# Believing the Simplest Course of Events

James Delgrande

Simon Fraser University

Canada

jim@cs.sfu.ca

(Joint work with Hector Levesque, University of Toronto)

# Introduction

*Goal:*

> We are interested in an approach for reasoning in a dynamic domain with nondeterministic actions in which an agent's (categorical) beliefs correspond to the simplest course of events consistent with the agent's observations.

Here *simplest* corresponds to the most likely or plausible explanation.

Consider the following situation:

- There is a light switch.
- Toggling the switch turns a light *on* if it is *off*, and *off* if it is *on*.
- As in the real world, we are never absolutely certain that pressing the switch will have the expected result.

# Introduction

An agent knows that the light is *on* and toggles the switch twice.

- With no other information, the agent would believe the light is *on* and both actions succeeded.

- If it *senses* that the light is *on*, it would *not* believe that perhaps both actions failed (even though this also accounts for the light being *on*).

- If the agent *senses* that the light is *off*, it would believe that a toggling action failed.

# Introduction

Consider what this requires:

- An agent will have a set of *beliefs* concerning the real world.
  - These beliefs may be incomplete or incorrect.
- An agent may execute *actions*
  - The agent's beliefs will evolve as actions are executed
  - Actions may fail, or have unintended consequences
  - Thus we will need to keep track of actions that the agent *believes* it executed, and those *actually* executed.
  - So one way or another we will need an account of *nondeterminism*.
- An agent may *sense*, or be told, information about the world.
  - This information may conflict with the agent's beliefs, so an account of *revision* is needed.
  - It may also conflict with the actions the agent believed it executed, so beliefs about actions may also need to be revised.

# Introduction

**Overall Approach**: Augment an epistemic extension to the *situation calculus* with ranking functions, as used in *belief revision*, along with a formalization of nondeterminism.

**Very Roughly:**

- An agent's beliefs will be represented by *situations* (think: *possible worlds*), encoded in FOL (rather than a modal logic).
- Situations are assigned a *plausibility ranking*. Those with rank $= 0$ characterise categorical beliefs and those $> 1$ characterise counterfactual states of affairs.
- These plausibilities are modified following sensing and action execution.

**(Claimed) Result**: A general, qualitative model of an agent that is able to reason and maintain its stock of beliefs in via sensing in a nondeterministic domain.

# Overview

- Introduction
- Background:
  - the situation calculus
  - belief revision
- The Approach:
  - intuitions
  - (some) details
  - properties
- Conclusion

# (A bit more) Introduction

We would like to handle sequences such as the following:

- An agent believes that lights $l_1$ and $l_2$ are off.
  It believes that it turns on $l_1$, but in fact switches on $l_2$.
  It believes $l_1$ is on and $l_2$ off.
  Via sensing it learns that $l_2$ is on.
  It then believes that $l_1$ is off, and that originally it turned on $l_2$ and not $l_1$.

- An agent believes that a light is on.
  It toggles the switch twice
  It believes that the light is on.
  It senses that the light is off.
  It then believes that one toggle action failed.

# (A bit more) Introduction

To handle situations such as the preceding:

- We require a theory of action and belief.

    ☞ We adopt the Scherl-Levesque extension to Reiter's *basic action theories* expressed in the situation calculus.

- An agent must keep track of not just its beliefs, but other (non-believed) possibilities.

    ☞ We use ranking functions, as a representation of an agent's *epistemic state*, to keep track of counterfactual situations.

- We require a theory of actions with unexpected or unpredictable outcomes.

    ☞ To this end, we develop a theory of qualitative nondeterminism

- These notions need to be integrated to allow for sequences of (possibly mistaken) actions, sensing, and (not covered here) revisions.

# Background: The Situation Calculus (SC)

- The SC is a FOL theory for reasoning about action.
    - Idea: Actions take the world from one state to another.

