpha'$ . Would we want to conclude that  $\alpha$  counts as  $\alpha'$ ? An affirmative suggests we should look at a stronger type of transitivity, of the form

$$C_{\mathcal{I}}^+(s, X, Y) \text{ and } C_{\mathcal{I}}^+(s, Y', Z) \text{ implies } C_{\mathcal{I}}^+(s, X, Z), \text{ if } Y \subseteq Y'$$

or possibly, of the even stronger form

$$C_{\mathcal{I}}^+(s, X, Y) \text{ and } C_{\mathcal{I}}^+(s, Y', Z) \text{ implies } C_{\mathcal{I}}^+(s, X, Z), \text{ if } \tau_R(s, Y) \subseteq \tau_R(s, Y')$$

which depends on the *actual* transitions  $R$  of an ATS and not just on truth-functional properties of action formulas. These, and other forms, remain to be investigated, and are not included in the present version of  $(\mathcal{C}+)^+$ .

We do not see this as an important limitation of  $(\mathcal{C}+)^+$ . The significance of the very restricted form of reflexivity is negligible. Transitivity of ‘counts as’ is adopted in (Jones and Sergot, 1996) in the absence of convincing counter-examples to the contrary, but is not a critical feature of the characterisation adopted there, and neither is it here.

## 8.10 Defeasibility of ‘counts as’

There is one further remark that we wish to make, concerning the defeasibility of ‘counts as’. We can expect that in many practical applications the laws specifying ‘counts as’ relationships will be formulated most naturally as general default rules that are subject to implicit exceptions, exceptions to exceptions, and so on. We do not see anything particularly unusual or problematic about the formalisation of such rules in  $(\mathcal{C}+)^+$ .  $(\mathcal{C}+)^+$  inherits from  $\mathcal{C}+$  and the formalism of causal theories a general treatment of non-monotonic inference, and can be used to formulate defeasible general rule and exception structures using devices familiar from the literature on knowledge representation. (See (Giunchiglia et al., 2004, Sect. 7) for a comparison of causal theories with other well-known formalisations of non-monotonic reasoning.)

It may seem however that there is a different, more fundamental, kind of defeasibility inherent in the treatment of ‘counts as’, and this deserves some comment.

It is a cornerstone of the treatment given here of ‘counts as’ relations that when  $\alpha$  counts\_as  $\beta$  then any transition of type  $\alpha$  is also a transition of type  $\beta$ . This is expressed by the constraint  $(C_{\varepsilon}^+)$  required of an ATS. The effect is that when  $\alpha$  counts\_as  $\beta$  holds in a state  $s$  there can be no transitions of type  $\alpha \wedge \neg\beta$  from the state  $s$ . The same constraint is implicit in formalisations of ‘counts as’ in other action formalisms, such as (Jones and Sergot, 1996), where it emerges in discussion of ‘exercise of power’.

It might appear that  $(C_{\varepsilon}^+)$  as it stands is too strong. Suppose, for example, that  $\alpha$  counts\_as  $\beta$  and that  $\beta$  causes  $F$ .  $F$  holds in all states following the occurrence of a transition of type  $\beta$ ; and since all occurrences of type  $\alpha$  are also occurrences of type  $\beta$ , we get in effect that  $\alpha$  ‘causes’  $F$ . This is as intended: we want to be able to infer that actions  $\alpha$  have effects by virtue of being actions of type  $\alpha$  and that they also have effects by virtue of being counted, in institution  $\mathcal{I}$ , as actions of type  $\beta$ . By the same argument, when  $\alpha$  counts\_as  $\beta$  and  $\beta$  is non-executable, then  $\alpha$  is also non-executable. For if there are no transitions of type  $\beta$  in some given state  $s$ , and all transitions of type  $\alpha$  are also transitions of type  $\beta$ , then there can be no transitions of type  $\alpha$  in state  $s$  either.

But suppose now that *raise\_hand* counts\_as *start\_race*. And suppose further that *start\_race* is non-executable in some state—the runners are not ready, the track is too wet, or any number of other conditions hold. We must infer, as above, that in these circumstances *raise\_hand* is also non-executable. It seems that

whenever the race is unable to start, the starter is unable to raise his hand, in the sense that his action of doing so is non-executable.

Now there is nothing strange or problematic about this inference if the purpose of the action description is to act as a system *specification*. The fact that there are circumstances in which raising a hand is non-executable when we expect that it should be executable indicates that there is a flaw in the specification. The provision of computational tools to help identify such flaws is one of the main aims of the development.

But what if the purpose of the action description is to provide an accurate formal model of some real-world fragment? Just because a race is not ready to start we do not infer that the starter is (physically) unable to raise his hand. It might seem that some adjustment is necessary, and that this adjustment must concern condition  $(C_\varepsilon^+)$ . Perhaps the constraint that when  $\alpha$  counts<sub>as</sub>  $\beta$ , every transition of type  $\alpha$  is a transition of type  $\beta$  is too strong, and condition  $(C_\varepsilon^+)$  should be made *defeasible* and subject to exceptions in some way. This seems to be the position adopted by e.g. (Gelati et al., 2002, 2004) who have expressed the view (in a different formal framework) that the ‘counts as’ connection is inherently defeasible.

It seems to us that this is the wrong way of looking at the problem. It is not the condition  $(C_\varepsilon^+)$  that needs to be made defeasible but the formalisation of what counts as what. In the example, if starting the race is non-executable unless circumstances  $G$  hold, then it is simply not true to say that raising a hand counts as starting the race; raising a hand counts as starting the race if  $G$  holds, but not otherwise.

Schematically, we have  $\alpha$  counts<sub>as</sub>  $\beta$ ;  $\beta$  causes  $F$ , but when  $\neg F$  holds,  $\beta$  is non-executable, and so too is  $\alpha$ . However,  $\alpha$  counts<sub>as</sub>  $\beta$  is usually conditional on some other factual circumstances: schematically, the pattern is  $\alpha$  counts<sub>as</sub>  $\beta$  if  $G$  and  $\beta$  causes  $F$ . When  $\neg F$  holds,  $\beta$  is non-executable, and  $\alpha$  is non-executable when  $G \wedge \neg F$  holds. But it will often be the case that  $G \wedge \neg F$  can never be true because of other features of the specification or of the real-world fragment being modelled. One purpose of a formalisation is precisely to detect such features.

In summary: it is not the condition  $(C_\varepsilon^+)$  that needs to be made defeasible. Viewed as a kind of conditional,  $\alpha$  counts<sub>as</sub>  $\beta$  is not a *defeasible* conditional, though in practice the laws defining when  $\alpha$  counts<sub>as</sub>  $\beta$  holds will often be naturally formulated as defeasible general rules subject to exceptions.



