

**«Уральский федеральный университет
имени первого Президента России Б.Н. Ельцина»**

**Институт социальных и политических наук
Департамент философии**

Кафедра онтологии и теории познания

Проблема фреймов и пути её решения в нормативных контекстах

Допустить к защите:

Зав. кафедрой
к.ф.н.,
профессор,
Кислов А. Г.

Магистерская диссертация
студента 2 курса
Касаткина А. В.

Научный руководитель:
д.ф.н.,
профессор,
Анкин Д. В.

Екатеринбург

2014

Оглавление

ВВЕДЕНИЕ.....	3
1. Определение проблемы фреймов.....	9
1.1. История вопроса.....	9
1.2. Постановка проблемы.....	14
2. Проблема ветвления.....	24
2.1. Минимизация изменений.....	24
2.2. Причинные связи.....	31
3. Проблема квалификации.....	38
3.1. Семантика возможных миров.....	38
3.2. Состояния по умолчанию: нормы предположений.....	61
ЗАКЛЮЧЕНИЕ.....	67
Список литературы.....	69

ВВЕДЕНИЕ

Актуальность темы исследования.

Тема данной работы находится в области одной из самых молодых и динамично развивающихся наук — Искусственного Интеллекта (ИИ). И хотя еще в 70-е гг., после знаменитого отчета Джеймса Лайтхилла¹, появились пессимистические настроения относительно будущего этой молодой и еще не окрепшей науки, она продолжает развиваться. Можно даже сказать, что только в последнее время она вышла на траекторию размеренного поступательного движения. До начала 90-х годов история этой науки представляла череду взлетов и падений, окрыляющих побед, и унижительных поражений. И только после этого началась настоящая работа, без сверх ожиданий, но и без сверх разочарований.

Однако, даже без сильных потрясений, динамичность развития ИИ, как и в целом Computer Science, никогда не вызывала сомнения. Постоянно ставятся все новые проблемы, находят новые решения и подходы. Однако некоторые проблемы прочно заняли пьедестал классических, и остаются там по сей день. Одна из таких проблем — так называемая «проблема фреймов». Эта проблема была поставлена еще в 1969 г.², на заре развития компьютерных технологий и непосредственно связанного с ними ИИ, однако удовлетворительного решения не нашла до сих пор. И тому видят несколько причин³. Это и особенности публикаций по данной тематике — они направлены на решение «игрушечных проблем» (обвинение в стиле сэра Д. Лайтхилла); и в целом особенность исследований в компьютерных науках — оно направлено не на решение общей проблемы, а на то, чтобы программа заработала, поэтому часто страдает излишней конкретикой и решениями ad hoc. Но главной причиной

¹ Lighthill, J: Artificial Intelligence: A General Survey // Artificial Intelligence: a paper symposium, Science Research Council

² Впервые сформулирована в McCarthy J. and Hayes P. J. Some philosophical problems from the standpoint of Artificial Intelligence Интернет источник. Режим доступа: <http://www-formal.stanford.edu/jmc/mcchay69.pdf> (дата обращения 15.05.2014 г.)

³ Morgenstern L. The Problem with Solutions to the Frame Problem Интернет источник. Режим доступа: <http://www-formal.stanford.edu/leora/fp.pdf> (дата обращения 15.05.2014 г.), pp.3-5

(приобретающей особую значимость на фоне ранее приведенных) является то, что это одна из ключевых проблем темпорального размышления. Она шире, чем кажется, и поиски её решения становятся вынужденным шагом при любой попытке размышлять о времени, изменениях и действиях. По крайней мере в рамках логицистского подхода к ИИ⁴.

Это действительно философская проблема (поставленная с точки зрения искусственного интеллекта). В явной форме здесь выражаются древние философские проблемы, связанные с понятием причинности и познания. Проблема фреймов ставит вопрос о том, как описать изменяющийся под влиянием действий мир так, чтобы не упустить никаких причинных связей, но и не считать все в мире ежесекундно меняющимся (старый спор Парменида с Гераклитом). В связи с попытками объяснить наш мир непонятливому компьютеру, эти философские проблемы выступают в четкой и ясной форме, и во многом становятся яснее нам самим.

Поэтому решение данной проблемы является весьма актуальной задачей, которую ставят перед собой многие исследователи, но зачастую не достигают желаемого результата. Хотя справедливости ради стоит отметить, что «неудачные попытки» решения проблемы фреймов — это не безнадежные провалы. Это по-своему удачные подходы, в которых, однако находят изъяны, зачастую — за счет переформулировки самой проблемы (подробнее о толкованиях проблемы фреймов в гл. 1). Так, например, Джон МакКарти (сформулировавший проблему) дал одно из решений проблемы, которое не удовлетворило научное сообщество, казалось несостоятельным и неосуществимым. Однако ряд специалистов, сотрудничавших с ним в этом проекте, затем использовали идеи МакКарти в создании программы RCS Advisor — программы, проверяющей способность шаттла к маневрированию в случае неполадок⁵. Поэтому вполне логично ожидать и возможность

⁴ Манифестом логицизма в ИИ считается работа McCarthy J. Programs with common sense. / M. L. Minsky (Ed.), Semantic information processing. Cambridge, MA: MIT Press, 1968, 403-409.).

⁵ Lifschitz V. The Frame Problem, Then and Now Интернет источник. Режим доступа: <http://www.cs.utexas.edu/users/vl/papers/jmc.pdf> (дата обращения 15.05.2014 г.)

практического применения от решения проблемы фреймов, даже неудачного.

Тот нормативный подход, который предлагается в данной работе, кажется актуальным данной проблеме, в первых, по причине того, что основывается на семантике возможных миров, которая очень хорошо подходит для рассуждений в условиях неполной информации, а также потому, что соответствует интуициям относительно того, как с подобными проблемами справляются люди.

Степень разработанности проблемы.

Литературы, посвященной проблеме фреймов в том или ином аспекте настолько много, что упомянуть в этой работе все представляется невозможным. Поэтому мы обратим внимание только на ключевые работы, представляющие глубоко проработанные и зарекомендовавшие себя подходы. Однако стоит отметить, что упомянутая литература в основном существует в виде статей в периодических изданиях и материалах конференций, в больших монографических исследованиях на эту тему наблюдается определенный недостаток.

Начинается история проблемы фреймов с уже упомянутой работы МакКарти [35]. Им же были предложены варианты решения проблемы на основе немонотонной логики, например в [33] и [32] Суть данного подхода заключается в возможности признания исключений и возможности последующего пересмотра выводов при признании чего-то «аномалией». Критика предложенного автором решения была представлена в знаменитой работе [21], где предлагается сценарий «Йельской стрельбы». Этот сценарий не поддается немонотонной логике, как и конкурентные действия.

Для преодоления этого затруднения были предложены так называемые хронологические подходы, которые вводят строгий временной порядок, и тем самым избегают затруднений, возникающих в Йельском сценарии. Это, например, [25], [28], [48].

Другим решением стали подходы, основанные на понятии причинности:

[26], [23],[13]. Но такие подходы, опять же, не позволяют конкурентных и неизвестных действий

Не использующий немонотонные рассуждения, но остающийся в рамках близкого к ситуационному исчислению формализма подход представлен в [16], где проводится категоризация функций на устойчивые и текущие, действия рассматриваются как элементарные, но порождающие цепи не прямых эффектов при определенных обстоятельствах. Похожий подход был независимо предложен в [47].

Однако объясняющих аксиом осталось все равно довольно много, около $2F$ (где F – количество функций). Модификация данного подхода в [43] привела к снижению количества аксиом до $F+A$, где A – это количество действий. В этой системе также возможно рассмотрение конкурентных действий, однако с весьма ограниченными характеристиками.

Процедурный подход, предложенный в [19], широко известный STRIPS, является по своей сути в первую очередь системой для планирования (поиска решения в графе состояний). В отношении проблемы фреймов там принимается посылка, что если действие не известно явным образом как меняющее значение какой-либо функции, то оно этого и не делает. Этот подход очень имеет очевидные преимущества: он интуитивно понятен и эффективен с вычислительной стороны, однако он не может моделировать условные действия, что отмечено в [39], [40]. Также он с трудом воспринимает конкурентные действия и в целом все, что не отражено явным образом — в соответствии со своим базовым принципом.

Если говорить о более общем взгляде на проблему фреймов как проблему темпорального рассуждения, то можно упомянуть здесь работы логиков, направленные на выявление адекватного описания временной структуры: [31], [41], [42], [44], [45]. Данные подходы относятся к тематике модальных логик, и потому в особенности стоят упоминания в данной работе, где подход также по своей сути является модальным. Однако здесь разговор идет только о времени.

Есть ряд работ, посвященных построению модальных логик действия: [1], [14], [22]. Однако здесь сосредоточенность на действии делает его главным объектом исследования, и потому главный вопрос проблемы фреймов — описание среды, которое мы должны вывести из описания состояния её в другое время и описания действия — не поднимается.

Модальных подходов к самой проблеме фреймов нами не было обнаружено, поэтому та узкая область, в которой развивается наша работа, остается не разработанной в литературе.

Объект исследования:

формальные системы, моделирующие изменяющуюся среду и действия в ней.

Предмет исследования:

Способы преодоления некоторых затруднений, встречающихся в такого рода системах, общее название которых — проблема фреймов.

Цель и задачи исследования.

Конечной целью данного исследования является демонстрация возможностей модально-логического подхода к решению проблемы фреймов, конкретно с точки зрения логики деонтических (нормативных) модальностей. Достижение поставленной цели осуществляется решением следующих задач:

1. Дать четкое рабочее определение проблемы фреймов для базового формализма, описывающего действие;
2. Рассмотреть способы решения в нем проблемы ветвления, как под-проблемы проблемы фреймов;
3. Рассмотреть применение семантики возможных миров как удобной модели для рассмотрения действий в условиях неполной информации;
4. Наметить способ построения на базе такой семантики нормативно-модальных правил вывода.

Теоретико-методологическая основа исследования.

Основная методология данного исследования — это логический анализ

темпоральных рассуждений.

Структура и объем работы.

Магистерская диссертация состоит из введения, трех глав, содержащих 6 подразделов, и заключения. В первой главе проводится историческое рассмотрение проблемы фреймов, дается строгое её определение для введенного базового формализма (как двух взаимосвязанных проблем — проблемы ветвлений и проблемы квалификации). Во второй главе рассматривается проблема ветвлений, дается вариант её решения за счет расширения базового формализма (добавляются понятия ограничений состояний, информации о влияниях, и, самое главное, законов причинных связей). В третьей главе рассматривается проблема квалификации, базовый формализм расширяется посредством введения для него семантики возможных миров, рассматриваются свойства этой модели и предлагается вариант решения проблемы квалификации за счет оценки литералов по степени ожидаемости, и вводящихся на этой основе модальных нормативных операторов. Список литературы включает 51 наименование. Общий объем работы – 72 страницы.

1. Определение проблемы фреймов

1.1. История вопроса

Проблема фреймов была впервые озвучена в знаменитой работе Джона МакКарти и Патрика Хайеса «Некоторые философские проблемы с точки зрения Искусственного Интеллекта»⁶. Данная статья нацелена на рассмотрение ряда философских вопросов, которые, во первых, должны быть решены для построения «разумных машин», во-вторых, могут быть прояснены с помощью моделирования для таких машин. Под «разумной машиной» в данном случае понимается программа (выполняемая на компьютере), которая из спецификации некоторой реальной ситуации, целевой ситуации, и стратегии действий, может сделать вывод относительно того, сможет ли данная стратегия достичь цели. Для адекватной работы такой программы необходимы формальные спецификации таких понятий, как *может*, *знает* и *влечет*. Философским вопросам о знании, причинности и возможности, выраженным в формальном языке, пригодном для перевода в программный, и посвящена данная статья. В ней предлагается формализм, позволяющий рассуждать об упомянутых понятиях, а в конце указывается ряд пока не решенных проблем. Среди них — проблема фреймов⁷.

Предложенный формализм — это *ситуационное исчисление* (подмножество логики предикатов первого порядка). Опишем основные моменты этого исчисления.

Ситуации — это некоторые мгновенные снимки реальности, описывающие мир в какой-либо момент времени. Например — 20:17:39 UTC 20 июля 1969 г. лунный модуль «Орел» коснулся поверхности Луны, на его борту были Нилом Армстронг и Эдвин Олдрин, Ричард Никсон был президентом США, Генеральным секретарем ЦК КПСС был Леонид Ильич Брежнев и т. д.

⁶ McCarthy J. and Hayes P. J. Some philosophical problems from the standpoint of Artificial Intelligence Интернет источник. Режим доступа: <http://www-formal.stanford.edu/jmc/mcchay69.pdf> (дата обращения 15.05.2014 г.)

⁷ McCarthy J. and Hayes P. J. Some philosophical problems from the standpoint of Artificial Intelligence Интернет источник. Режим доступа: <http://www-formal.stanford.edu/jmc/mcchay69.pdf> (дата обращения 15.05.2014 г.), pp. 30-31

Состояние — это коллекция ситуаций, например состояние, когда Леонид Ильич Брежнев являлся Генеральным секретарем ЦК КПСС - это коллекция ситуаций с 1966 по 1982 год. Состояние $On(BlockA, BlockB)$ – это коллекция ситуаций, когда блок А находится на блоке В. Состояние — это пример *флюента*. Интуитивно это понятие можно определить как нечто, чье значение изменяется со временем. Состояния являются Булевыми флюентами, то есть их возможные значения - «истина», «ложь». Есть и другие типы флюентов, например — $President(USA)$. Его значение в 1969 году (после 20 января) — Ричард Никсон, а в 2013 — Барак Обама. Специализированный предикат $Holds$ соотносит флюенты и ситуации, например $Holds(S1, On(BlockA, BlockB))$, где $S1$ — имя ситуации. Действия определяются как функции на множестве состояний. Так, действие $PutOn(BlockA, BlockB)$ отображает состояния, где на блоках А и В ничего не стоит, в состояния, где блок А стоит на блоке В. Функция $Result$ отображает ситуацию до производства действия в ситуацию после производства действия. Таким образом, если мы имеем $Holds(S0, Clear(BlockA))$ и $Holds(S0, Clear(BlockB))$, то мы можем говорить о $Result(PutOn(BlockA, BlockB), S0)$. И для того, чтобы сказать, что действие $PutOn(BlockA, BlockB)$ приводит к тому, что блок А теперь находится на блоке В, мы можем сказать: $Holds(Result(PutOn(BlockA, BlockB), S0), On(BlockA, BlockB))$.