- There are 2 distinguished sorts:
    - *actions*: e.g. $put(r, x, y)$ for robot $r$ putting object $x$ on $y$.
    - *situations*: these denote possible world histories.
        - $S_0$ denotes the initial state of the real world.
        - $do(a, s)$ denotes the situation that results from $s$ after executing action $a$.

- A predicate whose truth value is situation dependent is called a *fluent*.
    - E.g. $Holding(r, x, s)$

- In a basic action theory the truth of a fluent $\phi(do(a, s))$ is defined in terms of $a$ and fluents true at $s$ (next slide).

# The Situation Calculus

Examples:

- Definitions:

$$Init(s) \doteq \neg\exists a \exists s'.\ s = do(a, s')$$

- Foundational axioms:

$$Init(S_0)$$
$$do(a_1, s_1) = do(a_2, s_2) \supset a_1 = a_2 \wedge s_1 = s_2$$

- Blocks world:

$$On(A, B, S_0),\ On(B, Table, S_0)$$
$$Holding(x, do(a, s)) \equiv$$
$$((\neg Holding(x, s) \wedge a = PickUp)\ \vee$$
$$(\ Holding(x, s) \wedge a \neq PickUp))$$

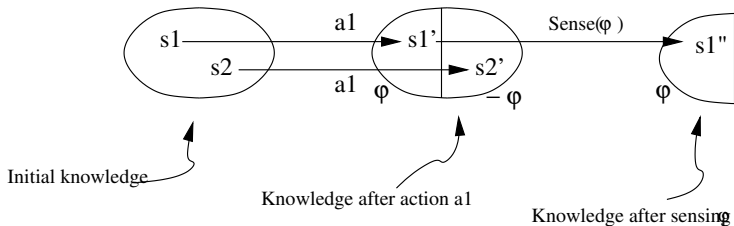# Knowledge and the Situation Calculus

Scherl and Levesque provide a possible worlds account of knowledge in the SC:

- A $B$ fluent gives the belief accessibility relation.

    - $B(s', s)$ holds when the agent in $s$ thinks that $s'$ might be the actual situation.

- $SF(a, s)$ holds when sensing action $a$ returns value 1 in $s$.

- Successor state axiom for $B$:

    $B(s'', do(a, s)) \equiv$
    $\quad \exists s'[s'' = do(a, s') \land B(s', s) \land (SF(a, s') \equiv SF(a, s))].$

- Belief is defined in terms of $B$:
    $Bel(\phi, s) \doteq \forall s'.B(s', s) \supset \phi[s'].$

# Knowledge and the Situation Calculus



- The first oval represents situations that are $B$ related to $S_0$;
  - I.e. the sitations characterising the agent's initial beliefs.
- The next oval represents situations that are $B$ related to $do(a, S_0)$;
- The last oval represents those related to $do(sense_\phi, do(a, S_0))$.

# Background: Belief Revision

- Next we extend this account to deal with situations with differing plausibilities where the agent's beliefs may be *revised*.

- First, we review key notions in *belief revision*

# Belief Revision

In revision, an agent

- incorporates a new belief $\phi$,
- while maintaining consistency (unless $\vdash \neg\phi$).

We'll use the standard semantic construction of faithful rankings.

- A faithful ranking is a *total preorder* over possible worlds
  - Lower-ranked worlds are more plausible
- We'll use non-negative integers to indicate plausibility values
  - This is slightly more general and easier to work with.
- Agent's beliefs given by the set of worlds with plausibility 0.
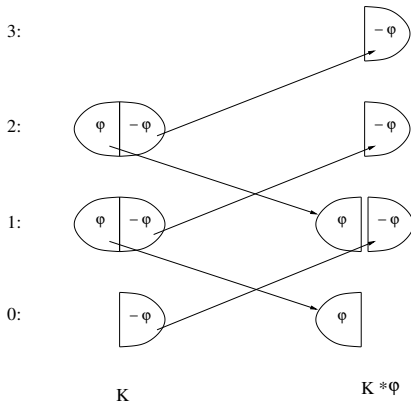
# Belief Revision: Characterization

- We adopt the approach suggested in [DarwichePearl97].
- In revising by $\phi$:
  - $\phi$ worlds retain their relative ranking, as do $\neg\phi$ worlds, but
  - the $\phi$ worlds have their ranking reduced so that a $\phi$-world has ranking 0.
  - The ranking of $\neg\phi$ worlds is increased by 1.