## 9 The language $(\mathcal{C}+)^{++}$ : Permission

We now consider a second extension to the language  $\mathcal{C}+$ , to provide a means of specifying the ‘permitted’, ‘acceptable’, ‘ideal’, ‘legal’ states and the ‘permitted’, ‘acceptable’, ‘ideal’, ‘legal’ actions and transition paths. This extension is independent of the extension to accommodate ‘counts as’/conventional generation, but since we will usually want both features in the intended applications we present it here as an extension to the language  $(\mathcal{C}+)^+$  of the previous section. The semantic structure is a ‘coloured ATS’ which is a structure of the form:

$$\langle \sigma^f, \sigma^\gamma, S, \mathbf{A}, R, C_I^+, S_{\text{green}}, R_{\text{green}} \rangle$$

(or a ‘coloured LTS’ of the form

$$\langle \sigma^f, S, \mathbf{A}, R, S_{\text{green}}, R_{\text{green}} \rangle$$

if we do not need the ‘counts as’ features). The two new components are

- $S_{\text{green}} \subseteq S$ , the set of ‘permitted’ (‘acceptable’, ‘ideal’, ‘legal’) states—we call  $S_{\text{green}}$  the ‘green’ states of the system;
- $R_{\text{green}} \subseteq R$ , the set of ‘permitted’ (‘acceptable’, ‘ideal’, ‘legal’) transitions—we call  $R_{\text{green}}$  the ‘green’ transitions of the system.

We refer to the complements  $S - S_{\text{green}}$  and  $R - R_{\text{green}}$  as the ‘red states’ and ‘red transitions’, respectively. It is also possible to consider a more elaborate structure, of *partially coloured* transition systems in which states and transitions can be green, red, or uncoloured, but we shall not present that version here.

Note that  $R_{\text{green}} \subseteq R$  means that we build in the property that ‘permission’ implies ‘can’. We impose the following additional constraint:

- if  $(s, \varepsilon, s') \in R_{\text{green}}$  and  $s \in S_{\text{green}}$  then  $s' \in S_{\text{green}}$ .

which we refer to as the *green-green-green* constraint. Performing a permitted (green) action (transition) in a permitted (green) state must always lead to a permitted (green) state. All other possible combinations of green/red states and green/red transitions are allowed. In particular, and *contra* the assumptions underpinning John-Jules Meyer’s (1988) construction of ‘dynamic deontic logic’, a non-permitted (red) transition can result in a permitted (green) state. Similarly, it is easy to devise examples in which a permitted (green) transition can lead to a non-permitted (red) state. Some illustrations will arise in the examples to be considered later in the section. The only combination that cannot occur is the one eliminated by the ‘green-green-green’ constraint: a permitted (green) transition from a permitted (green) state cannot lead to a non-permitted (red) state.

Semantic devices such as  $S_{\text{green}}$  and  $R_{\text{green}}$  are familiar in the field of deontic logic. For example, Carmo and Jones (1996) employ a similar structure (though not for the purpose of constructing a language for defining them, as we attempt here) which has ideal/sub-ideal states and ideal/sub-ideal transitions (unlabelled). van der Meyden’s (1996) ‘dynamic logic of permission’ employs a structure in which transitions, but not states, are classified as ‘permitted/non-permitted’. van der Meyden’s version was constructed as a response to problems

of Meyer’s ‘dynamic deontic logic’ (Meyer, 1988) which classifies transitions as ‘permitted/non-permitted’ by reference only to the state resulting from a transition. ‘Deontic interpreted systems’ (Lomuscio and Sergot, 2003) classify states as ‘green’/‘red’, where these states have further internal structure to model the local states of agents in a multi-agent context. In all of these examples (and others) the task has been to find axiomatisations of such structures in one form of deontic logic or another. Here we are concerned with a different task, that of devising a language for *defining* coloured transition systems of the form described above.

## 9.1 The language $(\mathcal{C}+)^{++}$

The language  $(\mathcal{C}+)^{++}$  extends the language  $(\mathcal{C}+)^+$  with two new forms of laws.

- A *state permission law* (or rule) is an expression of the form

$$\text{not-permitted } F \tag{15}$$

where  $F$  is a fluent formula (i.e., a formula of signature  $\sigma^f$ ).

- An *action permission law* (or rule) is an expression of the form

$$\text{not-permitted } \alpha \text{ if } \psi \tag{16}$$

where  $\alpha$  is an action formula (i.e., a formula of signature  $\sigma^a$ ) and  $\psi$  is a formula of signature  $\sigma^f \cup \sigma^a$ . *not-permitted*  $\alpha$  is shorthand for *not-permitted*  $\alpha$  if  $\top$ .

It is also convenient to allow two variants of rule forms (15) and (16): *oblig*  $F$  is an abbreviation for *not-permitted*  $\neg F$  and *oblig*  $\alpha$  is an abbreviation for *not-permitted*  $\neg\alpha$ .

Informally, in the transition system defined by an action description  $D$ , a state  $s$  is red whenever  $s \models F$  for any state permission law *not-permitted*  $F$ . All other states are green by default. A transition  $(s, \varepsilon, s')$  is red whenever  $s \cup \varepsilon \models \psi$  and  $\varepsilon \models \alpha$  for any action permission law *not-permitted*  $\alpha$  if  $F$  after  $\psi$ . All other transitions are green, *subject to* the ‘green-green-green’ constraint which may impose further conditions on the possible colouring of a given transition. Examples will follow shortly to illustrate the effects of the ‘green-green-green’ constraint. It is possible to construct variants of the language with the defaults working the other way round (that is, a variant where the permission laws specify what is permitted (green), with other states and transitions not permitted (red) by default), or more elaborate forms allowing defaults to be mixed, but these all turn out to be more awkward and rather unnatural. We will stick to one basic form in this present account.

Let  $D$  be an action description of  $(\mathcal{C}+)^{++}$ .  $D_{\text{basic}}$  refers to the subset of laws of  $D$  that are also laws of  $(\mathcal{C}+)^+$  (and hence also laws of  $\mathcal{C}+$ ). The coloured ATS defined by  $D$  has the augmented states  $S$  and transitions  $R$  that are defined by its  $(\mathcal{C}+)^+$  component  $D_{\text{basic}}$ . In the case where we are not using the ‘counts as’ features, the coloured LTS defined by  $D$  has the states  $S$  and transitions  $R$  of the labelled transition system defined by the  $\mathcal{C}+$  action description  $D_{\text{basic}}$ . In both cases, coloured ATS and coloured LTS, the green states  $S_{\text{green}}$  and green transitions  $R_{\text{green}}$  are defined as follows:

$$S_{\text{green}} =_{\text{def}} S - S_{\text{red}} \quad R_{\text{green}} =_{\text{def}} R - R_{\text{red}}$$

where

$$\begin{aligned}
S_{\text{red}} &=_{\text{def}} \{s \mid s \models F \text{ for some state permission law not-permitted } F \text{ in } D\} \\
R_{\text{red}} &=_{\text{def}} \{(s, \varepsilon, s') \mid s \cup \varepsilon \models \psi, \varepsilon \models \alpha, s' \models F \text{ for some action} \\
&\quad \text{permission law not-permitted } \alpha \text{ if } F \text{ after } \psi \text{ in } D\} \\
&\quad \cup \{(s, \varepsilon, s') \mid s \in S_{\text{green}} \text{ and } s' \notin S_{\text{green}}\}
\end{aligned}$$

The last component of the  $R_{\text{red}}$  definition ensures that the ‘green-green-green’ constraint is satisfied.