Вопрос, который ведет к появлению проблемы фреймов, заключается в следующем: если имеет место ситуация $Holds(S0, Clear(BlockA))$; $Holds(S0, Clear(BlockB))$; $Holds(S0, Red(BlockB))$, то что мы можем сказать относительно значения $Holds(Result(PutOn(BlockA, BlockB), S0), Red(BlockB))$?

Когда мы описываем действия, мы должны принимать посылку о том, что описания мира изменяются с течением времени. Описание действия как раз и заключается в описании такого изменения: действие задается как функция перехода от одной ситуации к другой, изменившейся. Но некоторые данные о мире нерелевантны по отношению к действию: действие не зависит от них, и они не зависят от действия. Однако мы каким-то образом должны указать, что с

ними происходит при переходе от одной ситуации к другой.

Это формирует проблему предсказания результата действия. В авторской формулировке МакКарти и Хайеса не это имелось в виду под проблемой фреймов, но в расширенном толковании сейчас и это тоже относят к общей теме проблемы фреймов.

Сама же проблема фреймов появляется в связи с тем вариантом решения проблемы предсказания, который предлагают авторы статьи. Они предлагают прописать аксиомы, указывающие на неизменность функций при совершении действия. Такие аксиомы называются *аксиомами фреймов*, так как задают те рамки⁸, в которых действия влияют на среду. Пример аксиомы фреймов:

$$\text{Holds}(s, \text{Red}(b1)) \rightarrow \text{Holds}(\text{Result}(\text{PutOn}(b1, b2), s), \text{Red}(b1))$$

Или, в более общем виде:

$$\text{color}(s, b1) = \text{color}(\text{Result}(\text{PutOn}(b1, b2), s), b1)$$

Но если при описании каждого действия указывать в качестве правил результаты функции перехода для каждого флюента, то таких правил потребуется очень много: при n флюентов и m действий их количество будет равно nm . Именно эта необходимость столь значительного количества аксиом для того, чтобы получить вывод о том, что большинство вещей в мире остаются теми же самыми после того, как мы предприняли какое-либо действие, и называется *проблемой фреймов*.

Проблема в данном случае связана с тем, что составление списка аксиом, а затем вычисление значений флюентов в согласии с ним, отнимают слишком много времени и ресурсов. Это невыгодно. Кроме того, это контринтуитивно. МакКарти является основателем логицистского подхода в области Искусственного Интеллекта (ИИ), и цель его заключается в том, чтобы смоделировать в терминах формальной логики мышление человека. Причем он отмечал, что программы неплохо научились справляться со сложным для человека умственным трудом (математические операции с большими числами,

⁸ Frame (англ.) - рамка

доказательство теорем и т. д.), однако простейшие процессы размышления, доступные каждому человеку, еще только предстоит симулировать программами⁹. Это и было его целью — моделирование обыденного сознания. И модель в данном случае кажется сомнительной: для определения того, чтобы определить, какой будет ситуация после завершения действия, нам обычно не приходится долго вспоминать, как это действие воздействует на те или иные свойства объектов.

Еще одна проблема возникает в связи с тем, что аксиомы фреймов часто оказываются ложными¹⁰, в частности в случае, если в системе допустимы конкурентные действия: если кто-то распылит краску на блок В в то время, пока мы ставим на него блок А, цвет не сохранится. То же может произойти и в более простых системах: наше действие может повлечь за собой разного рода не прямые эффекты, и будучи выполнено одинаково в разных условиях, приведет к различным последствиям. Например то, загорится лампа или нет после нажатия на выключатель, зависит от состояния электропроводки.

Итак, в наиболее буквальном смысле, проблема фреймов заключается в том, что для вывода о том, что большинство особенностей мира остаются неизменными при выполнении какого-либо действия, требует огромного количества аксиом. Но уже в указанной статье авторы предлагают альтернативную формулировку — как предсказывать состояния, следующие за выполнением действия внутри ситуационного исчисления, но без использования аксиом фреймов. То есть как мы можем выразить в строгом виде, что кроме того, что явно известно как изменяющееся, все остальное остается неизменным?

Эта классическая проблема фреймов в последующем переинтерпретировалась множеством способов, обобщающих изначальную формулировку. Первой ступенью¹¹ обобщения можно считать рассмотрение

⁹ MMcCarthy J. Programs with common sense. / M. L. Minsky (Ed.), Semantic information processing. Cambridge, MA: MIT Press, 1968, 403-409., p. 403

¹⁰ См.: McDermott D. A temporal logic for reasoning about processes and plans. Cognitive Science, 6, 101-155.

¹¹ Рассмотрение степеней обобщения проблемы в разных формулировках см.: Morgenstern L. The Problem with

проблемы фреймов как *проблемы устойчивости*¹²: общей проблемы предсказания свойств, которые остаются неизменными при совершении действий (являются *устойчивыми* по отношению к действию). Здесь нет привязки ни к аксиомам фреймов, ни к ситуационному исчислению. Это общая проблема рассуждения об изменяющемся с течением времени посредством действий мире.

Далее к проблеме устойчивости добавилась проблема определения того, как вещи изменяются со временем¹³. Такая формулировка проблемы фреймов погружает нас в область общих проблем *темпорального проектирования*. Эти две части общей проблемы рассуждения о изменениях во времени заметно отличаются: если проблема устойчивости — это проблема такого представления изменяющегося мира, которое не потребует множества аксиом фреймов; то в случае предсказания изменений нет никаких возражений против введения аксиом, специфицирующих влияние действий на среду. Проблема тут становится вычислительной: как задать все без исключения изменения, вызванные действием, и сделать это максимально эффективно. Вычисление результатов действий может занять очень много времени в мирах, богатых причинными взаимосвязями. Так, если принести в комнату сумку, то в комнате будут находиться и все вещи, которые были в сумке, а не только она сама. Появление таких *непрямых (продленных) эффектов* действия порождает свой подраздел проблемы фреймов — *проблему ветвлений*¹⁴. Подробнее о ней мы поговорим ниже.

Еще более общей формулировкой проблемы является включение в её рассмотрение и направленного в прошлое темпорального проектирования¹⁵ -

Solutions to the Frame Problem Интернет источник. Режим доступа: <http://www-formal.stanford.edu/leora/fp.pdf> (дата обращения 15.05.2014 г.), pp. 10-13.

¹² См.: Shoham Y. (1988). Reasoning about change: Time and causation from the standpoint of artificial intelligence. Cambridge, MA: MIT Press.

¹³ См.: Morgenstern L. and Stein L. A. Why things go wrong: A formal theory of causal reasoning. Proceedings of AAAI 1988, 518-523.

¹⁴ См.: Finger J. J. Exploiting constraints in design synthesis. Department of Computer Science STAN-CS-88-1204, Stanford University.

¹⁵ Также предлагается в Morgenstern L. and Stein L. A. Why things go wrong: A formal theory of causal reasoning. Proceedings of AAAI 1988, 518-523.

как по результатам действия мы можем вычислять причины? Тем самым проблема фреймов становится общей проблемой рассуждения о времени. В неё включается проблематика объяснения — как объяснить то, что наши предсказания о результатах действия оказались ошибочными? И также здесь появляется *проблема квалификации*¹⁶ — грубо говоря, проблема спецификации условий, при которых действие достигнет предполагаемого результата. Её мы также рассмотрим отдельно.

Самой общей интерпретацией проблемы фреймов является философская интерпретация её как следствия проблемы индукции¹⁷: нас заставляет ожидать тот или иной эффект действия (или то, что останется неизменным поле него) исключительно повторение таких случаев в опыте. А так как индукция — это недостоверное знание, поэтому и появляются проблемы при попытках описать правила, полученные индуктивным путем, в виде посылок дедукции.

Таким образом, проблема фреймов является довольно сложным явлением: будучи поставлена как конкретная проблема в рамках конкретного подхода, она вобрала в себя как подпроблемы множество сложностей, возникающих в целом при рассуждении об изменяющемся во времени мире и действиях в нем. Трендом в поиске ее решения стало предварительное установление формулировки, которая будет рассматриваться именно в данном исследовании. И наше исследование последует тем же путем.

1.2. Постановка проблемы.

Для строгой формулировки проблемы мы введем предварительный вариант формализма, описывающего действия, четко сформулируем проблемы, возникающие в нем, а затем попытаемся добавить к базовой системе необходимые расширения. Общий вид формализма во многом следует за

¹⁶ См.: McCarthy J. Applications of circumscription to formalizing common-sense knowledge. *Artificial Intelligence*, 28, 86-116.

¹⁷ См.: Fetzer J.H. The frame problem: *Artificial Intelligence meets David Hume*. In K.M. Ford and P.J. Hayes (Eds.), *Reasoning agents in a dynamic world: The frame problem*, 55-69. Greenwich, CT: JAI Press.

предложенной Майклом Тьельшером¹⁸ теорией действия. Поэтому иногда мы будем опускать некоторые доказательства, ссылаясь на проделанную им работу. Однако мы будем использовать более удобную для нас нотацию, и в ходе решения поставленных проблем постепенно будем все далее удаляться от авторского варианта.

Первое, что нам необходимо — это описание *состояний*. Состояние мы будем понимать как мгновенный снимок фрагмента реальности, который мы моделируем, в определенный момент времени¹⁹. Существенным является то, что данное описание должно быть составлено из атомарных пропозиций. Это необходимо нам для того, чтобы затем, когда в рассмотрение будут введены действия, мы могли их описывать как воздействующие только на малую часть ситуации. Если же описание состояния не имеет внутренней структуры, (как, например, в теории автоматов), нам необходимо будет описывать действие как изменение одной ситуации (в целом) на другую, при этом описывая результат выполнения действия для каждой ситуации отдельно.

Атомарные пропозиции описывают свойства, или, более обще — отношения между *индивидами* (или *объектами*). Истинностные значения таких пропозиций могут меняться с течением времени, тем самым меняя состояние. Мы назовем их *функциями*²⁰.

Индивиды мы обозначим множеством $O = \{o_1, o_2, \dots\}$, в качестве переменных будем использовать строчные буквы, написанные курсивом (по возможности из конца алфавита). Функции — $F = \{F_1(x_1, x_2, \dots, x_n), \dots\}$, где x_1, x_2, \dots, x_n являются объектами. Аридность функций также может быть нулевой. *Функциональный литерал* (или просто *литерал*) — это функция, или ее отрицание, записываемое как $\sim F_1(x_1, x_2, \dots, x_n)$. Множество литералов

¹⁸ Thielsher M. Challenges for action theories. Springer, 2000

¹⁹ Состояние здесь оказывается синонимом ситуации, в том виде, как она понимается в ситуационном исчислении МакКарти. Пока мы действительно позволим себе использовать эти понятия как синонимы, однако позднее различие между ними будет проведено.

²⁰ Автор называет их “fluents” или “fluent literals”, но так как достойный перевод для fluents сложно предложить, а транслитерация «флюенты» не слишком хорошо звучит, мы будем использовать слово «литералы» (благодаря второму варианту оригинала), понимая под ними именно такие «текучие», изменчивые литералы.

противоречиво, если содержит функцию вместе с ее отрицанием, в других случаях оно непротиворечиво. Состояние — это максимальное непротиворечивое множество литералов (то есть содержащее каждую из функций в том или ином значении, но только в одном). Множество состояний будет записываться как $S = \{S_1, S_2, \dots\}$; каждое состояние будет иметь вид $S_i = \{\sim F_1(x_1, x_2, \dots, x_n), F_2(x_1, x_2, \dots, x_n), \dots\}$.

Вторым важным компонентом любой теории действий является само действие. Действия, будучи выполнены, вызывают изменение состояний. Но мы будем описывать их не как функцию на множестве состояний, а как функцию на множестве литералов. Основной нашей посылкой здесь является то, что действие влияет только на маленький фрагмент реальности, точнее — на один литерал. Действие меняет его истинностное значение на противоположное. Мы будем записывать это так:

$\alpha(x)$ меняет $F(x)$ на $\sim F(x)$

Это означает, что действие α , будучи применено к объекту x меняет его состояние с $F(x)$ на $\sim F(x)$. Действие может также применяться не к одному объекту, а к нескольким. Тогда вместо (x) будет стоять последовательность (x_1, x_2, \dots, x_n) . Однако последовательность эта в любом случае совпадает в действии и в манипулируемом литерале. Также действия могут применяться к литералам с арностью 0, тогда оно вообще не будет зависеть от объектов. Такая запись будет называться *законом действия*, причем одно действие может описываться несколькими законами, например:

$\alpha(x)$ меняет $F(x)$ на $\sim F(x)$

$\alpha(x)$ меняет $\sim F(x)$ на $F(x)$

Множество имен действий мы обозначим $A = \{\alpha(x), \dots\}$, множество законов действий — $AL = \{\alpha(x) \text{ меняет } C(x) \text{ на } E(x), \dots\}$ (где $C(x)$ и $E(x)$ — литералы, разные значения одной функции). Закон действия применим к некоторому состоянию тогда, когда в состоянии выполняется $C(x)$. Результатом применения закона действия в состоянии S_i будет новое состояние S'_i ,

представляющее собой $(S_i \setminus C(x)) \cup E(x)$. S_i мы назовем *предварительным состоянием-наследником* S_i .