☞ However, any approach to iterated revision can be used in the framework.

# Revision in [DarwichePearl97]

Think of the total preorder as giving the agent's *epistemic state*.

For revising by $\phi$ we have:

# Expressing Plausibilities in the Situation Calculus

We can (tentatively) express plausibility using

$$B(s', n, s) \quad \text{where } n \geq 0$$

to indicate that in $s$ the agent considers $s'$ to have plausibility $n$.

- The agent's beliefs at $s$ are given by situations $s'$ where $B(s', 0, s)$.

☞ More later

# The Approach: Nondeterminism

Our stance:

- Nondeterminism is an *epistemic* notion reflecting an agent's limited knowledge and perception.
- The world is deterministic
  - ☞ Each state of the world is uniquely determined by its predecessor and the action executed.
- Examples
  - Flipping a coin
  - Inadvertently pressing the wrong light switch
  - An action failing for no known reason

# Nondeterminism

- We introduce predicate *Alt* where

  $$Alt(a_1, a_2, p, s)$$

  expresses that an agent intending to execute action $a_1$ may in fact execute $a_2$ with plausibility $p$ in situation $s$.

- Most often, for action $a$,

  $Alt(a, a, 0, s)$
  will hold.

# Nondeterminism

Examples:

- Toggling ($t$) a light switch:
  $Alt(t, x, p, s) \equiv (x = t \land p = 0) \lor (x = null \land p = 1)$

- Flipping ($f$) a coin:
  $Alt(f, x, p, s) \equiv (x = fH \land p = 0) \lor (x = fT \land p = 0)$

  ☞    $f$ is a *virtual action*; it is never executed in the real world.

- Throwing a dart
  $tB$ is the action of throwing a dart so it hits the dartboard;
  $tW$ is the action where the dart hits the adjacent wall.

  $$Alt(tB, x, p, s) \equiv$$
  $$\neg Dim(s) \supset ((x = tB \land p = 0) \lor (x = tW \land p = 1)) \quad \land$$
  $$Dim(s) \supset ((x = tB \land p = 0) \lor (x = tW \land p = 0))$$

# Nondeterminism and Belief

Example:

- There are two switches, left and right, both *off*.
- If the agent flips the left switch, it will believe the left switch is *on*.
- If the agent attempts to flip the right switch, but instead flips the left one, it will believe the left switch is *off*.

Conclusion:

> *When there can be nondeterministic actions, the physical actions that actually occur are insufficient to determine the situations the agent considers possible*

☞ This leads us to adopt a four-place fluent $B(s', n, \sigma, s)$ where $\sigma$ represents the sequence of actions that the agent *believed* it was performing at the time.

# Nondeterminism and Belief

$B(s', n, \sigma, s)$ expresses that:

if:

- the agent believes it executed action sequence $\sigma$,
- but actually executed the actions in $s$,

then

- situation $s'$ has plausibility $p$ according to the agent.

# Nondeterminism and Belief

*Alt* actions bear on an agent's beliefs in two ways

1. For $B(s', n, \sigma, s)$, the actions in $\sigma$ and $s$ are pairwise *Alt*-related.

2. Assume that $B(s', n, \sigma, s)$ and $Alt(a_1, a_2, p, s)$ hold.
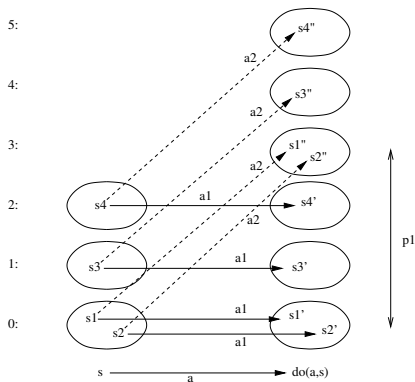
   if:
   - the agent believes it executed $a_1$ in $s$

   then:
   - situation $do(s', a_2)$ would have plausibility $n + p$ in the resultng epistemic state.