## 9.2 Example

Suppose it is forbidden for a man and a woman to be alone together in a room. Let  $a$  and  $b$  be men and  $c$  be a woman, and suppose they inhabit a world in which there is one dwelling consisting of two rooms with a connecting internal door and one external door, as depicted in Figure 9.1. For the sake of simplicity, suppose further that exactly one of the three persons  $a$ ,  $b$  and  $c$  moves from one room to another at any one time, and ignore the possibility that  $a$ ,  $b$  or  $c$  move out of the dwelling or other persons move in from the outside. These simplifications are imposed merely to keep the diagrams of the resulting transition system manageable. The example works just as well without the simplifications but the diagrams become too large to exhibit here.

Let fluents  $loc(x) = p$  represent the location of person  $x$ , for  $x$  ranging over  $a$ ,  $b$ , and  $c$ , and  $p$  ranging over the locations ‘left’ and ‘right’ (say). Let  $move(x) = p$  represent the action in which  $x$  moves to location  $p$ . The action description for the example specifies that the  $loc(x)$  fluents are inertial, that the  $move(x)$  actions are exogenous, that exactly one  $move(x)$  action takes place at any transition, that  $move(x) = p$  is nonexecutable when  $loc(x) = p$ , and that  $move(x) = p$  ‘causes’  $loc(x) = p$ . These parts are just as for the language  $\mathcal{C}+$ ; the details are straightforward and omitted. The state permission laws for the example may be stated in various ways, of which one simple formulation is as follows (where  $p$  ranges over the two locations ‘left’ and ‘right’):

$$\begin{aligned}
&\text{not-permitted } loc(a) = p \wedge loc(c) = p \wedge \neg loc(b) = p \\
&\text{not-permitted } loc(b) = p \wedge loc(c) = p \wedge \neg loc(a) = p
\end{aligned}$$

There are more general and more elegant ways of formulating the required laws; this simple form is adequate for present purposes and does not distract from the main point of the example.

The transition system defined by this action description is shown in Figure 9.1, first with just the states coloured as determined by the state permission laws, and then with the colours of the transitions shown as determined by the ‘green-green-green’ constraint.

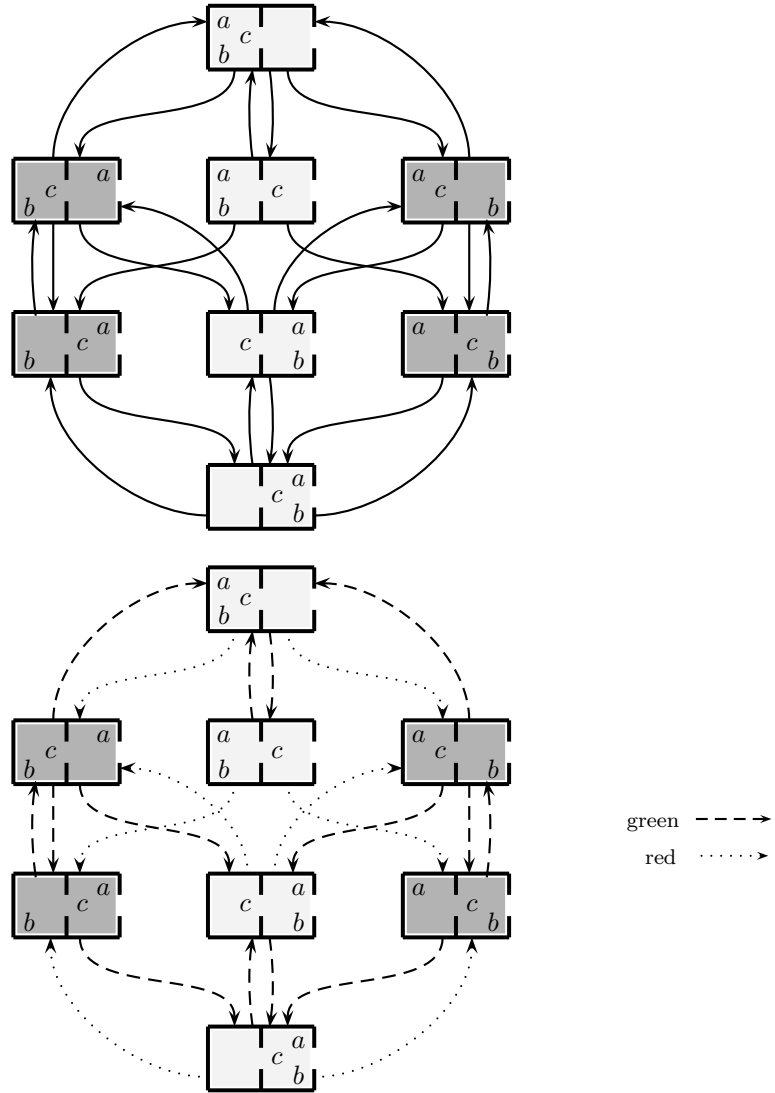


Figure 9.1: The transition system for the connecting rooms example. The top picture shows the colouring of the states: the darker shade indicates the red states as determined by the state permission laws, the others are green by default. The bottom diagram shows the colouring of the transitions, as determined by the 'green-green-green' constraint.

### 9.3 Translation to causal theories

The translation of an action description  $D$  in  $(\mathcal{C}+)^{++}$  to a causal theory  $\Gamma_m^D$  is exactly as for the translation of  $(\mathcal{C}+)^+$  theories, with the following extension to deal with the state and action permission laws.

Let **status** and **trans** be distinguished fluent and action constants, respectively, in the signature of  $\Gamma_m^D$ , both with possible values **green** and **red**. They will be used to represent the colour of a state and the colour of a transition, respectively.