Итак, наш базовый домен (универсум рассуждения) D состоит из четверки (O, F, A, AL) , где O - множество индивидов, F - множество функций, A - множество действий и AL - множество законов действий. Ситуации мы не включаем в домен, так как они формируются из базовых элементов посредством определения. *Моделью изменений* Σ домена D будет отображение пары состояние-действие на множество состояний (являющихся предварительными наследниками).

Каждый домен предлагает общее, ситуационно-независимое знание об эффектах действий. Эта информация может использоваться для заключения о результатах действий в той или иной реальной ситуации.

Сценарий действия задается конкретным состоянием, информацию о котором мы получаем благодаря наблюдению. Мы предполагаем, что у нашего агента имеются сенсоры, позволяющие получить информацию о состоянии того или иного объекта в терминах литералов. Она выражается в виде *восприятия*, которое задается множеством литералов. Множество восприятий мы обозначим $B = \{B_1, \dots\}$, где каждое восприятие имеет вид $B_i = \{F_i, \dots\}$. Пара (B, D) задает сценарий.

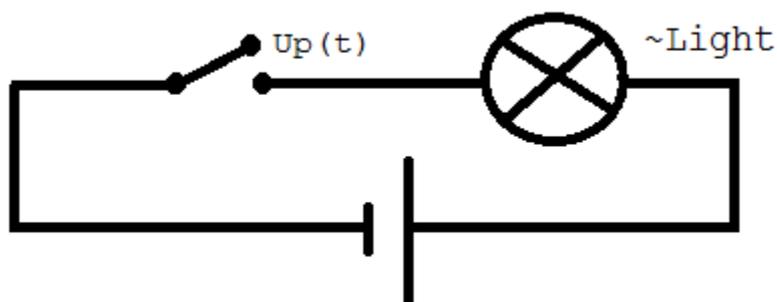
Функция $Res([\])$ от последовательности действий (возможно пустой), определенная на множестве состояний, будет обозначать результирующее состояние, которое мы получаем, произведя последовательно указанные действия.

Сценарий может быть изображен в виде дерева, где корнем будет изначальное состояние, данное нам восприятием, а ветви будут образовываться последовательным применением всех имеющихся в описании домена действий сначала к изначальному состоянию, а затем к наследникам. Дуги можно пометить именем действия, выполнение которого приводит к данной ситуации. $Res([\])$ будет иметь в качестве своего значения ту вершину дерева, которая

достигается при переходе от начального состояния по дугам, пометки на которых соответствуют тем действиям, которые результат которых мы пытаемся найти.

Это будет нашей базовой теорией действия. Сейчас мы проиллюстрируем её на примере, который послужит формулировке проблемы фреймов для указанной теории действий. Нужно отметить, что изначальной формулировки проблемы фреймов мы уже избежали: описав состояния как составные структуры, мы смогли указать в законах действий то, что они влияют только на малую область состояния. А описывая правила формирования состояний-наследников (значения функции $Res([\])$), указали, что все остальное в ситуации остается неизменным. Однако такое решение порождает новые проблемы, которые мы и рассмотрим далее.

Рассмотрим электрическую цепь²¹:



Домен мы можем описать следующим образом:

$$O = \{t\};$$

$$F = \{Up(x), Light\};$$

$$A = \{\alpha(x)\};$$

$$AL = \{\alpha(x) \text{ меняет } \sim Up(x) \text{ на } Up(x), \alpha(x) \text{ меняет } Up(x) \text{ на } \sim Up(x)\}.$$

То есть у нас есть один индивид t , представляющий собой триггер (переключатель). Есть две функции: одна, $Up(x)$, означает положение триггера (вверху он или внизу), вторая, $Light$ (с арностью 0), означает состояние лампы:

²¹ Примеры с электрическими цепями активно применяются в работе М. Тьельшера, и мы, увидев в этом удачный прием, следуем в этом за автором, предлагая, в том числе, и свои цепи. Данная цепь приведена в Thielsher M. Challenges for action theories. Springer, 2000, p.18

горит она или нет. И есть действие, представляющее собой переключение триггера: оно меняет верхнее положение на нижнее, и наоборот.

Наблюдение предоставляет нам следующее восприятие:

$$\mathbf{B0} = \{\sim Up(t), \sim Light\}$$

Это задает нам сценарий. Итак, изначальная ситуация такова:

$$\mathbf{S0} = \{\sim Up(t), \sim Light\}$$

Каков будет результат действия $Res([\alpha(t)])$?

Согласно имеющимся законам действия, мы должны сделать вывод, что в итоге получится следующее:

$$\mathbf{S1} = \{Up(t), \sim Light\}$$

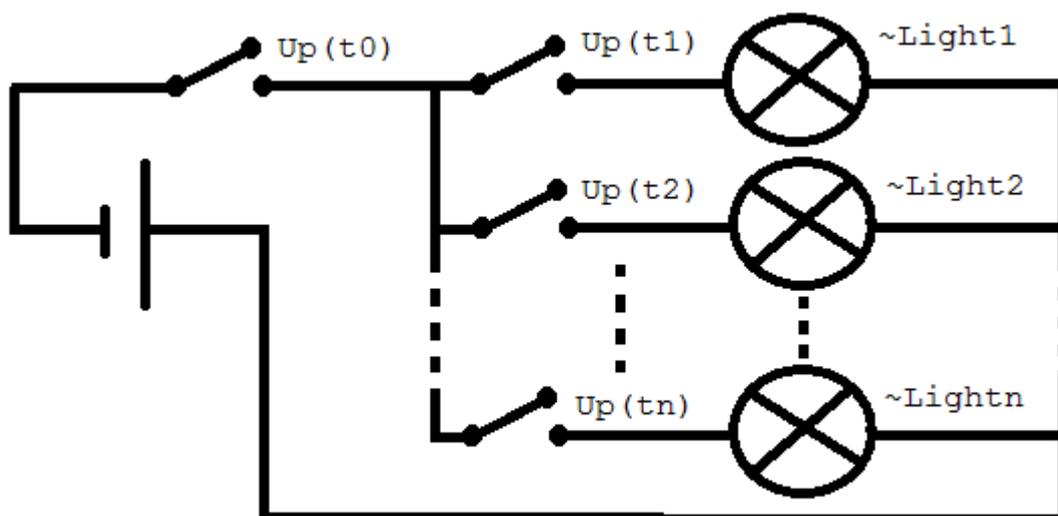
Однако это противоречит интуиции, что при опускании переключателя цепь замкнется, и лампа загорится. Как можно ввести в нашу систему это соображение?

Самый очевидный вариант — расширить законы действия. Мы можем переписать их так:

$AL = \{\alpha(x) \text{ меняет } \{\sim Up(x), \sim Light\} \text{ на } \{Up(x), Light\}, \alpha(x) \text{ меняет } \{Up(x), Light\} \text{ на } \{\sim Up(x), \sim Light\}\}$.

Тогда действие $\alpha(t)$ в указанной ситуации $\mathbf{S0} = \{\sim Up(t), \sim Light\}$ приведет к ожидаемому результату $\mathbf{S1} = \{Up(t), Light\}$. Однако всегда ли возможно такое решение?

Рассмотрим следующий пример:



$$O = \{t_0, t_1, t_2, \dots, t_n\};$$

$$F = \{Up(x), Light_1, Light_2, \dots, Light_n\};$$

Мы предполагаем здесь, что нам также нужно ввести действие по переключению триггера в иное положение. Но тогда как описать его законы? Для всех случаев, кроме t_0 , все не так сложно. В общем виде законы для них будут выглядеть так:

$$\alpha(t_i) \text{ меняет } \{\sim Up(t_0), Up(t_i), \sim Light_i\} \text{ на } \{\sim Up(t_0), \sim Up(t_i), Light_i\};$$

$$\alpha(t_i) \text{ меняет } \{\sim Up(t_0), \sim Up(t_i), Light_i\} \text{ на } \{\sim Up(t_0), Up(t_i), \sim Light_i\};$$

$$\alpha(t_i) \text{ меняет } \{Up(t_0), Up(t_i), \sim Light_i\} \text{ на } \{Up(t_0), \sim Up(t_i), \sim Light_i\};$$

$$\alpha(t_i) \text{ меняет } \{Up(t_0), \sim Up(t_i), \sim Light_i\} \text{ на } \{Up(t_0), Up(t_i), \sim Light_i\}.$$

Заметим, что при n подцепей с переключателем и лампой, таких законов потребуется $4n$. А как быть с законами для t_0 ? Если явно прописывать состояния, то законов только для включения его потребуется 2^n (по одному для каждой комбинации включения и выключения n триггеров), и еще столько же для выключения. Такое количество законов действий ($2^n + 4n$) заставляет нас вспомнить об оригинальной формулировке проблемы фреймов. Здесь, конечно, речь идет не об аксиомах фреймов, а о законах действия, но по сути мы получаем то же самое: чтобы описать, что меняется в мире посредством действия, а что нет, мы должны прописать огромное количество правил.

Такая формулировка проблемы и называется *проблемой ветвления*: не прямые эффекты действия, которые могут появиться при некоторых условиях, заставляют нас на графе, изображающем сценарий, добавлять дополнительные разветвления для одного действия, а попытка преодоления этих трудностей вызывает добавление к нашей модели огромного количества правил.

Это будет первая проблема, которую мы будем рассматривать. Вторая же проблема — проблема квалификации. Ее рассмотрение будет главной нашей задачей, ей уделяется наибольшее внимание в нашей работе, и именно в решении ее мы вводим семантику возможных миров и нормативные модальные

операторы, вынесенные в название работы. В основе этой проблемы лежит то, что действия могут не привести к желаемым результатам вследствие того, что какие-то условия этого действия не соблюдены. Данная проблема была озвучена после появления первых решений проблемы фреймов, связанных с *минимизацией состояний* (стратегия, предпринятая нами в нашей базовой теории действий — ограничение изменений, производимых действием). Вследствие такой минимизации, выводы могут стать неоднозначными: два и более состояния удовлетворяют условиям минимизации. Одним из первых примеров таких сценариев стала «Йельская стрельба»²². Там речь идет о действии заряжания ружья и выстрела через некоторое время. Будет ли жив Фред после выстрела?

В терминах введенной нами теории действия, можно выразить это так. Домен состоит из:

$$O = \{\text{turkey}^{23}, \text{gun}\};$$

$$F = \{\text{Alive}(x), \text{Loaded}(x)\};$$

$$A = \{\text{LOAD}(x), \text{FIRE}(x, y)\};$$

$AL = \{\text{LOAD}(x) \text{ меняет } \sim\text{Loaded}(x) \text{ на } \text{Loaded}(x)\}, \text{FIRE}(x, y) \text{ меняет } \{\text{Alive}(x), \text{Loaded}(y)\} \text{ на } \{\sim\text{Alive}(x), \sim\text{Loaded}(y)\}, \text{FIRE}(x, y) \text{ меняет } \{\text{Alive}(x), \sim\text{Loaded}(y)\} \text{ на } \{\text{Alive}(x), \sim\text{Loaded}(y)\}\}.$

Наблюдение, задающее нам сценарий: $B0 = \{\text{Alive}(\text{turkey}), \sim\text{Loaded}(\text{gun})\}.$

Итак, $S0 = \{\text{Alive}(\text{turkey}), \sim\text{Loaded}(\text{gun})\}.$

$Res([\text{LOAD}(\text{gun})]) = \{\text{Alive}(\text{turkey}), \text{Loaded}(\text{gun})\}.$

Пока все ясно А теперь мы предполагаем, что прошло какое-то время. В нашей модели выразить это довольно трудно. Мы сделаем это с помощью пропуска в последовательности действий. Итак, наш вопрос:

$$Res([\text{LOAD}(\text{gun}), \dots, \text{FIRE}(\text{turkey}, \text{gun})]) = ?$$

²² См.: Hanks S. and McDermott D. Nonmonotonic logic and temporal projection. *Artificial Intelligence*, 33(3):379-412.

²³ В оригинальном примере речь шла о Фреде, со временем Фред эволюционировал в индейку (например, в Thielsher M. *Challenges for action theories*. Springer, 2000). Мы последуем этой гуманистической традиции, хотябы для того, чтобы не путать записываемое с большой буквы имя, которое играет у нас роль индивида, с именем функции.

Не зная состояния, из которого мы начинаем действие, мы не можем точно предсказать результат. Так как у нас есть два закона для действия $FIRE(turkey, gun)$, точно предсказать его результат мы не можем. Более того, возможна ситуация, где ни один из законов не применим. Например, если второе, неизвестное действие, было таким же: $FIRE(turkey, gun)$, то

$$Res([LOAD(gun), FIRE(turkey, gun)]) = \{\sim Alive(turkey), \sim Loaded(gun)\}$$

А в такой ситуации мы не можем применить наше действие, так как нами не рассмотрены законы действия $FIRE(x, y)$ для случая, когда оно производится в ситуации, содержащей $\sim Alive(x)$.

Конечно, вариант проблемы Йельской стрельбы для нашего базового формализма выглядит не вполне естественно. Однако порожденная им *проблема квалификации* может рассматриваться в любом формализме. Она обращает внимание на то, что условия действий не всегда могут быть соблюдены. Однако проверка всех условий — дело затруднительное, и зачастую не всегда нужное: конечно, проводка может испортиться. Но это редкая ситуация, и мы не проверяем каждый раз, перед тем как включить свет, все ли в порядке. То же самое можно сказать относительно автомобиля, который мы пытаемся завести: причин, по которым двигатель может не работать, может быть множество. Однако если мы будем проверять, все ли в порядке, перед тем как завести машину, мы рискуем никуда не уехать в ближайшую неделю. Однако эти данные о возможных проблемах очень важны в случае, если мы не смогли завести двигатель или включить свет. Нам нужно знать, что может пойти не так, чтобы при возникновении проблем знать, на что обратить внимание, как объяснить себе неудачу и выйти из затруднительного положения. Проблема адекватного описания условий действия — это и есть проблема квалификации. Для её введения нужно отказаться от полноты информации. Например, в приведенном примере с Йельской стрельбой, мы можем сделать это просто посредством того, что сделаем недоступным восприятию состояние ружья. И если не зная, заряжено оно или нет, мы будем вынуждены стрелять, мы должны

иметь возможность предположить результат действия, и в зависимости от исхода принять верное решение (например если ружье не выстрелило — зарядить его). В общем и целом это проблема эффективного действия в условиях неполной информации. Точнее — проблема нахождения формализма, позволяющего реализовать такое действие.