# Alternative Actions

Agent believes it executes $a_1$ in $s$; $Alt(a_1, a_2, p_1, s)$, $Alt(a_1, a, p_2, s)$ are true.



Note: There can be many Alt actions to $a_1$.

# Evolution of the B Fluent

- $B(s', n, \sigma, s)$ means that at $s$, where the agent believes it has executed actions in $\sigma$, $s'$ has plausibility n.

- Beliefs are characterised by the most plausible accessible situations:

$$Bel(\phi, \sigma, s) \doteq \forall s'. B(s', 0, \sigma, s) \supset \phi[s'].$$

- The agent's initial beliefs are characterised by $B$ instances of the form $B(s', n, \langle \rangle, S_0)$.

- We wish to characterise $B$ following the execution of action $a$, for physical actions and sensing actions.

- This leads to a somewhat daunting successor-state axiom for the B fluent (see the paper!).

- We next sketch the intuitions for the two types of actions.

# Change in Plausibility: Sensing Actions

Sensing actions are handled via revision as sketched earlier.

☞ Sensing actions are assumed to always succeed.

Consider $B(s', n, \sigma, s)$. Let $a$ be the action of sensing $\phi$.

- $a$-successors to $B$ look like

  $$B(do(a, s'),\ n',\ \sigma \cdot a,\ do(a, s))$$

  where
  - if the sensing result of $\phi$ at $s$ and $s'$ agree then
    $n' = n - MinPlaus(\phi, s)$
  - otherwise
    $n' = n + 1$.

# Change in Plausibility: Physical Actions

Consider $B(s', n, \sigma, s)$. Let $a$ be a physical action and assume that $Alt(a_i, a, p_1, s)$ and $Alt(a_i, a^*, p_2, s)$ are true. There are two cases.

1: An $a_i, a$-successor to $B$ looks like

$$B(do(a_i, s'), \ n, \ \sigma \cdot a_i, \ do(a, s))$$

- The agent intends to execute $a_i$; in fact it executes $a$.
- The plausibility of the $a_i$-successor of $s'$ is unchanged.
- Note that $a = a_i$ is Scherl-Levesque, extended to plausibilities.

2: An $a^*, a_i, a$-successor to $s$ looks like:

$$B(do(a^*, s'), \ n + p_2, \ \sigma \cdot a_i, \ do(a, s))$$

- The agent intends to execute $a_i$; in fact it executes $a$;
  $a^*$ is an alternative to $a_i$.
- Thus the plausibility of $do(a^*, s')$ is increased by $p_2$ at $do(a, s)$.

# Example: Toggling a Light Switch

- A light is initially *on*, and this is known by the agent.

- Toggling the switch changes the state of the light from *on* to *off* or vice versa.

- The agent toggles the light switch twice.
  It believes the light is on

- It observes that the light is *off*.
  It concludes that one of the toggling actions must have failed.

# Example: Formalization

There is just one initial situation, $S_0$.
A basic action theory is given as follows:

- $On(S_0)$

- $B(S_0, 0, \langle\rangle, S_0)$

- $On(do(a, s)) \equiv (a = t \wedge \neg On(s)) \vee (a \neq t \wedge On(s))$

- $SF(a, s) \equiv On(s) \vee a \neq sL$

- $Alt(t, x, p, s) \equiv (x = t \wedge p = 0) \vee (x = null \wedge p = 1)$
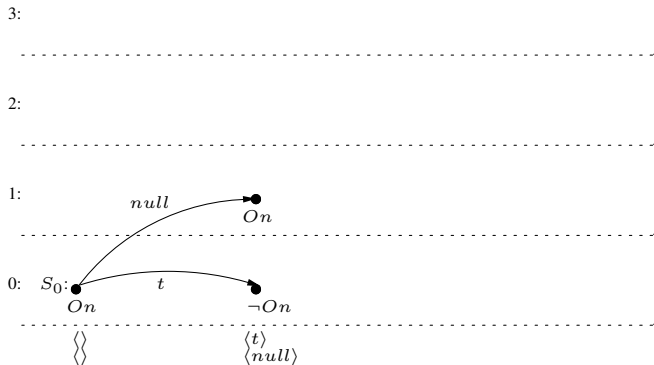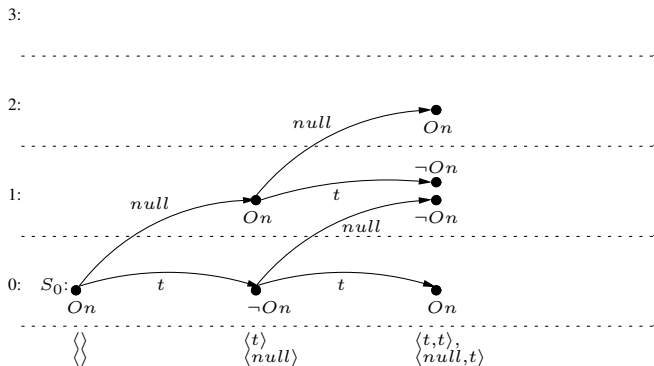
# Example

Initially:

3:

2:

1:

0:   $S_0$: ●
        $On$

        ⟨⟩
        ⟨⟩

Following a failed toggling action:



3:

2:

1:      $null$     $On$

0:   $S_0$:   $On$    $t$     $\neg On$

$\langle \rangle$       $\langle t \rangle$

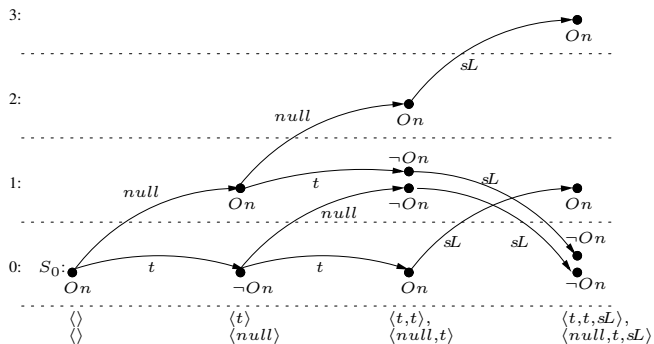$\langle \rangle$       $\langle null \rangle$

# Example

Next following a successful toggling action:

# Example

After sensing the light:

# The Approach: Properties

We obtain the following results.

- If an agent intends to execute $a_i$ but in fact executes $a$, it will believe the action effects of $a_i$.

- The agent believes the result of a sensing action.

- If an agent believes $\phi$ to hold, then it believes it will believe $\phi$ after sensing $\phi$.

  - Of course, if $\phi$ is false then it will believe $\neg\phi$ after sensing $\phi$.

- For revision defined in the obvious fashion, the AGM postulates hold.

# Conclusion

We have developed a general model of an agent that

- may execute (apparently) nondeterministic actions
- and may sense its environment.

The agent's beliefs evolve according to

- the sequence of actions it believes it executes
- and the results of sensing actions.

Notably, the agent believes those actions occurred which give the simplest explanation of its observations.

# Conclusion

The approach

- is developed within an epistemic extension of the situation calculus, incorporating plausibility orderings,
- in order to integrate reasoning about (nondeterministic) actions with sensing and (not covered here) belief revision.

As well:

- We retain the results of basic action theories, and so inherit the formal results of such theories.
- While we present a specific approach, the framework is readily generalisable.