For every state permission law **not-permitted**  $F$  and time index  $i \in 0..m$ , include in  $\Gamma_m^D$  a causal rule of the form

$$\text{status}[i] = \text{red} \Leftarrow F[i] \quad (17)$$

and for every  $i \in 0..m$ , a causal rule of the form

$$\text{status}[i] = \text{green} \Leftarrow \text{status}[i] = \text{green} \quad (18)$$

to specify the default colour of a state. A state permission law of the form **oblig**  $F$  produces causal rules of the form

$$\text{status}[i] = \text{red} \Leftarrow \neg F[i]$$

For every action permission law **not-permitted**  $\alpha$  if  $F$  after  $\psi$  and time index  $i \in 0..m-1$ , include in  $\Gamma_m^D$  a causal rule of the form

$$\text{trans}[i] = \text{red} \Leftarrow F[i+1] \wedge \alpha[i] \wedge \psi[i] \quad (19)$$

and for every  $i \in 0..m-1$ , a causal rule of the form

$$\text{trans}[i] = \text{green} \Leftarrow \text{trans}[i] = \text{green} \quad (20)$$

to specify the default colour of a transition. An action permission law of the form **oblig**  $\alpha$  if  $F$  after  $\psi$  produces causal rules of the form

$$\text{trans}[i] = \text{red} \Leftarrow F[i+1] \wedge \neg \alpha[i] \wedge \psi[i]$$

Finally, to capture the ‘green-green-green’ constraint, include for every  $i \in 0..m-1$  a causal rule of the form

$$\text{trans}[i] = \text{red} \Leftarrow \text{status}[i] = \text{green} \wedge \text{status}[i+1] = \text{red} \quad (21)$$

It is straightforward to show that models of the causal theory  $\Gamma_m^D$  correspond to all paths of length  $m$  through the coloured transition system defined by  $D$ , where the fluent constant **status** and the action constant **trans** encode the colours of the states and transitions, respectively.

Note that when an action description  $D$  in  $(\mathcal{C}+)^{++}$  is definite, it translates to a causal theory  $\Gamma_m^D$  which is also definite.

Causal rules (17), (18) and (20) can be expressed in the language  $\mathcal{C}+$ , that is to say, state permission laws and the default colouring of states and transitions can be expressed as laws in  $\mathcal{C}+$ , with **status** and **trans** treated as ordinary fluent and action constants. The causal rules (19) and (21), however, corresponding to action permission laws and the green-green-green constraint, cannot be expressed

in  $\mathcal{C}+$ . The form of these rules, where there is a time index in the antecedent greater than the time index in the consequent, conflicts with the very strong reading of ‘causes’ that underpins the language  $\mathcal{C}+$ . Rules corresponding to (19) and (21) are deliberately excluded from the language  $\mathcal{C}+$ , by design.  $\mathcal{C}+$  could express a restricted form of action permission law, of the form

$$\text{not-permitted } \alpha \text{ if } \top \text{ after } \psi \quad (22)$$

since this translates to the causal rule

$$\text{trans}[i] = \text{red} \Leftarrow \alpha[i] \wedge \psi[i]$$

which is the translated form of an action dynamic law of  $\mathcal{C}+$ . However, there is no similar simplification that allows the ‘green-green-green’ constraint to be expressed satisfactorily as a law in  $\mathcal{C}+$ .  $\mathcal{C}+$  can express:

$$\perp \text{ if status} = \text{red after trans} = \text{green} \wedge \text{status} = \text{green}$$

which translates to causal rules:

$$\perp \Leftarrow \text{status}[i+1] = \text{red} \wedge \text{trans}[i] = \text{green} \wedge \text{status}[i] = \text{green} \quad (23)$$

and it can express

$$\text{status} = \text{green after trans} = \text{green} \wedge \text{status} = \text{green}$$

which translates to causal rules:

$$\text{status}[i+1] = \text{green} \Leftarrow \text{trans}[i] = \text{green} \wedge \text{status}[i] = \text{green} \quad (24)$$

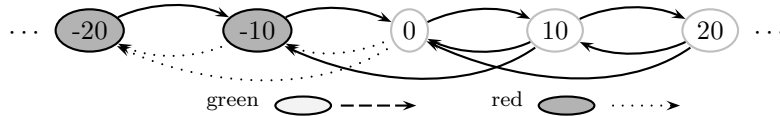
but both of these give quite the wrong effect. Both can be interpreted as defining transition systems. But although these transition systems satisfy the green-green-green constraint, they give the wrong interaction between state and action permission laws and the associated defaults.

The following example provides an illustration.

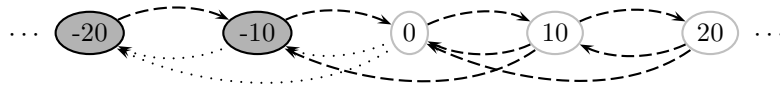
#### 9.4 The ‘green-green-green’ constraint: Example

Consider a simple example of a bank (with one customer): states of the bank are represented by the value of a single fluent constant *balance* which represents the current balance of the customer’s account. Suppose values of *balance* are multiples (including negative multiples) of 10 euro. There are three actions that the customer can perform: withdraw 10 euro, withdraw 20 euro, deposit 10 euro. A state is forbidden (not permitted, red) if the balance is less than 0. Suppose further, for the sake of the example, that a withdrawal (of any amount) is forbidden (red) when the balance is zero or negative. A suitable action description for this example in  $(\mathcal{C}+)^{++}$  is very straightforward; we omit the details.

The transition system, colouring first only the states and transitions that are explicitly specified to be red, is as follows:

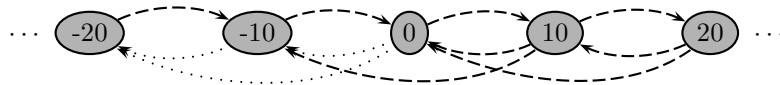


The labels on transitions are obvious from context so we omit them for clarity. Suppose first that we apply the defaults and the green-green-green constraint in a different order to that specified in Section 9.3 above: suppose we first assume that all transitions not explicitly forbidden (red) are permitted (green), then apply the green-green-green constraint, and only then the default that states not red are green. At the first step we obtain the following partially coloured transition system:



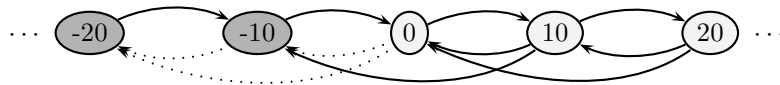
Now applying the green-green-green constraint forces us to conclude that the state 10 is red. This is because the transition from 10 to  $-10$  is green, and the state  $-10$  is red, so the state 10 cannot be green; it must be red.

Furthermore, state 0 must now also be red, because there is a green transition from 0 to a red state, 10. And the state 20 must be red too, because there is a green transition from a red state, 10, to 20. And so on, for all states greater than 20. We are left with the following picture.

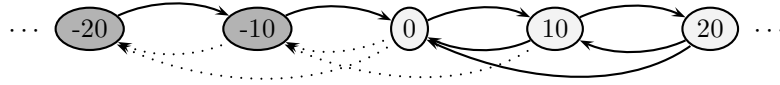


The result is surely quite counter-intuitive: we conclude that all balances, positive as well as negative, are not permitted, which is surely wrong.