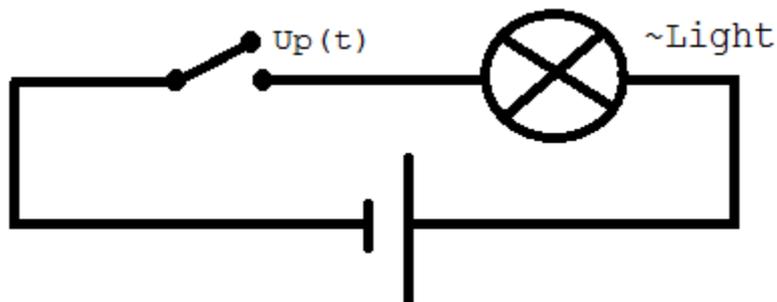
Эти две проблемы — проблема ветвления и проблема квалификации и будут рассматриваться нами как две стороны проблемы фреймов. Первая из них заключается в эффективной формализации информации обо всех не прямых эффектах, вызываемых действием; вторая — нахождение способов описания среды и действий в ней, в условиях неполной информации.

2. Проблема ветвления.

2.1. Минимизация изменений.

В данной главе мы подробно рассмотрим первую из двух, поднятых в предыдущей главе проблем. Остановились мы на том, что выявили необходимость в каком-то способе более краткого описания цепей причинных связей, порождаемых действием, чем полные описания всех возможных эффектов для всех возможных ситуаций.

Первым шагом на этом пути станет введение *ограничений состояний*²⁴. Вспомним нашу цепь с переключателем и лампой:



Тот факт, что замыкание цепи вызывает зажигание лампы, является некоторым общим характерным свойством данного домена. Оно следует из его структуры. Состояние, когда цепь замкнута, а лампа не горит является просто неприемлемым для данного домена (если, конечно, сама лампа, проводка и источник тока не повреждены). Это соображение мы и предлагаем выразить в виде нового элемента домена - ограничений, которые должны выполняться в любом состоянии, чтобы оно было приемлемым. Они будут выражаться в виде логической формулы. Все состояния должны удовлетворять этой формуле, то есть подстановка конкретных значений литералов на место соответствующих функций в формуле ограничения должны сохранять формулу истинной. Если состояние не удовлетворяет ограничению, то оно считается неприемлемым. В

²⁴ Ограничения состояний и их применение для предложенной системы рассмотрено в Thielscher M. Challenges for action theories. Springer, 2000, pp. 19-23. Вообще же идея минимизации изменений на основе каких-либо законов (называемых по-разному) предложена неформально в Winslett M. Reasoning about action using a possible models approach. In *Proceedings of the AAAI National Conference on Artificial Intelligence*, pp. 89-93, Saint Paul, MN, August 1988; формализованно — в Lifschitz V. Frames in the space of situations. *Artificial Intelligence*, 46:365-376, 1990.

данном домене можно ввести такое ограничение:

$$\text{Light} \equiv \sim \text{Up}(t)$$

В ситуации, заданной $\mathbf{B0} = \{\sim \text{Up}(t), \sim \text{Light}\}$ предварительное состояние-наследник, согласно законам действия, как уже отмечалось, будет таким:

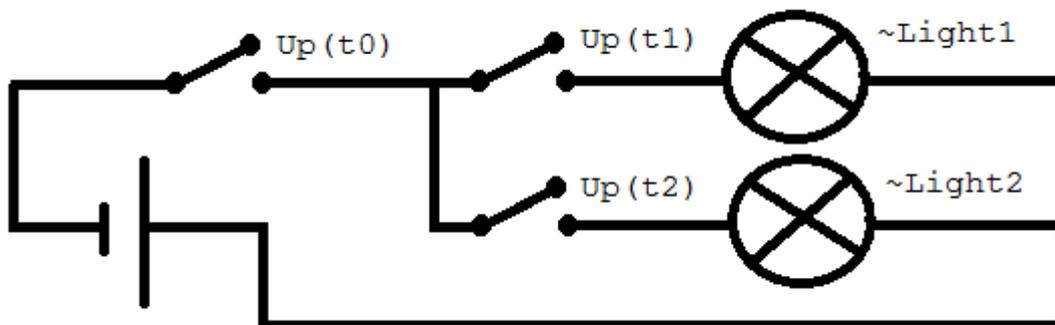
$$\text{Res}([\alpha(t)]) = \{\sim \text{Up}(t), \sim \text{Light}\}$$

При подстановке получившихся значений в наше ограничение, мы получаем ложь. Значит это состояние неприемлемо. Что же делать далее? Отбросить его и заменить состоянием $\{\text{Up}(t), \sim \text{Light}\}$ или $\{\sim \text{Up}(t), \text{Light}\}$?

Далее введем понятие *минимизации изменений*. Пусть $\mathbf{S0}$ — состояние, совместимое с ограничениями, α — некоторое действие. Оно приводит к некоторому предварительному состоянию наследнику \mathbf{T} . Если это состояние-наследник неприемлемо, то мы рассмотрим приемлемые состояния $\mathbf{T1}, \mathbf{T2}, \dots, \mathbf{Tn}$, которые содержат прямой эффект действия α , но являются приемлемыми. То из них, которое отличается от $\mathbf{S0}$ на меньшее количество литералов (это будет \mathbf{Ti} , если верно $(\mathbf{Ti} \setminus \mathbf{S0} \subset \mathbf{Tj} \setminus \mathbf{S0})$ для всех \mathbf{Tj}), будет состоянием в минимальными изменениями. Такое состояние мы и поставим на место наследника $\mathbf{S0}$ (уже не предварительного, т.к. мы проверили его на совместимость с ограничениями).

Состоянием с минимальными изменениями будет $\{\sim \text{Up}(t), \text{Light}\}$ (так как только в нем выполняется прямой эффект α). Его же мы ожидали получить и интуитивно.

Но поможет ли нам тот же подход решить нашу проблему в более богатом домене, где требуется слишком много аксиом? Например, рассмотрим такой (похожий на приведенный в первой главе, но ограниченный для наглядности):



$$O = \{t_0, t_1, t_2\};$$

$$F = \{Up(x), Light1, Light2\};$$

$$A = \{\alpha(x)\};$$

$$AL = \{\alpha(x) \text{ меняет } \sim Up(x) \text{ на } Up(x), \alpha(x) \text{ меняет } Up(x) \text{ на } \sim Up(x)\}.$$

Множество ограничений мы обозначим $C = \{\Psi, \dots\}$, где Ψ — формулы логики предикатов. В данном случае мы можем ввести следующие ограничения:

$$C = \{Light1 \equiv \sim Up(t_1) \ \& \ \sim Up(t_0), Light2 \equiv \sim Up(t_2) \ \& \ \sim Up(t_0)\}$$

В более сложных доменах, с большим количеством подцепей с триггером и лампой, мы можем задать ограничения формулой с квантором:

$$\forall i (Light_i \equiv Up(t_i) \wedge Up(t_0))$$

$$\text{Предположим, } \mathbf{B0} = \{Up(t_0), Up(t_1), \sim Up(t_2), \sim Light1, \sim Light2\}.$$

Вычислим $Res([\alpha(t_0)])$. Согласно законам действия, получаем:

$$\{\sim Up(t_0), Up(t_1), \sim Up(t_2), \sim Light1, \sim Light2\}$$

Данное состояние неприемлемо, т. к. нарушает условие $Light2 \equiv \sim Up(t_2) \ \& \ \sim Up(t_0)$ — его правая часть истинна, а левая — ложна. Построим удовлетворительные состояния, содержащие $\sim Up(t_0)$ — прямой эффект действия $\alpha(t_0)$. Это будут:

$$(1) \{\sim Up(t_0), Up(t_1), Up(t_2), \sim Light1, \sim Light2\}$$

$$(2) \{\sim Up(t_0), \sim Up(t_1), \sim Up(t_2), Light1, \sim Light2\}$$

$$(3) \{\sim Up(t_0), \sim Up(t_1), \sim Up(t_2), Light1, Light2\}$$

$$(4) \{\sim Up(t_0), Up(t_1), \sim Up(t_2), \sim Light1, Light2\}$$

$$S_0 \setminus (1) = \{\sim Up(t_0), Up(t_2)\}$$

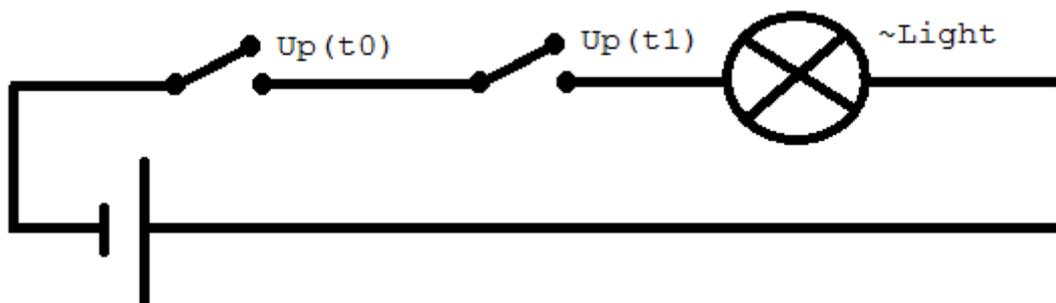
$$S_0 \setminus (2) = \{\sim Up(t_0), \sim Up(t_1), Light_1\}$$

$$S_0 \setminus (3) = \{\sim Up(t_0), \sim Up(t_1), Light_1, Light_2\}$$

$$S_0 \setminus (4) = \{\sim Up(t_0), Light_2\}$$

Мы получили два ближайших к S_0 случая. А ожидаемый был ясен — это случай (4). Изменение положения триггера, как в случае (1) не является интуитивно приемлемым следствием действия $\alpha(t_0)$.

На самом деле ту же проблему мы получили бы и в цепи с одной лампой, но двумя переключателями:



$$O = \{t_0, t_1\};$$

$$F = \{Up(x), Light\};$$

$$A = \{\alpha(x)\};$$

$$AL = \{\alpha(x) \text{ меняет } \sim Up(x) \text{ на } Up(x), \alpha(x) \text{ меняет } Up(x) \text{ на } \sim Up(x)\}.$$

$$C = \{Light \equiv \sim Up(t_1) \ \& \ \sim Up(t_0)\}$$

$$B_0 = \{Up(t_0), \sim Up(t_1), \sim Light\}.$$

$$Res'([\alpha(t_0)]) = \{\sim Up(t_0), \sim Up(t_1), \sim Light\}$$

$$(1) \{\sim Up(t_0), Up(t_1), \sim Light\}$$

$$(2) \{\sim Up(t_0), \sim Up(t_1), Light\}$$

$$S_0 \setminus (1) = \{\sim Up(t_0), Up(t_1)\}$$

$$S_0 \setminus (2) = \{\sim Up(t_0), Light\}$$

Мы получили всего два состояния, но они одинаково близки к

изначальному.

Почему мы ожидаем изменения состояния лампы, а не переключателя? На переключатели мы можем оказать прямое воздействие, а на лампу — нет. Но при этом её состояние зависит от переключателей. Для преодоления сложившихся затруднений нам нужно каким-то образом выделить те функции, которые вызывают изменения, и те, которые его претерпевают.

Возможным решением может быть *категоризация функций*²⁵. Мы можем выделить два типа функций: те, которые изменяются под влиянием действий и изменяют значения других как своих не прямых эффектов; и те, которые более изменчивы. Назовем их, соответственно, *первичными* и *вторичными*. И тогда мы сможем указать, что более близким к предыдущему приемлемому состоянию будут такие состояния — наследники, где меньше отличий в первичных функциях, и только на втором этапе оценивания мы будем обращать внимание на вторичные функции.

Первичные функции мы обозначим *Fp*, а вторичные - *Fs*. Для большей наглядности первичные функции также будем подчеркивать в выражениях, где это разделение играет роль. Тогда можно переписать условия нашего последнего примера следующим образом:

$$Fp = \{\underline{Up}(x)\};$$

$$Fs = \{\text{Light}\};$$

И тогда, получив в результате

$$S0 \setminus (1) = \{\sim\underline{Up}(t0), \underline{Up}(t1)\}$$

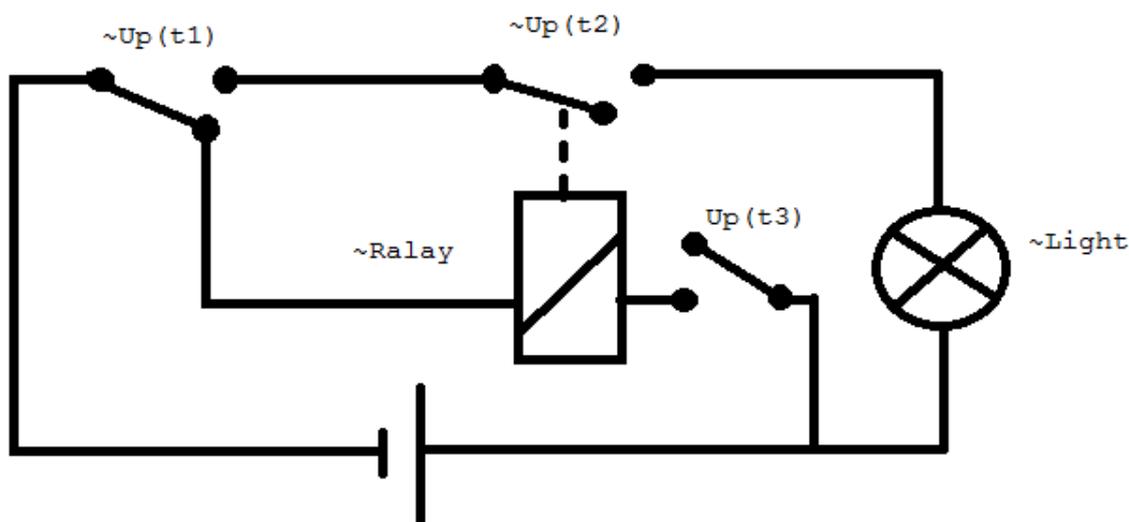
$$S0 \setminus (2) = \{\sim\underline{Up}(t0), \text{Light}\}$$

Мы выберем в качестве наследника $\{\sim\underline{Up}(t0), \sim\underline{Up}(t1), \text{Light}\}$, т. к. это состояние отличается от *S0* только одним первичным литералом (и одним вторичным), а $\{\sim\underline{Up}(t0), \underline{Up}(t1), \sim\text{Light}\}$ — двумя.