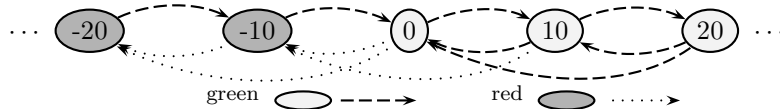
Contrast this with what is obtained by treating constraints and defaults in the manner specified in Section 9.3. First, states that are not explicitly red are green by default:



Now apply the green-green-green constraint: the transition from 10 to  $-10$  must be red because 10 is a permitted (green) state while  $-10$  is a forbidden (red) one.



Finally all other transitions are green by default, leaving the transition system coloured as follows.



And surely this is what we would expect: 10, 20, ... are all permitted states, and it is the transition from 10 to -10, i.e., the withdrawal of 20 euro when the balance is at 10 euro, that is not permitted.

Notice also that in this example, *contra* Meyer's definition of forbidden actions in 'dynamic deontic logic' (Meyer, 1988), there is a green transition (from -20 to -10) which results in a red state. And it is surely right that this transition should be green: it is surely permitted, given the permission laws stated above and the declared reading of defaults, that the customer is permitted to deposit 10 euro even though this still leaves his account in the negative. Although not shown in this example, it is also easy to find cases where a forbidden (red) transition can lead to a green (permitted) state. Suppose, for instance, that the bank decides to specify that withdrawals of less than 20 euro are no longer permitted (because of the administrative inconvenience, say): such withdrawals are still possible (they are not 'nonexecutable') but from the bank's point of view they are undesirable, that is to say, not permitted ('red').

Or: suppose now there are two customers, *custA* and *custB*. On Meyer's account, whenever *custA*'s account balance is negative, it is forbidden for *custB* to deposit or withdraw from his own account. This is surely wrong.

## 9.5 Relaxation of syntax

Permission laws of the form (15) and (16) are convenient but rather restrictive. For more flexibility, the  $(\mathcal{C}+)^{++}$  language also allows distinguished fluent and action constants **status** and **trans** to be used explicitly in formulas and causal laws. The atoms **status = red** and **trans = red** can then be regarded as instances of what are sometimes called 'violation constants' in deontic logic. It is also easy to allow more 'shades' of red and green to allow different notions of permitted/legal/acceptable to be mixed. We will not employ that device in the examples discussed in this paper.

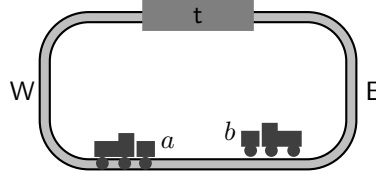
## 9.6 Example (trains)

The following example is used in (van der Hoek et al., 2004; Jamroga et al., 2004) to illustrate the use of alternating-time logic (ATL) (Alur et al., 2002) for determining the effectiveness of 'social laws' designed to co-ordinate the actions of agents in multi-agent system. The term is from (Shoham and Tennenholtz, 1992a,b, 1995; Moses and Tennenholtz, 1995) We will use the example for a different purpose, to illustrate how the permission component of  $(\mathcal{C}+)^{++}$  can be used to analyse variants of the example in which agents may fail to obey social laws. More elaborate versions of the example, in which additional components are introduced in order to enforce compliance with the 'social laws', are discussed in (Sergot, 2005).

There are two trains, *a* and *b*, with *a* running clockwise round a double track, and *b* running anti-clockwise. There is a tunnel in which the double track becomes a single track. If the trains are both inside the tunnel at the same time they will collide. The tunnel can thus be seen as a kind of critical section, or as a resource for which the trains must compete.



There are obviously many ways in which the example can be formulated. The following will suffice for present purposes. Let simple fluent constants  $loc(a)$  and  $loc(b)$  represent the locations of trains  $a$  and  $b$  respectively. They both have possible values  $\{W, t, E\}$ . For action constants, we take  $a$  and  $b$  with possible values  $\{go, stay\}$ . (Action constants  $act(a)$  and  $act(b)$  may be easier to read but we choose  $a$  and  $b$  for brevity.) Both are exogenous.



The  $\mathcal{C}+$  action description representing the possible movements of the trains is as follows.

$$\text{inertial } loc(a), loc(b) \quad (25)$$

train  $a$  moves clockwise:

$$\begin{aligned} a = go & \text{ causes } loc(a) = t \text{ if } loc(a) = W \\ a = go & \text{ causes } loc(a) = E \text{ if } loc(a) = t \\ a = go & \text{ causes } loc(a) = W \text{ if } loc(a) = E \end{aligned} \quad (26)$$

train  $b$  moves anti-clockwise:

$$\begin{aligned} b = go & \text{ causes } loc(b) = t \text{ if } loc(b) = E \\ b = go & \text{ causes } loc(b) = W \text{ if } loc(b) = t \\ b = go & \text{ causes } loc(b) = E \text{ if } loc(b) = W \end{aligned} \quad (27)$$

collisions:

$$\begin{aligned} collision & \text{ iff } loc(a) = t \wedge loc(b) = t \quad \% \text{ for convenience} \\ \text{nonexecutable } a = go & \text{ if } collision \\ \text{nonexecutable } b = go & \text{ if } collision \end{aligned} \quad (28)$$

$collision$  is a statically determined Boolean fluent constant, introduced for convenience. Here and in the rest of this section,  $F$  iff  $G$  is used as shorthand for the pair of laws  $F$  if  $G$  and  $\text{default } \neg F$ .

Suppose, for the sake of an example, that we impose additional norms (social laws), as follows: no train is permitted to enter the tunnel unless the other train has just emerged. (We assume that this will be observed by the train that is preparing to enter. One of the variants of the example in (Sergot, 2005) introduces a communication mechanism—a semaphore—to allow the trains to communicate their positions to each other.) Will such a law be effective in avoiding collisions?

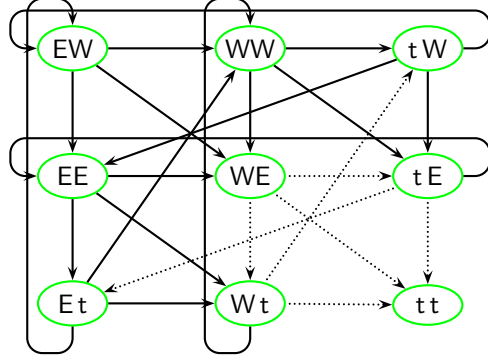


Figure 9.2: Coloured transition system for the trains example. A state label such as EW is short for  $\{loc(a) = E, loc(b) = W\}$ . Transition labels are omitted for clarity. Horizontal edges are transitions in which train  $a$  moves and  $b$  does not. Vertical edges are transitions in which train  $b$  moves and  $a$  does not. Diagonal edges are transitions in which both trains move. Reflexive edges, corresponding to transitions in which neither train moves, are omitted from the diagram to reduce clutter. Dotted lines indicate red transitions. All states and all other transitions, including the omitted reflexive ones, are green.

First, it is convenient to define the following auxiliary action constants (all Boolean and non-exogenous):

$$\begin{aligned}
 enter(a) &\text{ iff } a = go \wedge loc(a) = W \\
 exit(a) &\text{ iff } a = go \wedge loc(a) = t \\
 enter(b) &\text{ iff } b = go \wedge loc(b) = E \\
 exit(b) &\text{ iff } b = go \wedge loc(b) = t
 \end{aligned} \tag{29}$$

Again, these are introduced merely for convenience; the example can be constructed easily enough without them. Now we formulate the social laws:

$$\begin{aligned}
 \text{not-permitted } enter(a) &\text{ if } loc(b) \neq W \\
 \text{not-permitted } enter(b) &\text{ if } loc(a) \neq E
 \end{aligned} \tag{30}$$

Let  $D_{\text{trains}}$  be the  $(\mathcal{C}+)^{++}$  action description consisting of laws (25)–(30). The coloured transition system defined by  $D_{\text{trains}}$  is shown in Figure 9.2.

How do we test the effectiveness of the social laws (30)? Since the causal theory  $\Gamma_1^{D_{\text{trains}}}$  encodes the transitions defined by  $D_{\text{trains}}$ , the following captures the property that if both trains comply with the social laws, no collisions will occur.

$$comp(\Gamma_1^{D_{\text{trains}}}) \models \neg collision[0] \wedge trans[0] = \text{green} \rightarrow \neg collision[1] \tag{31}$$

This can be checked, as in C<sub>2</sub>ALC, by using a standard sat-solver to determine that the formula  $comp(\Gamma_1^{D_{\text{trains}}}) \cup \{\neg collision[0] \wedge trans[0] = \text{green} \wedge collision[1]\}$  is not satisfiable. The property (31) is equivalently expressed as:

$$comp(\Gamma_1^{D_{\text{trains}}}) \models \neg collision[0] \wedge collision[1] \rightarrow trans[0] = \text{red} \tag{32}$$

which says that a collision occurs only following a transition in which either one train or both violate the norms.

Notice that  $\text{comp}(\Gamma_1^{D_{\text{trains}}}) \not\models \text{trans}[0] = \text{green} \rightarrow \neg \text{collision}[1]$ : as formulated by  $D_{\text{trains}}$ , the transition from a collision state to itself is green.

One major advantage of taking  $\mathcal{C}+$  as the basic action formalism, compared to other action formalisms in AI, is its explicit transition system semantics. This enables a wide range of other analytical techniques to be applied. In particular, system properties can be expressed in the branching time temporal logic CTL and verified on the transition system defined by a  $\mathcal{C}+$  or  $(\mathcal{C}+)^{++}$  action description using standard model checking systems. (See e.g. (Clarke et al., 2000).)

We will say that a formula  $\varphi$  of CTL is *valid* on a (coloured) transition system  $\langle S, I(\sigma^a), R, S_{\text{green}}, R_{\text{green}} \rangle$  defined by  $(\mathcal{C}+)^{++}$  action description  $D$  when  $s \cup \varepsilon \models \varphi$  for every  $s \cup \varepsilon$  such that  $(s, \varepsilon, s') \in R$  for some state  $s'$ . The definition is quite standard, except for a small adjustment to allow action constants in  $\varphi$  to continue to be evaluated on transition labels  $\varepsilon$ . (And we do not distinguish any particular set of initial states; all sets in  $S$  are initial states.) We will also say in that case that formula  $\varphi$  is valid on the action description  $D$ .

In CTL, the formula  $\text{AX} \varphi$  expresses that  $\varphi$  is satisfied in the next state in all future branching paths from now.<sup>5</sup>  $\text{EX}$  is the dual of  $\text{AX}$ :  $\text{EX} \varphi \equiv \neg \text{AX} \neg \varphi$ .  $\text{EX} \varphi$  expresses that  $\varphi$  is satisfied in the next state of some future branching path from now. The properties (31) and (32) can thus be expressed in CTL as follows:

$$\neg \text{collision} \wedge \text{trans} = \text{green} \rightarrow \text{AX} \neg \text{collision} \quad (33)$$

or equivalently  $\neg \text{collision} \wedge \text{EX} \text{collision} \rightarrow \text{trans} = \text{red}$ . It is easily verified by reference to Figure 9.2 that these formulas are valid on the action description  $D_{\text{trains}}$ . Also valid is the CTL formula  $\text{EX} \text{trans} = \text{green}$  which expresses that there is always a permitted action for both trains. This is true even in collision states, since the only available transition is then the one where both trains remain idle, and that transition is green. The CTL formula  $\text{EF} \text{collision}$  is also valid on  $D_{\text{trains}}$ , signifying that in every state there is at least one path from then on with *collision* true somewhere in the future.<sup>6</sup>

(Sergot, 2005) goes on to consider more elaborate versions of the example: in general, we want to be able to verify formally whether the introduction of additional control mechanisms—additional controller agents, communication devices, restrictions on agents' possible actions—are effective in ensuring that agents comply with the norms ('social laws') that govern their behaviour. These more elaborate versions can be analysed in similar fashion.

---

<sup>5</sup>  $s_0 \cup \varepsilon_0 \models \text{AX} \varphi$  if for every infinite path  $s_0 \varepsilon_0 s_1 \varepsilon_1 \dots$  we have that  $s_1 \cup \varepsilon_1 \models \varphi$ .

<sup>6</sup>  $s_0 \cup \varepsilon_0 \models \text{EF} \varphi$  if there is an (infinite) path  $s_0 \varepsilon_0 \dots s_m \varepsilon_m \dots$  with  $s_m \cup \varepsilon_m \models \varphi$  for some  $m \geq 0$ .

## 10 Conclusion

The language(s)  $(\mathcal{C}+)^{++}$  provide an executable formalism for representing the institutional aspects of agent societies (specifically, the ‘counts as’ relation between actions) and a simple treatment of permission for expressing norms or ‘social laws’ that govern their behaviour.

Definite, atom-conjunctive action descriptions of  $(\mathcal{C}+)^+$ , which are the action descriptions of practical interest, can be implemented straightforwardly by translation to definite action descriptions of  $\mathcal{C}+$ . We thus have direct routes to implementation, using the established computational techniques for  $\mathcal{C}+$ , such as *CCALC* and translations to logic programs. And although in the paper we only illustrated the use of CTL and model checking techniques in connection with the permission component of  $(\mathcal{C}+)^{++}$ , it is clear that they can also be used to verify system properties of action descriptions of  $(\mathcal{C}+)^+$ , such as the auction protocol presented in Section 8.7.