²⁵ Идеи выделения типов функций представлены уже в Lifschitz V. Frames in the space of situations. *Artificial Intelligence*, 46:365-376, 1990., также они предлагаются (с различными названиями для типов) в Val del A. and Shoham Y. Deriving properties of belief update from theories of action (II). In R. Bajcsy, editor, *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 732-737, Chambery, France, August 1993. Morgan Kaufmann. и Sandewall E. Reasoning about actions and change with ramication. In *Computer Science Today*, volume 1000 of LNCS. Springer, 1995.

Однако всегда ли можно определить категорию функции? Является ли разделение функций на прямо и косвенно манипулируемые подходящим для всех случаев?

Мы, к сожалению, вынуждены дать отрицательный ответ на это вопрос. Покажем это на следующем примере:



Здесь в цепь мы включили реле, которое может воздействовать на триггер $t2$. Попробуем описать этот домен:

$$O = \{t1, t2, t3\};$$

$$A = \{\alpha(x)\};$$

$$AL = \{\alpha(x) \text{ меняет } \sim Up(x) \text{ на } Up(x), \alpha(x) \text{ меняет } Up(x) \text{ на } \sim Up(x)\}.$$

$$C = \{Light \equiv Up(t1) \& Up(t2), Relay \equiv \sim Up(t1) \& \sim Up(t3), Relay \rightarrow \sim Up(t2)\}$$

Что происходит с функциями? Разделим их по принципу возможности прямого манипулирования:

$$Fp = \{Up(x)\};$$

$$Fs = \{Light, Relay\};$$

Переключать мы можем любые триггеры, на лампу и реле же воздействовать напрямую мы не можем, только посредством реле. Предположим, наше восприятие будет следующим:

$$B0 = \{\sim Up(t1), Up(t2), Up(t3), \sim Light, \sim Relay\}$$

Каков будет результат выполнения действия $\alpha(t1)$? Прямой эффект приведет нас к предварительному состоянию наследнику

$$Res'([\alpha(t1)]) = \{Up(t1), Up(t2), Up(t3), \sim Light, \sim Ralay\}.$$

Оно не выполняет ограничения, (конкретно - $Light \equiv Up(t1) \& Up(t2)$), поэтому строим приемлемые состояния, содержащие $Up(t1)$:

$$(1) \{Up(t1), Up(t2), Up(t3), Light, \sim Ralay\}$$

$$(2) \{Up(t1), \sim Up(t2), Up(t3), \sim Light, \sim Ralay\}$$

$$(3) \{Up(t1), Up(t2), \sim Up(t3), Light, \sim Ralay\}$$

$$(4) \{Up(t1), \sim Up(t2), \sim Up(t3), \sim Light, \sim Ralay\}$$

$$S0 \setminus (1) = \{\underline{Up(t1)}, Light\}$$

$$S0 \setminus (2) = \{\underline{Up(t1)}, \sim \underline{Up(t2)}\}$$

$$S0 \setminus (3) = \{\underline{Up(t1)}, \sim \underline{Up(t3)}, Light\}$$

$$S0 \setminus (4) = \{\underline{Up(t1)}, \sim \underline{Up(t2)}, \sim \underline{Up(t3)}\}$$

Четвертое состояние отличается на три первичных литерала, третье — на два (и один вторичный), второе также на два, первое же — всего на один. Поэтому последнее и будет нашим состоянием-наследником, что полностью соответствует интуиции. Отметим, что роль в выборе между (1) и (2) сыграло то, что $Up(t2)$ мы признали первичным литералом.

А что будет, если произвести действие $\alpha(t3)$?

$$Res'([\alpha(t3)]) = \{\sim Up(t1), Up(t2), \sim Up(t3), \sim Light, \sim Ralay\}$$

Этот предварительный наследник нарушает ограничение $Ralay \equiv \sim Up(t1) \& \sim Up(t3)$, поэтому рассмотрим приемлемые состояния с $\sim Up(t3)$:

$$(1) \{\sim Up(t1), \sim Up(t2), \sim Up(t3), \sim Light, Ralay\}$$

$$(2) \{Up(t1), Up(t2), \sim Up(t3), Light, \sim Ralay\}$$

$$(3) \{Up(t1), \sim Up(t2), \sim Up(t3), \sim Light, \sim Ralay\}$$

Сравним их с начальным состоянием:

$$S0 \setminus (1) = \{\sim \underline{Up(t2)}, \sim \underline{Up(t3)}, Ralay\}$$

$$S0 \setminus (2) = \{\underline{Up(t1)}, \sim \underline{Up(t3)}, Light\}$$

$$S0 \setminus (3) = \{\underline{Up(t1)}, \sim \underline{Up(t2)}, \sim \underline{Up(t3)}\}$$

Выбрать здесь состояние-наследник невозможно, так как и в (1) и в (2) по два первичных и по одному вторичному литералу. (3) мы можем отбросить, а выбрать между первыми двумя не в силах. Хотя интуитивно ясно, что выбрать нужно первый вариант: триггер t_2 под действием реле должен опуститься. Отметим, что именно то, что мы посчитали $Up(t_2)$ первичным литералом (то, что помогло нам сделать верный вывод относительно $Res([\alpha(t_1)])$), не дало нам сделать верный выбор здесь.

Проблема заключается в том, что триггер t_2 может подвергаться как прямому воздействию, так и непрямому — под действием реле (которое само лишь вторично). Это показывает нам, что причинные взаимосвязи не так просты. Но что если попытаться ввести в нашу систему сами причинные связи?

2.2. Причинные связи.

Вместо того, чтобы пытаться минимизировать изменения в состояниях по тому или иному принципу, мы зададим вопрос напрямую: какова форма причинных связей²⁶? Ведь именно они порождают не прямые эффекты действий: триггер t_2 в приведенном ранее примере изменял свое положение не потому, что мы переключили t_3 , а потому, что появился ток в обмотке реле. А ток там появился не только потому, что мы опустили переключатель t_3 , но и потому, что в районе t_1 цепь тоже была замкнута. Иначе наше действие ни к чему бы не привело (кроме прямого эффекта). Наше действие было лишь *поводом* для того, чтобы сложившиеся *условия* смогли повлиять на реле, а затем и на переключатель. Это рассуждение дает нам форму закона причинной связи:

А вызывает смену $\sim B$ на B при условии C ;

«А» здесь выступает поводом, переключающим эффектом. Это либо прямой эффект действия, либо результат другого, аналогичного по форме,

²⁶ Впервые понятие причинности было использовано в Elkan C. Reasoning about action in rst-order logic. In *Proceedings of the Conference of the Canadian Society for Computational Studies of Intelligence (CSCSI)*, pp/ 221-227, Vancouver, Canada, May 1992. Morgan Kaufmann. для классического ситуационного исчисления. Мы следуем здесь за автором Thielsher M. *Challenges for action theories*. Springer, 2000. Кроме того, что это причинные связи для рассматриваемой нами теории действия, они еще, в отличии от первого упомянутого варианта, позволяют циклическую зависимость.

закона причинной связи. «В» — это та функцию, на которую оказывается воздействие. Мы, как и в законах действий, полагаем, что воздействие производится точно, меняя только один литерал. «С» же — это условие выполнения закона причинной связи, та почва, посеяв в которую переключающий эффект, мы получаем в результате изменение.

Выпишем законы причинных связей для нашей цепи с реле:

- (1) Relay вызывает смену $Up(t2)$ на $\sim Up(t2)$ при условии \top ;
- (2) $\sim Up(t1)$ вызывает смену $\sim Relay$ на Relay при условии $\sim Up(t3)$;
- (3) $\sim Up(t3)$ вызывает смену $\sim Relay$ на Relay при условии $\sim Up(t1)$;
- (4) $Up(t1)$ вызывает смену Relay на $\sim Relay$ при условии \top ;
- (5) $Up(t3)$ вызывает смену Relay на $\sim Relay$ при условии \top ;
- (6) $Up(t1)$ вызывает смену $\sim Light$ на Light при условии $Up(t2)$;
- (7) $Up(t2)$ вызывает смену $\sim Light$ на Light при условии $Up(t1)$;
- (8) $\sim Up(t1)$ вызывает смену Light на $\sim Light$ при условии \top ;
- (9) $\sim Up(t2)$ вызывает смену Light на $\sim Light$ при условии \top ;

Знак тавтологии \top мы использовали для описания ситуации, когда изменение происходит при любых условиях.

Применяются законы причинных связей следующим образом. Мы будем записывать предварительные состояния-наследники в виде пары (S', E) , где S' будет означать само состояние (полное непротиворечивое множество литералов), а E — множество эффектов. В это множество будут входить как прямой эффект действия, так и непрямые, полученные вследствие применения законов причинных связей. Сначала мы будем перебирать все применимые к прямому эффекту законы (те, где он выступает переключающим эффектом), затем к первому непрямому, затем к второму непрямому и т. д. Отметим, что порядок применения законов причинных связей не имеет значения²⁷.

Итак, попробуем вывести результат действия $\alpha(t3)$. После применения закона действия $\alpha(x)$ мы получим

²⁷ Доказательство приведено в Thielsher M. Challenges for action theories. Springer, 2000, pp. 33-34

$(\{\sim\text{Up}(t1), \text{Up}(t2), \sim\text{Up}(t3), \sim\text{Light}, \sim\text{Ralay}\}, \{\sim\text{Up}(t3)\})$;

Далее, применением (3) получим

$(\{\sim\text{Up}(t1), \text{Up}(t2), \sim\text{Up}(t3), \sim\text{Light}, \text{Ralay}\}, \{\sim\text{Up}(t3), \text{Ralay}\})$;

Других законов причинных связей, где переключаящим эффектом выступает $\sim\text{Up}(t3)$, нет. Но есть закон (1), где таковым является полученный на предыдущем шаге не прямой эффект Ralay . Отсюда

$(\{\sim\text{Up}(t1), \sim\text{Up}(t2), \sim\text{Up}(t3), \sim\text{Light}, \text{Ralay}\}, \{\sim\text{Up}(t3), \text{Ralay}, \sim\text{Up}(t2)\})$;

Ralay более не встречается среди наших законов причинных связей на месте «повода», встречается $\sim\text{Up}(t2)$, но для применения закона (9) у нас нет необходимости, т. к. вызываемый им эффект $\sim\text{Light}$ уже присутствует в нашем состоянии-наследнике (а необходимый для применения закона (9) литерал Light , соответственно, нет). Поэтому мы останавливаемся на достигнутом, и получаем:

$\text{Res}([\alpha(t3)]) = \{\sim\text{Up}(t1), \sim\text{Up}(t2), \sim\text{Up}(t3), \sim\text{Light}, \text{Ralay}\}$

Это полностью соответствует нашим ожиданиям. Не возникнет ли проблем с выводением $\text{Res}([\alpha(t1)])$?

Применением к $\mathbf{B0} = \{\sim\text{Up}(t1), \text{Up}(t2), \text{Up}(t3), \sim\text{Light}, \sim\text{Ralay}\}$ закона действия $\alpha(x)$ получаем

$(\{\text{Up}(t1), \text{Up}(t2), \text{Up}(t3), \sim\text{Light}, \sim\text{Ralay}\}, \{\text{Up}(t1)\})$;

А затем, согласно (6), выводим

$(\{\text{Up}(t1), \text{Up}(t2), \text{Up}(t3), \text{Light}, \sim\text{Ralay}\}, \{\text{Up}(t1), \text{Light}\})$;

Закон (4) неприменим, т. к. у нас уже есть $\sim\text{Ralay}$, Light не производит никаких не прямых эффектов, поэтому мы делаем вывод

$\text{Res}([\alpha(t1)]) = \{\text{Up}(t1), \text{Up}(t2), \text{Up}(t3), \text{Light}, \sim\text{Ralay}\}$

И это именно тот вывод, который мы ожидали получить.