The permission component of  $(\mathcal{C}+)^{++}$  provides a means of distinguishing between permitted, acceptable, desirable (‘green’) states and transitions. It can be used as an extension of  $\mathcal{C}+$ , or in combination with the ‘counts as’ features, as an extension of  $(\mathcal{C}+)^+$ . Definite action descriptions of  $(\mathcal{C}+)^{++}$  cannot be translated to  $\mathcal{C}+$  because of the ‘green-green-green constraint’ but they can be translated to (definite) non-monotonic causal theories, requiring only very minor modification to existing computational tools such as *CCALC*. They can also be used with model checking systems, as illustrated briefly with the ‘trains’ example in Section 9.6. In addition to verifying system properties that hold if all agents/system components behave in accordance with norms/‘social laws’, we are able to analyse system properties that hold when agents/system components fail to comply with norms, and to analyse formally the effectiveness of introducing additional control, enforcement, and recovery mechanisms. This has been demonstrated on some small examples (see e.g. (Sergot, 2005)) but there is no reason to think that the techniques would not scale up to realistic examples, within the limits of existing model checking technology.

The semantical device employed in  $(\mathcal{C}+)^{++}$ , of partitioning states, and here also transitions, into ‘ideal’ (‘green’) and ‘sub-ideal’ (‘red’), is familiar in the field of deontic logic. It is essentially the same device as is employed in ‘Standard Deontic Logic’ (SDL). It has very well known inadequacies, particularly in regard to the representation of ‘contrary-to-duty’ structures. To some extent these problems are mitigated in  $(\mathcal{C}+)^{++}$  because it is a richer representational formalism than SDL and because it can deal with temporal aspects of ‘contrary-to-duty’ where these are present. However, there are benchmark problems in the deontic logic literature that will not receive adequate treatment in  $(\mathcal{C}+)^{++}$ . A detailed evaluation of the permission component of  $(\mathcal{C}+)^{++}$ , and a discussion of further possible refinements, is outside the scope of this paper and will be covered separately elsewhere.

The language(s)  $(\mathcal{C}+)^{++}$  inherit from  $\mathcal{C}+$  a number of very desirable features. Besides implementation methods, and support for concurrent actions, non-determinism, and indirect effects of actions, the underlying formalism of causal theories provides a general purpose non-monotonic framework for expressing defeasible general rules and defeasible effects of actions.

The language(s)  $(\mathcal{C}+)^{++}$  also inherit a number of important limitations. One limitation is that causal rules can refer only to states and transitions at most one time-step away from each other. This makes it impossible to refer directly to the values of fluents further in the past, or to features of earlier transitions. In the auction example of Section 8.7, for example, it was necessary to introduce additional fluents to represent the last bid made in order to encode a simple revocation mechanism. More complicated examples soon become unwieldy. We might wish to say, for example, that three attempted invalid bids by an agent in an auction are permitted, but a fourth is not. Relying on additional fluents to record the relevant fragment of history not only increases the number of states (exponentially) but also tends to make the representation obscure. It is also impossible to say, directly, that the execution of one action causes the execution of another action at the next time-step, or at some other time in the future. These are limitations of  $\mathcal{C}+$  rather than of the causal theories into which  $\mathcal{C}+$  is translated: Craven and Sergot (2005) present a generalisation of  $\mathcal{C}+$  which removes this restriction, though it provides only a partial solution and the corresponding  $(\mathcal{C}+)^{++}$  extensions remain to be developed.

Other limitations of  $\mathcal{C}+$  and  $(\mathcal{C}+)^{++}$  are limitations of transition systems generally. States and transitions are global, and so lack the structure required to model large numbers of interacting components, such as multi-agent systems. Transitions have no (explicit) duration, and there is no support for time metrics: encoding delays and deadlines in  $\mathcal{C}+$  is awkward at best. Naturally there are other structures in computer science that have been designed to overcome the limitations of transition systems. One can imagine devising similar representational formalisms based on these other structures. The details are far from trivial of course.

Besides addressing these various issues, current development of  $(\mathcal{C}+)^{++}$  includes the following topics:

- technical investigation of logical properties of  $\mathcal{C}+$  and of  $(\mathcal{C}+)^{++}$  (see (Sergot and Craven, 2005) for a preliminary study);
- experiments with the use of the CCALC-based implementation of  $(\mathcal{C}+)^{++}$  as the reasoning engine for the ‘society simulator’ developed by Alexander Artikis (Artikis et al., 2002) in place of the event calculus and  $\mathcal{C}+$  currently employed;
- refinement of the representational formalism to provide more structure in states, in two respects:
  - support for multiple institutions (multiple ‘counts as’ relations)
  - the introduction of local states, and local actions, for agents (and environment) for modelling multi-agent systems;
- development of the  $\mathcal{EC}+$  implementation described in (Craven and Sergot, 2003) which provides an event calculus style of computation for queries on (a subset of)  $\mathcal{C}+$  action descriptions;
- investigations of the use of  $\mathcal{C}+$  and  $(\mathcal{C}+)^{++}$  as input languages for model checking systems, specifically NuSMV (Cimatti et al., 2002).

## Acknowledgements

A preliminary version of  $(\mathcal{C}+)^{++}$  was produced in the course of the EU funded FET project ALFEBIITE (IST-1999-10298). The author wishes to express his thanks to Alex Artikis and Rob Craven for many valuable discussions and suggestions, to Joohyung Lee for helpful clarifications of features of the  $\mathcal{C}+$  language, and to Jeremy Pitt for his patience.

## References

- V. Akman, S. T. Erdoğan, J. Lee, V. Lifschitz, and H. Turner. Representing the Zoo World and the Traffic World in the language of the Causal Calculator. *Artificial Intelligence*, 153:105–140, 2004.
- R. Alur, T.A. Henzinger, and O. Kupferman. Alternating-time temporal logic. *Journal of the ACM*, 49(5):672–713, September 2002.
- A. Artikis, J. Pitt, and M. J. Sergot. Animated specification of computational societies. In Cristiano Castelfranchi and W. Lewis Johnson, editors, *Proc. 1st International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS'02), Bologna*, pages 1053–1062. ACM Press, July 2002.
- A. Artikis, M. J. Sergot, and J. Pitt. An executable specification of an argumentation protocol. In *Proc. 9th International Conference on Artificial Intelligence and Law (ICAIL'03), Edinburgh*, pages 1–11. ACM Press, 2003a.
- A. Artikis, M. J. Sergot, and J. Pitt. Specifying electronic societies with the Causal Calculator. In Fausto Giunchiglia, James Odell, and Gerhard Weiss, editors, *Agent-Oriented Software Engineering III. Proc. 3rd International Workshop (AOSE 2002), Bologna, July 2002*, LNCS 2585, pages 1–15. Springer, 2003b.
- José Carmo and Andrew J. I. Jones. Deontic database constraints, violation and recovery. *Studia Logica*, 57(1):139–165, 1996.
- A. Cimatti, E. Clarke, E. Giunchiglia, F. Giunchiglia, M. Pistore, M. Roveri, R. Sebastiani, and A. Tacchella. NuSMV Version 2: An OpenSource Tool for Symbolic Model Checking. In *Proc. International Conference on Computer-Aided Verification (CAV 2002), Copenhagen, July 2002*, LNCS 2404. Springer, 2002. See <http://nusmv.irst.itc.it>.
- E. M. Clarke, O. Grumberg, and D. A. Peled. *Model Checking*. MIT Press, Cambridge, 2000.
- R. A. Craven and M. J. Sergot.  $\mathcal{EC}+$ —efficient computation of queries for the language  $\mathcal{C}+$ . Technical report, Department of Computing, Imperial College, London, UK, 2003.
- Robert Craven and Marek Sergot. Distant causation in  $\mathcal{C}+$ . *Studia Logica*, 79(1):73–96, February 2005.
- Thomas Eiter, Wolfgang Faber, Nicola Leone, Gerald Pfeifer, and Axel Polleres. System description: The DLV-K planning system. In Thomas Eiter, Wolfgang Faber, and Mirosław Truszczyński, editors, *Proc. 6th International Conference on Logic Programming and Nonmonotonic Reasoning*, LNAI 2173, pages 413–416. Springer-Verlag, 2001.
- J. Gelati, G. Governatori, A. Rotolo, and G. Sartor. Declarative power, repre-