Заметим, что ряд выписанных нами законов действий проявляют некоторую симметрию. Имеются в виду пары (2)-(3) и (6)-(7). Может быть нам стоило записать их иначе, например так:

$\{\sim\text{Up}(t1), \sim\text{Up}(t3)\}$ вызывает смену $\sim\text{Ralay}$ на Ralay ;

$\{Up(t1), Up(t2)\}$ вызывает смену $\sim Light$ на $Light$;

Это позволило бы нам уменьшить количество законов причинных связей, что ускорило бы вывод: ведь получив какой-либо эффект нам необходимо перебрать все эти законы для проверки, нет ли среди них применимого в данной ситуации. Заметим, что в приведенном примере это бы ничего не изменило. Но в некоторых случаях различие между условием и поводом очень значительно. В таких ситуациях, где нам важно различать, какое состояние уже имело место, и потому является лишь условием; а какое появилось только недавно, и потому может вызвать новые изменения. Поэтому важным нововведением является сохранение информации о полученных эффектах во множестве E , и добавление его к предварительным наследникам. Это позволяет нам получать разумные выводы в случаях *циклической зависимости*. Простым примером здесь может служить пара переключателей, соединенных тугой пружиной. Когда мы меняем положение одного, второй также меняет положение, и наоборот. Причинные связи здесь будут такими:

(1) $\sim Up(t1)$ вызывает смену $Up(t2)$ на $\sim Up(t2)$ при условии \top ;

(2) $\sim Up(t2)$ вызывает смену $Up(t1)$ на $\sim Up(t1)$ при условии \top ;

(3) $Up(t1)$ вызывает смену $\sim Up(t2)$ на $Up(t2)$ при условии \top ;

(4) $Up(t2)$ вызывает смену $\sim Up(t1)$ на $Up(t1)$ при условии \top ;

Рассмотрим два варианта восприятия:

$\mathbf{B0} = \{\sim Up(t1), \sim Up(t2)\}$;

$\mathbf{B'0} = \{Up(t1), Up(t2)\}$;

В ситуации, выводимой из $\mathbf{B0}$ переключим $t1$. Предварительным наследником (без указания эффекта) будет

(a) $\{Up(t1), \sim Up(t2)\}$.

В ситуации же $\mathbf{B'0}$ переключим $t2$. Предварительным наследником, как мы можем заметить, снова будет (a). Минимизация изменений никак не может тут помочь определить, каким же будет результат каждого действия. Но может помочь информация о полученных эффектах. В случае $\mathbf{B0}$ и действия $\alpha(t1)$ мы

получим

$$(\{Up(t1), \sim Up(t2)\}, \{Up(t1)\}),$$

Откуда по (3):

$$(\{Up(t1), Up(t2)\}, \{Up(t1), Up(t2)\}),$$

Итоговое состояние будет

$$S1 = \{Up(t1), Up(t2)\}.$$

В случае же $B'0$ и действия $\alpha(t2)$ будет

$$(\{Up(t1), \sim Up(t2)\}, \{\sim Up(t2)\});$$

$(\{\sim Up(t1), \sim Up(t2)\}, \{\sim Up(t2)\})$ по (2);

$$S'1 = \{\sim Up(t1), \sim Up(t2)\}.$$

Такие возможности, конечно, являются плюсом подхода, основанного на причинных связях. Но мы вынуждены признать, что множество причинных связей является довольно громоздким. Кроме того, перечислить их все довольно сложно, нужно ничего не упустить.

На помощь нам здесь приходит понятие *информации о влияниях*. Это общая информация, которая лежит в основе причинных связей, показывающая, какая функция в качестве переключающего эффекта может воздействовать на какую-либо другую, выступающую в роли непрямого эффекта. Задать информацию о влияниях можно в виде упорядоченных пар функций. Например в нашем последнем примере это будет следующее множество пар:

$$\{(Up(t1), Up(t2)), (Up(t2), Up(t1))\}.$$

Мы обозначим это множество $I = \{(Fi, Fj), \dots\}$.

С помощью информации о влияниях можно вывести законы причинных связей из ограничений состояний. Для этого нужно привести формулу, описывающую ограничение к конъюнктивной нормальной форме (КНФ), затем в каждом конъюнкте для каждой пары подлежащих литералам функций проверить, не принадлежит ли она множеству I . И если да, то ко множеству законов причинных связей (обозначим его CL) нужно добавить следующую запись:

$\neg l_j$ вызывает смену $\neg l_k$ на l_k при условии

Где \prod^l означает конъюнкцию, а l_1, \dots, l_m – литералы, дизъюнкция которых образует конъюнкт, а $l_j \text{ b } l_k$ – те литералы, чьи подлежащие формулы входят в информацию о влияниях в виде (l_j, l_k)

Для пояснения разберем вывод законов причинной связи для нашего последнего рассмотренного домена.

$$I = \{(Up(t1), Up(t2)), (Up(t2), Up(t1))\};$$

$$C = \{Up(t1) \equiv Up(t2)\}.$$

Выведем КНФ $Up(t1) \equiv Up(t2)$:

$$\begin{aligned} Up(t1) \equiv Up(t2) &= (Up(t1) \rightarrow Up(t2)) \& (Up(t2) \rightarrow Up(t1)) = \\ &= (\sim Up(t1) \vee Up(t2)) \& (\sim Up(t2) \vee Up(t1)) \end{aligned}$$

Рассмотрим первый конъюнкт. Из входящих в него подлежащих литералам $(\sim Up(t1) \text{ и } Up(t2))$ формул $(Up(t1) \text{ и } Up(t2))$ можно получить две упорядоченные пары: $(Up(t1), Up(t2))$ и $(Up(t2), Up(t1))$. Обе они принадлежат информации о влияниях. Поэтому в соответствии с формулой, получаем два закона действия:

$Up(t1)$ вызывает смену $\sim Up(t2)$ на $Up(t2)$ при условии \top ;

$\sim Up(t2)$ вызывает смену $Up(t1)$ на $\sim Up(t1)$ при условии \top ;

(\top мы используем потому, что других дизъюнктов не осталось)

Из второго дизъюнкта, по тем же соображениям, получаем

$\sim Up(t1)$ вызывает смену $Up(t2)$ на $\sim Up(t2)$ при условии \top ;

$Up(t2)$ вызывает смену $\sim Up(t1)$ на $Up(t1)$ при условии \top ;

Таким образом, мы получили все четыре введенных ранее интуитивно закона причинных связей. В несколько измененном виде этот алгоритм также позволяет выводить законы причинных связей и для случаев, когда в ограничениях состояний встречаются кванторы²⁸.

В заключение, опишем получившуюся модель. Домен D теперь задается шестеркой (O, F, A, AL, C, I) , где O - множество индивидов, F - множество

²⁸ См.: Thielsher M. Challenges for action theories. Springer, 2000, алгоритм без кванторов — p. 37; алгоритм с кванторами — p. 42

функций, A - множество действий, AL - множество законов действий, C - множество ограничений состояний и I - множество упорядоченных пар, задающих информацию о влияниях. Из C и I могут быть выведены законы причинных связей, обозначаемые CL , и имеющие форму

A вызывает смену $\sim B$ на B при условии C ;

где A означает переключающий эффект, B — это та функцию, на которую оказывается воздействие, C — условие выполнения закона причинной связи.

Предварительные состояния-наследники записываются в виде пары (S', E) , где S' будет означать само состояние, а E — множество эффектов, приведших к этому состоянию.

Также, по соображениям, которые прояснятся в следующей главе, к описаниям состояния мы далее будем добавлять информацию о той последовательности действий, которая привела к нему. Мы будем записывать это так: $S_i = \{F_1, F_2, \dots, F_n\}$ после $[\alpha_1, \dots, \alpha_n]$.

К восприятиям мы добавим привязку к состоянию, в котором они имеют место, в виде $B_i = \{F_1, F_2, \dots, F_n\}$ в S_i .

Также мы хотели бы изменить форму закона действия, для унификации с законами причинных связей. Теперь мы будем их записывать их с условием, в виде

α меняет $\sim B$ на B при условии C .

Это позволит нам рассуждать о неудачах в выполнении действий без введения специального множества их ограничений. Отметим, что при этом выявляется интересная связь между действием и причинными связями: действие работает точно также, только в нем переключающим эффектом выступает не естественная причина, а акт, выполненный агентом.

3. Проблема квалификации

3.1. Семантика возможных миров

В данной главе мы рассмотрим то, как указанная теория действий может применяться в условиях неполной информации. Данный вопрос является одним из ключевых для теории действий потому, что большинство реальных человеческих действий выполняются именно в таких условиях. И проблемы, связанные с этим, составляют один из важнейших аспектов проблемы фреймов: аксиомы фреймов могут заставить нас сохранить для последующего состояния лишнюю информацию, которая уже успела измениться. И тогда выводы о результатах действий, и даже о возможности произвести то или иное действие, оказываются неверными.

Конкретно мы рассмотрим вопрос о том, как возможно рассуждение о результатах действий в случаях, когда у нас нет возможности получить полную информацию о среде; как можно посредством последовательности действий получить более полную информацию о среде; а также рассмотрим способы, позволяющие включить в наш универсум произвольные изменения среды (то есть происходящие не вследствие действий агента) и конкурентные действия (то есть произведенные другими агентами). Сразу отметим, что если раньше мы рассматривали понятие ситуации и состояния как синонимы, то сейчас будем в основном говорить в терминах ситуаций, а состоянию дадим несколько иное значение.

В целом мы начнем с некоторого интуитивного рассуждения, основанного на примерах, и после рассмотрения нескольких значимых на наш взгляд моментов, перейдем к более строгой формулировке полученных результатов.

Для начала рассмотрим пример²⁹. Агент находится в автомобиле, вставляет ключ в замок зажигания, и собирается завести двигатель. Ему известно действие «завести», обозначим его α , которое меняет состояние двигателя с «не работает» ($\sim A$) на «работает» (A). Но у этого действия есть

²⁹ Thielsher M. Challenges for action theories. Springer, 2000, pp. 86-87, авторская терминология не сохраняется

условие: в выхлопной трубе не должно быть никакого предмета, например картофелины («труба чиста» - В; «труба забита» - \sim В). Но в начальный момент времени нам известно только то, что двигатель не работает (восприятие в момент 1 состоит из одного элемента - $\{\sim A\}$). Какой вывод относительно результата мы можем сделать?

Очевидно, что отсутствие информации о состоянии выхлопной трубы заставляет нас рассмотреть два варианта: труба либо пуста, и тогда повернув ключ, мы получим работающий двигатель; либо труба забита, и тогда сколько бы мы не поворачивали ключ в замке, двигатель не заработает. Но как указать, что же ждать в следующий момент в рамках нашей модели представления действий?

Типичным способом решения данной проблемы является применение немонотонных рассуждений: мы предполагаем какой-либо вариант изначально, как вариант «по умолчанию», а какой-то оставляем как «аномальный». И если нам неизвестно, как обстоят дела с теми вещами, о которых нам ничего не известно, мы делаем вывод «по умолчанию». Если же нам вдруг становится известно, что выполнено какое-либо условие, делающее рассматриваемое действие невыполнимым, то мы пересматриваем наш вывод «по умолчанию» в свете новой информации. Однако проблемой такого подхода является то, что вывод «по умолчанию», хоть и может быть пересмотрен, является в определенной степени окончательным: мы не сомневаемся в нем. Это, в какой-то мере, является интуитивно понятным и близким к здравому смыслу человека: мы склонны ошибаться и не учитывать каких-либо факторов. Но нужно учесть, что немонотонные рассуждения предполагают, что нам известны условия, при которых действие не может быть выполнено (необходимо явно задать те наблюдения, которые могли бы заставить нас пересмотреть результат действия). То есть выходит, что мы знаем, что наше предприятие может не быть удачным, знаем из-за чего оно может не удалиться, не знаем точно, выполнены ли все условия, и все же ожидаем, что все получится. Такой подход, с точки зрения

здорового смысла, кажется слишком оптимистичным.

Мы хотим предложить другой подход, основанный на модальной логике. Точнее, он основан не на исчислении, а на семантике Возможных миров, которая очень удачно отражает сложившуюся ситуацию. И она позволяет ввести соответствующие модальные операторы, которые помогут нам выразить те дополнительные особенности наших ожиданий, которые возникают в ситуациях, когда нам не известна полная информация о среде. С их помощью мы сможем выразить недостоверность наших выводов относительно будущих состояний, а это, в свою очередь, даст нам возможность делать выводы относительно прошлых состояний. Причиной этого является тот факт, что неудавшееся действие мы не считаем аномалией, а считаем просто индикатором невыполнения условий данного действия, что позволяет нам получить больше информации о мире.

Сразу же обратим внимание на то, что неполнота информации о нашем универсуме проявляется только в отсутствии некоторых «эмпирических» данных, то есть данных о конкретных положениях дел. Общие же - «теоретические» - знания о мире у нас сохраняются. Мы предполагаем, что нам известны все законы мира: ограничения состояний, информация о влияниях, законы действий. Также изначально мы предполагаем, что нам известны все действия, выполненные после изначального состояния, и в нем нет никаких других изменений (не действует другой агент, не происходит каких-либо случайных изменений, не произведенных нами или не вызванных нашими действиями через цепочку причинных связей). Однако позднее мы обратим внимание на потенциал развития нашего подхода для рассуждения о моделях с конкурентными действиями.