- sentation, and mandate. A formal analysis. In *Proceedings of Conference on Legal Knowledge and Information Systems (JURIX)*, 2002.
- Jonathan Gelati, Antonino Rotolo, Giovanni Sartor, and Guido Governatori. Normative autonomy and normative co-ordination: Declarative power, representation and mandate. *Artificial Intelligence and Law*, 12:53–81, 2004.
- Michael Gelfond and Vladimir Lifschitz. Action Languages. *Electronic Transactions on AI*, 3(16), 1998. <http://www.ep.lin.se/ea/cis/1998/016>.
- Enrico Giunchiglia, Joohyung Lee, Vladimir Lifschitz, Norman McCain, and Hudson Turner. Nonmonotonic causal theories. *Artificial Intelligence*, 153: 49–104, 2004.
- Enrico Giunchiglia, Joohyung Lee, Vladimir Lifschitz, and Hudson Turner. Causal laws and multi-valued fluents. In *Proc. of the Fourth Workshop on Nonmonotonic Reasoning, Action, and Change*, Seattle, August 2001.
- Enrico Giunchiglia and Vladimir Lifschitz. An action language based on causal explanation: Preliminary report. In *Proc. AAAI-98*, pages 623–630. AAAI Press, 1998.
- Enrico Giunchiglia and Vladimir Lifschitz. Action languages, temporal action logics, and the situation calculus. In *Working Notes of the IJCAI-99 Workshop on Nonmonotonic Reasoning, Action, and Change*, 1999.
- A. I. Goldman. *A Theory of Human Action*. Prentice-Hall, New Jersey, 1970.
- W. Jamroga, W. van der Hoek, and M. Wooldridge. On obligations and abilities. In Alessio Lomuscio and Donald Nute, editors, *Proc. 7th International Workshop on Deontic Logic in Computer Science (DEON'04)*, Madeira, May 2004, LNAI 3065, pages 165–181. Springer, 2004.
- A. J. I. Jones and M. J. Sergot. A formal characterisation of institutionalised power. *Journal of the IGPL*, 4(3):429–445, 1996. Reprinted as (Jones and Sergot, 1997).
- A. J. I. Jones and M. J. Sergot. A formal characterisation of institutionalised power. In E. G. Valdés, W. Krawietz, G. H. von Wright, and R. Zimmerling, editors, *Normative Systems in Legal and Moral Theory. Festschrift for Carlos E. Alchourrón and Eugenio Bulygin*, pages 349–367. Duncker & Humboldt, Berlin, 1997.
- Vladimir Lifschitz. On the logic of causal explanation. *Artificial Intelligence*, 96:451–465, 1997.
- Vladimir Lifschitz and Hudson Turner. Representing transition systems by logic programs. In *Proc. Fifth International Conference on Logic Programming and Nonmonotonic Reasoning*, pages 92–106, 1999.
- A. Lomuscio and M. J. Sergot. Deontic interpreted systems. *Studia Logica*, 75 (1):63–92, October 2003.
- Norman McCain and Hudson Turner. Causal theories of action and change. In *Proc. AAAI-97, Providence, RI*, pages 460–465. AAAI Press, 1997.
- R. van der Meyden. The dynamic logic of permission. *Journal of Logic and Computation*, 6(3):465–479, 1996.
- John-Jules Ch. Meyer. A different approach to deontic logic: Deontic logic viewed as a variant of dynamic logic. *Notre Dame Journal of Formal Logic*, 29(1):109–136, 1988.
- Y. Moses and M. Tennenholtz. Artificial social systems. *Computers and Artificial Intelligence*, 14(6):533–562, 1995.

- Ilkka Niemelä. Logic programs with stable model semantics as a constraint logic programming paradigm. *Annals of Mathematics and Artificial Intelligence*, 25:241–273, 1999.
- H. Prakken. Formalising Robert’s rules of order. Technical Report 12, GMD – German National Research Center for Information Technology, 1998.
- H. Prakken and T. Gordon. Rules of order for electronic group decision making – a formalization methodology. In J. Padget, editor, *Collaboration between Human and Artificial Societies*, LNCS 1624, pages 246–263. Springer, 1999.
- John R. Searle. *Speech Acts*. Cambridge University Press, 1969.
- M. J. Sergot. Modelling unreliable and untrustworthy agent behaviour. In B. Dunin-Keplicz, A. Jankowski, A. Skowron, and M. Szczuka, editors, *Monitoring, Security, and Rescue Techniques in Multiagent Systems*, Advances in Soft Computing, pages 161–178. Springer-Verlag, 2005.
- Marek Sergot and Robert Craven. Some logical properties of nonmonotonic causal theories. In *Proc. Eighth International Conference on Logic Programming and Non-Monotonic Reasoning (LPNMR’05), Calabria, September 2005*, LNAI 3662. Springer, 2005. In press.
- Y. Shoham and M. Tennenholtz. Emergent conventions in multi-agent systems. In *Proceedings of Conference on Principles of Knowledge Representation and Reasoning (KR’92)*, pages 225–231, 1992a.
- Y. Shoham and M. Tennenholtz. On the synthesis of useful social laws for artificial agent societies. In *Proceedings of Tenth National Conference on Artificial Intelligence (AAAI’92), San Diego*, pages 276–281. The AAAI Press/ The MIT Press, 1992b.
- Y. Shoham and M. Tennenholtz. On social laws for artificial agent societies: Off-line design. *Artificial Intelligence*, 73(1-2):231–252, 1995.
- Hudson Turner. A logic of universal causation. *Artificial Intelligence*, 113: 87–123, 1999.
- W. van der Hoek, M. Roberts, and M. Wooldridge. Social laws in alternating time: Effectiveness, feasibility, and synthesis. Technical report, Dept. of Computer Science, University of Liverpool, 2004. Submitted.