Вернемся к нашему примеру. Мы обойдемся в данном случае без индивидов, в нашей модели будут только литералы $F=\{A,B\}$. У нас есть закон действия α :

α меняет $\sim A$ на A при условии B ;

Также в изначальной ситуации S_0 у нас есть восприятие $B_0 = \{\sim A\}$ в S_0 . Относительно значения литерала B нам ничего не известно, и мы предполагаем, что у нас нет возможности проверить его значение с помощью наблюдения. Есть два возможных значения данного литерала, и, соответственно, два возможных мира:

$$w_1 = \{\sim A, B\}; w_2 = \{\sim A, \sim B\}.$$

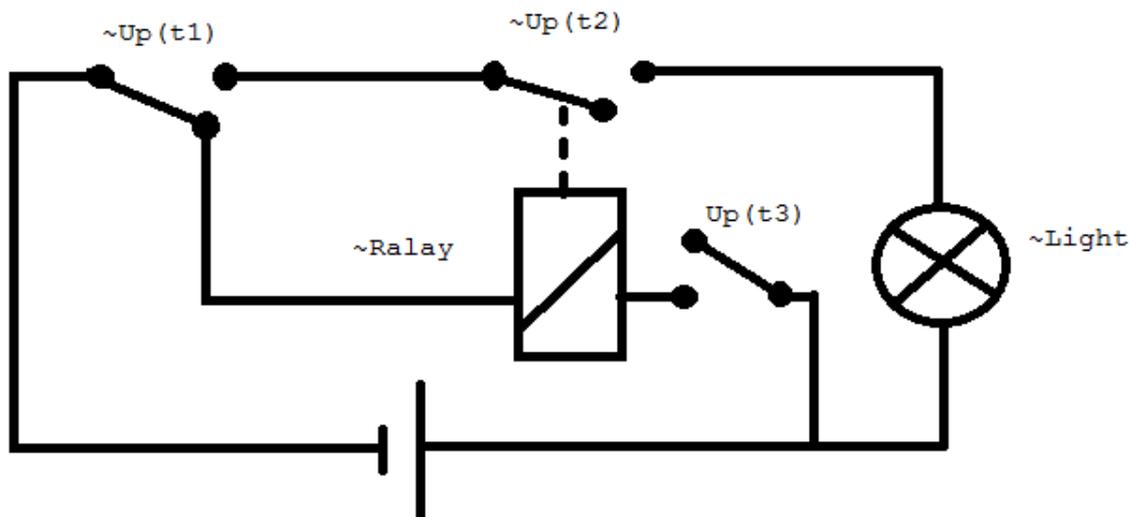
В первом случае действие α выполнимо и приведет к результату A . Во втором же случае действие α не приведет к ожидаемому результату. Соответственно, мы имеем два варианта последующих ситуаций:

$$w_1: S_1 = \{A, B\} \text{ после } [\alpha];$$

$$w_2: S_1 = \{\sim A, B\} \text{ после } [\alpha];$$

Здесь мы приводим действие α как выполненное даже во втором мире, так как планируем использовать информацию о том, что действие было произведено, но не привело к результату. Это ключевой момент для нашего подхода: информация о том, что действие не было выполнено, поможет нам сделать вывод о той ситуации, в которой оно было произведено. Так, мы можем теперь проверить, в каком же мире мы сейчас находимся, благодаря восприятиям: судя по данным в S_0 , информация о состоянии литерала A доступна нашим сенсорам. Определив то, в каком мире мы находимся, w_1 или w_2 , мы сможем сделать вывод о состоянии другого литерала, B , что мы не можем сделать напрямую. Причем его состояние определяется не только для S_1 , но и для S_0 , что является выводом «о прошлом». А возможность делать такие выводы - необходимая для удовлетворительного решения проблемы фреймов особенность, согласно тем требованиям, которые мы установили себе ранее.

Однако такая определенность возможна не во всех сценариях. Могут случаться такие ситуации, когда положение дел в текущем состоянии определено полностью, но на прошлое состояние это не проливает свет. Например рассмотрим электрическую цепь с тремя триггерами, лампой и реле:



В этой цепи замыкание подцепы с реле опускает триггер $t2$ в положение $\sim Up(t2)$. Это мы опишем законом причинной связи:

(1) Relay вызывает смену $Up(t2)$ на $\sim Up(t2)$ при условии \top ;

При этом относительно самого реле мы можем сказать:

(2) $\sim Up(t1)$ вызывает смену $\sim Relay$ на Relay при условии $\sim Up(t3)$;

(3) $\sim Up(t3)$ вызывает смену $\sim Relay$ на Relay при условии $\sim Up(t1)$;

(4) $Up(t1)$ вызывает смену Relay на $\sim Relay$ при условии \top ;

(5) $Up(t3)$ вызывает смену Relay на $\sim Relay$ при условии \top ;

Предположим, что нам не известно положение $t2$. Все остальное остается также, как изображено на рисунке: в ситуации $S0$ у нас имеется восприятие $B0 = \{\sim Up(t1), Up(t3), \sim Relay, \sim Light\}$ в $S0$. В соответствии с двумя возможными состояниями литерала $Up(t2)$ мы можем построить два возможных мира:

w1: $S0 = \{\sim Up(t1), Up(t3), \sim Relay, \sim Up(t2), \sim Light\}$ после $[\]$;

w2: $S0 = \{\sim Up(t1), Up(t3), \sim Relay, Up(t2), \sim Light\}$ после $[\]$;

Далее, предположим, мы произвели действие $\alpha(t3)$, подчиняющееся закону:

(0) $\alpha(x)$ меняет $Up(x)$ на $\sim Up(x)$ при условии \top ;

Рассмотрим **w1**:

Предварительное состояние-наследник S_0 , после применения закона действия $\alpha(x)$, будет представлять собой следующую пару предположений и эффектов:

$$(\{\sim Up(t1), \sim Up(t3), \sim Relay, \sim Up(t2), \sim Light\}, \{\sim Up(t3)\});$$

Далее мы видим, что, согласно имеющимся у нас знаниям о причинных связях (3), мы должны получить следующее:

$$(\{\sim Up(t1), \sim Up(t3), Relay, \sim Up(t2), \sim Light\}, \{\sim Up(t3), Relay\});$$

А вот эффект $Relay$, не смотря на наличие (1), не станет запускающим эффектом, потому что $Up(t2)$ не принадлежит нашему предварительному состоянию. Никакие другие причинные связи здесь также более не приложимы, поэтому мы получили окончательный вариант наследника для первого возможного мира:

w1: $S_1 = \{\sim Up(t1), \sim Up(t3), Relay, \sim Up(t2), \sim Light\}$ после $[\alpha]$

Теперь обратим внимание на второй возможный мир. Там мы получаем:

$$(\{\sim Up(t1), \sim Up(t3), \sim Relay, Up(t2), \sim Light\}, \{\sim Up(t3)\})$$
 по (0);
$$(\{\sim Up(t1), \sim Up(t3), Relay, Up(t2), \sim Light\}, \{\sim Up(t3), Relay\})$$
 по (3);
$$(\{\sim Up(t1), \sim Up(t3), Relay, \sim Up(t2), \sim Light\}, \{\sim Up(t3), Relay, \sim Up(t2)\})$$
 по (1);

Более никакие законы не применимы, и в итоге:

w2: $S_1 = \{\sim Up(t1), \sim Up(t3), Relay, \sim Up(t2), \sim Light\}$ после $[\alpha]$

Даже если теперь завесу тайны откроют, и мы сможем проверить состояние переключателя t_2 , это не поможет нам определить, в каком из миров мы находимся — они одинаковы. Разная у них только история. И никакой возможности определить, какова она была, при имеющейся информации нет.

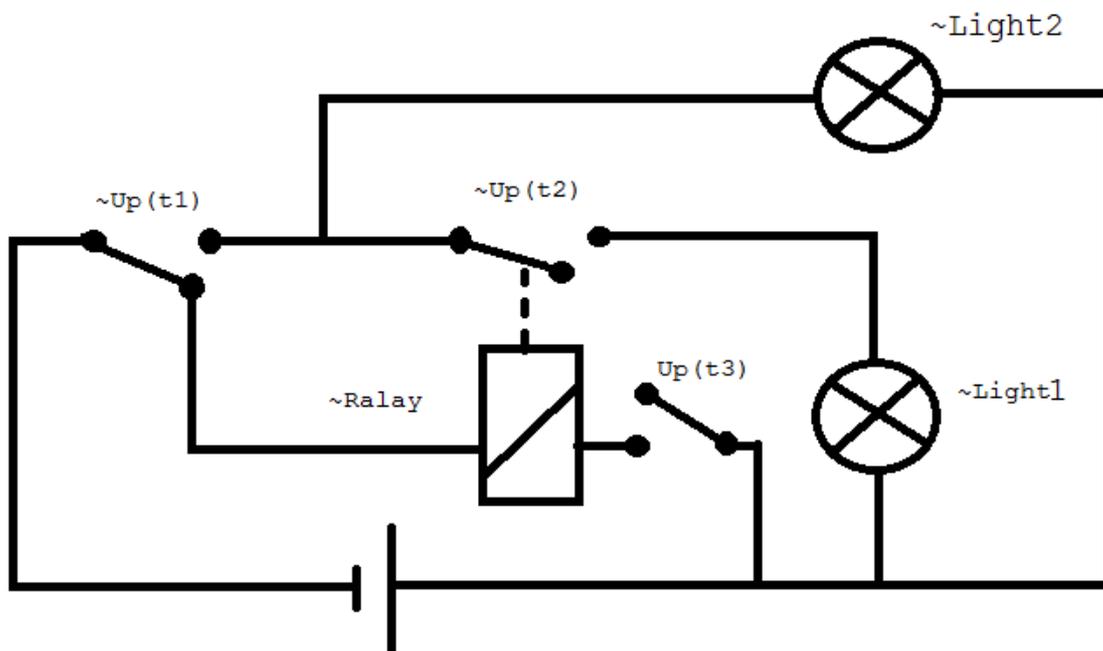
Сам по себе этот факт не компрометирует предложенный подход — интуитивные рассуждения нашего рационального сознания никак нам не помогут в разрешении вопроса об изначальном состоянии t_2 . Просто теперь становится ясным, что граф возможных миров, связанных отношением достижимости, не является деревом: ветви могут пересекаться с развитием истории нашего сценария.

Это требует особого рассмотрения вопроса о соотношении ситуации и возможного мира, а также ситуации и восприятия. Ситуации меняются с течением времени, но при этом история развивается в каком-либо из возможных миров. Однако истории, развивающиеся в нескольких возможных мирах, могут пересекаться в каком-либо состоянии. Что при этом происходит с возможными мирами?

Они сливаются в один. Так как все условия действий и причинных связей рассматриваются только в данной ситуации, разные истории, приведшие к одной и той же ситуации, уже не могут быть дифференцированы. И из объединяющего два (или более) мира состояния начинается новая история, оно выступает как новое изначальное состояние (и хотя за этим может последовать новое расщепление, оно будет производиться уже на другом основании). С одной стороны это может показаться минусом, т. к. мы не можем получить полной информации о сценарии. Однако с другой стороны, это является плюсом: мы можем специально искать такие последовательности действий, которые приведут нас к слиянию разных историй в одну, чтобы в дальнейшем не испытывать неудобств, могущих возникнуть из неполноты знаний о среде. Конечно, это возможно не в любом сценарии. Также необязательно срастание каких-либо возможных миров означает полную определенность в ситуации: могут остаться другие возможные миры, которые отличаются от объединившихся теми условиями и следствиями, которые не были затронуты предпринятыми действиями. Но это не мешает приобрести определенность хоть в чем-то (особенно если это наиболее важно и принципиально для достижения поставленной цели).

Такую стратегию можно назвать «обнулением»: она минимизирует неполноту информации о среде (различия между возможными мирами) за счет изменения изначального состояния. Когда мы подходим к шахматной доске, за которой кто-то играл, не закончил партию и не расставил, мы можем поступить двумя способами: продолжить начатую игру, или расставить фигуры по местам

и начать все с начала. Так же, оказавшись в незнакомой среде, не имея средств для её исследования (или имея, но более трудоемкие и неадекватные нашим целям), но будучи способными привести эту среду к некоему «начальному» состоянию, мы можем просто сделать это, а потом начать действовать в ней с чистого листа. Возможно это тогда, когда у нас в наличии есть какое-либо действие, имеющее «отменяющий»³⁰ эффект: его выполнение в одном случае (одном возможном мире) ничего не меняет (в интересующем нас плане, не делает хуже), а в другом (других) — меняет в нужную нам сторону. Добавим в нашу цепь еще одну лампу, и условимся, что наша задача — зажечь только её, при этом не зажигая первую лампу:



Приведенные ранее законы причинных связей продолжают работать и в этой цепи, но к ним нужно добавить законы относительно $Light1$ и $Light2$:

- (6) $Up(t1)$ вызывает смену $\sim Light2$ на $Light2$ при условии \top ;
- (7) $Up(t1)$ вызывает смену $\sim Light1$ на $Light1$ при условии $Up(t2)$;
- (8) $Up(t2)$ вызывает смену $\sim Light1$ на $Light1$ при условии $Up(t1)$;
- (9) $\sim Up(t1)$ вызывает смену $Light1$ на $\sim Light1$ при условии \top ;

³⁰ Использование слова «отменяющий» не предполагает, что эффект (обычно не прямой) от действия будет изменением значения литерала с истины на ложь. Отмена здесь может быть более сложной конструкции — какого-либо обстоятельства (выразимого не одним литералом), невыгодного нам.

(10) $\sim\text{Up}(t1)$ вызывает смену Light2 на $\sim\text{Light2}$ при условии \top ;

(11) $\sim\text{Up}(t2)$ вызывает смену Light1 на $\sim\text{Light1}$ при условии \top ;

Изначальная ситуация остается той же, с добавлением информации о второй лампе:

B0 = { $\sim\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \sim\text{Light1}, \sim\text{Light2}$ } в **S0**;

w1: **S0** = { $\sim\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ } после [];

w2: **S0** = { $\sim\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ } после [];

Если сразу применить действие $\alpha(t1)$, то в первом случае мы получим желаемый результат, где горит только вторая лампа:

({ $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ }, { $\text{Up}(t1)$ }) по (0);

({ $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \text{Light2}$ }, { $\text{Up}(t1), \text{Light2}$ }) по (6);

Литерал Light2 не связан ни с каким другим законами причинности, для другого закона с $\text{Up}(t1)$ в качестве переключающего эффекта не соблюдено условие $\text{Up}(t2)$. Поэтому мы получаем итоговое состояние:

w1: **S1** = { $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \text{Light2}$ } после [$\alpha(t1)$];

Но во втором возможном мире выполнение того же действия приведет к нежелательным последствиям — загорятся обе лампы:

({ $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ }, { $\text{Up}(t1)$ }) по (0);

({ $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \text{Up}(t2), \sim\text{Light1}, \text{Light2}$ }, { $\text{Up}(t1), \text{Light2}$ }) по (6);

({ $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \text{Up}(t2), \text{Light1}, \text{Light2}$ }, { $\text{Up}(t1), \text{Light2}, \text{Light1}$ }) по (7);

Более никакие законы не применимы.

w2: **S1** = { $\text{Up}(t1), \text{Up}(t3), \sim\text{Relay}, \text{Up}(t2), \text{Light1}, \text{Light2}$ } после [$\alpha(t1)$];

Если же перед тем, как совершить действие $\alpha(t1)$, произвести «обнуление» посредством $\alpha(t3)$, то мы достигнем цели в любом случае:

w1: ({ $\sim\text{Up}(t1), \sim\text{Up}(t3), \sim\text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ }, { $\sim\text{Up}(t3)$ }) по (0);

({ $\sim\text{Up}(t1), \sim\text{Up}(t3), \text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ }, { $\sim\text{Up}(t3), \text{Relay}$ }) по (3);

S1 = { $\sim\text{Up}(t1), \sim\text{Up}(t3), \text{Relay}, \sim\text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ } после [$\alpha(t3)$];

w2: ({ $\sim\text{Up}(t1), \sim\text{Up}(t3), \sim\text{Relay}, \text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ }, { $\sim\text{Up}(t3)$ }) по (0);

({ $\sim\text{Up}(t1), \sim\text{Up}(t3), \text{Relay}, \text{Up}(t2), \sim\text{Light1}, \sim\text{Light2}$ }, { $\sim\text{Up}(t3), \text{Relay}$ }) по (3);

$(\{\sim\text{Up}(t1),\sim\text{Up}(t3),\text{Relay},\sim\text{Up}(t2),\sim\text{Light1},\sim\text{Light2}\},\{\sim\text{Up}(t3),\text{Relay},\sim\text{Up}(t2)\})$

по (1);

$\mathbf{S1}=\{\sim\text{Up}(t1),\sim\text{Up}(t3),\text{Relay},\sim\text{Up}(t2),\sim\text{Light1},\sim\text{Light2}\}$ после $[\alpha(t3)]$;

Так как отличий в $\mathbf{S1}$ в различных мирах нет, мы можем продолжить без раздвоения историй. Состоянием-наследником $\mathbf{S1}$ после выполнения $\alpha(t1)$ будет наше целевое состояние:

$(\{\text{Up}(t1),\sim\text{Up}(t3),\text{Relay},\sim\text{Up}(t2),\sim\text{Light1},\sim\text{Light2}\},\{\text{Up}(t1)\})$ по (0);

$(\{\text{Up}(t1),\sim\text{Up}(t3),\sim\text{Relay},\sim\text{Up}(t2),\sim\text{Light1},\sim\text{Light2}\},\{\text{Up}(t1),\sim\text{Relay}\})$ по (4);

$(\{\text{Up}(t1),\sim\text{Up}(t3),\sim\text{Relay},\sim\text{Up}(t2),\sim\text{Light1},\text{Light2}\},\{\text{Up}(t1),\sim\text{Relay},\text{Light2}\})$ по (6);

$\mathbf{S2}=\{\text{Up}(t1),\sim\text{Up}(t3),\sim\text{Relay},\sim\text{Up}(t2),\sim\text{Light1},\text{Light2}\}$ после $[\alpha(t3),\alpha(t1)]$;

Проиллюстрированная стратегия имеет определенное значение для вопросов планирования действий, которые непосредственно связаны с проблемой квалификации, однако не являются основной темой данной работы (хотя немного о них еще будет сказано позднее). Сейчас же обратимся ко второму моменту, озвученному выше: соотношению ситуации и восприятия.

В рассматриваемых до данной главы сценариях разница между этими понятиями была не совсем понятна: и то и другое представляет из себя множество литералов, только ситуация дополнена последовательностью действий, приведших к ней; а восприятие — указанием на ситуацию, в которой оно наблюдалось. Восприятию мы уделяли внимание только в изначальной ситуации. Собственно, первое задавало последнюю, и эквивалентность их не позволяла четко выявить различие. Однако различие их интуитивных смыслов мы обозначали: восприятия — это данные сенсоров относительно среды. А ситуация — это полное описание среды в терминах нашей модели.

Неполнота информации, порождающая проблему квалификации, является неполнотой восприятия: мы по тем или иным причинам не можем получить актуальную информацию о значениях всех литералов. И тут появляется различие: восприятие оказывается неполным описанием среды, в отличии от

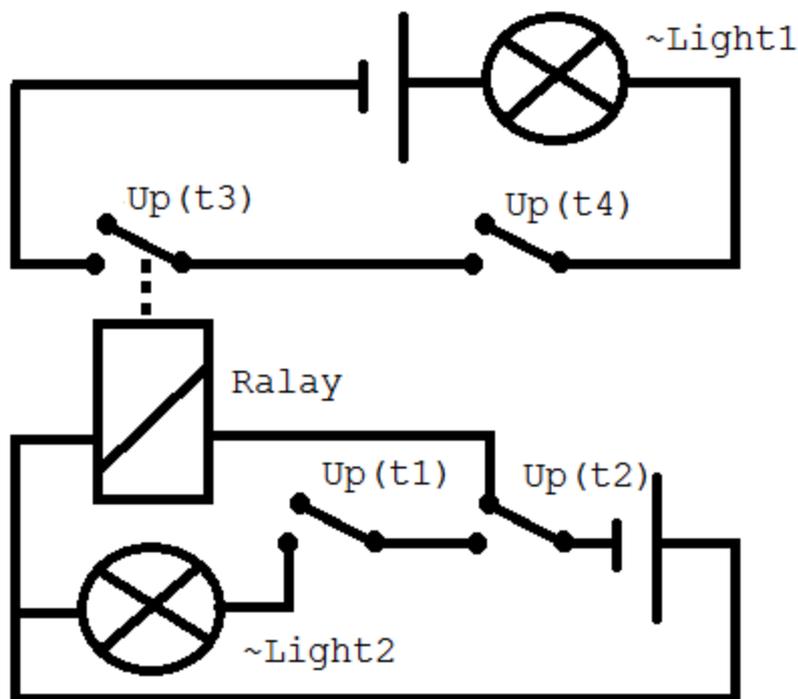
ситуации, для которой мы требовали, чтобы каждый литерал входил в нее, и только один раз, либо со значением «истина», либо со значением «ложь». Восприятие продолжает играть важную роль для формирования ситуаций — данным восприятия мы доверяем, и считаем это той сердцевиной ситуации, которая нам известна в полной мере. Второй важный момент — это ограничения состояний. Их мы считаем известными полностью (об этом будет еще сказано несколько слов в заключении) и полностью достоверными. Поэтому следствия, получаемые из подстановки данных восприятия в ограничения состояния мы также будем считать достоверными. А те литералы, которые не получили достоверного истинного значения после операций получения восприятия и выводов из них согласно ограничениям, заставляют нас формировать несколько ситуаций, перебирающих все согласующиеся с достоверными данными варианты полного описания среды. Так мы получаем возможные миры с различными ситуациями, соответствующими одному набору воспринятых данных.

Далее отметим, что ситуации же мы получали в результате рассуждения о том, каков будет результат того или иного действия. Именно получение ожидаемых полных описаний среды было нашей задачей. Как эти ситуации, получаемые посредством вывода, соотносятся с восприятиями? Ситуации испытывают определенное влияние и вывода, и восприятия. Но всегда ли они совпадают?

Вывод основан на законах причинности и законах действий. Мы полагаем их здесь приемлемыми в том смысле, что они никогда не могут дать нам вывода, противоречащего восприятию («фактам»). В сценариях с полной информацией, совпадение ситуации, выведенной из предыдущей согласно законам действия и причинных связей, и ситуации, заданной восприятием, должно иметь место (иначе мы сочтем нашу модель неудачной и требующей уточнения). В сценариях с неполной информацией посредством вывода мы можем получить столько же ситуаций, сколько имели возможных миров (для

одного действия). Результат восприятия же, как мы видели, также задает множество возможных ситуаций, которые принадлежат разным возможным мирам. Должны ли совпадать множества ситуаций, построенных посредством вывода, и заданных восприятием?

Нет, совпадать они не должны. Для того, чтобы не разочароваться в построенной модели среды (законах действий и причинных связей, введенных ограничениях состояний), нам достаточно того, чтобы эти множества пересекались. И пересечение даже окажется желательным — оно принесет больше информации о среде. Поясним сказанное на примере. Рассмотрим следующую пару электрических цепей:



Ограничениями состояний здесь будут выражения:

- (a) $Light1 \equiv \sim Up(t3) \ \& \ \sim Up(t4)$
- (b) $Light2 \equiv \sim Up(t1) \ \& \ \sim Up(t2)$
- (c) $Ralay \equiv Up(t2)$
- (d) $Ralay \rightarrow \sim Up(t3)$

Из них (с помощью информации о влияниях) мы можем получить следующие законы причинных связей:

- (1) Relay вызывает смену $Up(t3)$ на $\sim Up(t3)$ при условии \top ;
- (2) $Up(t2)$ вызывает смену $\sim Relay$ на $Relay$ при условии \top ;
- (3) $\sim Up(t2)$ вызывает смену $Relay$ на $\sim Relay$ при условии \top ;
- (4) $\sim Up(t3)$ вызывает смену $\sim Light1$ на $Light1$ при условии $\sim Up(t4)$;
- (5) $\sim Up(t4)$ вызывает смену $\sim Light1$ на $Light1$ при условии $\sim Up(t3)$;
- (6) $\sim Up(t1)$ вызывает смену $\sim Light2$ на $Light2$ при условии $\sim Up(t2)$;
- (7) $\sim Up(t2)$ вызывает смену $\sim Light2$ на $Light2$ при условии $\sim Up(t1)$;
- (8) $Up(t1)$ вызывает смену $Light2$ на $\sim Light2$ при условии \top ;
- (9) $Up(t2)$ вызывает смену $Light2$ на $\sim Light2$ при условии \top ;
- (10) $Up(t3)$ вызывает смену $Light1$ на $\sim Light1$ при условии \top ;
- (11) $Up(t4)$ вызывает смену $Light1$ на $\sim Light1$ при условии \top ;

Также мы используем закон действия $\alpha(x)$:

- (0) $\alpha(x)$ меняет $\sim Up(x)$ на $Up(x)$ при условии \top ;

Будем считать, что сенсоры агента способны воспринимать только состояние ламп и переключателя $t2$.

$\mathbf{B0} = \{\sim Up(t2), \sim Light1, Light2\}$ в $\mathbf{S0}$;

Учитывая ограничения состояний, мы должны признать $\sim Up(t1)$, и $\sim Relay$. Относительно $Up(t3)$ и $Up(t4)$ мы получаем только информацию, что они не могут быть оба ложны в силу $\sim Light1$ и (a). Отсюда:

$\mathbf{w1: S0} = \{\sim Up(t1), \sim Up(t2), Up(t3), \sim Up(t4), \sim Relay, \sim Light1, Light2\}$ после [];

$\mathbf{w2: S0} = \{\sim Up(t1), \sim Up(t2), \sim Up(t3), Up(t4), \sim Relay, \sim Light1, Light2\}$ после [];

$\mathbf{w3: S0} = \{\sim Up(t1), \sim Up(t2), Up(t3), Up(t4), \sim Relay, \sim Light1, Light2\}$ после [];

После действия $\alpha(t2)$ получаем:

$\mathbf{w1: } (\{\sim Up(t1), Up(t2), Up(t3), \sim Up(t4), \sim Relay, \sim Light1, Light2\}, \{Up(t2)\})$ по (0);

$(\{\sim Up(t1), Up(t2), Up(t3), \sim Up(t4), Relay, \sim Light1, Light2\}, \{Up(t2), Relay\})$ по (2);

$(\{\sim Up(t1), Up(t2), Up(t3), \sim Up(t4), Relay, \sim Light1, \sim Light2\},$

$\{Up(t2), Relay, \sim Light2\}$) по (9);

$(\{\sim Up(t1), Up(t2), \sim Up(t3), \sim Up(t4), Relay, \sim Light1, \sim Light2\},$

$\{Up(t2), Relay, \sim Light2, \sim Up(t3)\})$ по (1);

$(\{\sim Up(t1), Up(t2), \sim Up(t3), \sim Up(t4), Relay, Light1, \sim Light2\},$

$\{Up(t2), Relay, \sim Light2, \sim Up(t3), Light1\})$ по (4);

S1 = $\{\sim Up(t1), Up(t2), \sim Up(t3), \sim Up(t4), Relay, Light1, \sim Light2\}$ после $[\alpha(t2)]$;

w2: $(\{\sim Up(t1), Up(t2), \sim Up(t3), Up(t4), \sim Relay, \sim Light1, Light2\}, \{Up(t2)\})$ по (0);

$(\{\sim Up(t1), Up(t2), \sim Up(t3), Up(t4), Relay, \sim Light1, Light2\}, \{Up(t2), Relay\})$ по

(2);

$(\{\sim Up(t1), Up(t2), \sim Up(t3), Up(t4), Relay, \sim Light1, \sim Light2\},$

$\{Up(t2), Relay, \sim Light2\})$ по (9);

S1 = $\{\sim Up(t1), Up(t2), \sim Up(t3), Up(t4), Relay, \sim Light1, \sim Light2\}$ после $[\alpha(t2)]$;

w3: $(\{\sim Up(t1), Up(t2), Up(t3), Up(t4), \sim Relay, \sim Light1, Light2\}, \{Up(t2)\})$ по (0);

$(\{\sim Up(t1), Up(t2), Up(t3), Up(t4), Relay, \sim Light1, Light2\}, \{Up(t2), Relay\})$ по

(2);

$(\{\sim Up(t1), Up(t2), Up(t3), Up(t4), Relay, \sim Light1, \sim Light2\},$

$\{Up(t2), Relay, \sim Light2\})$ по (9);

$(\{\sim Up(t1), Up(t2), \sim Up(t3), Up(t4), Relay, \sim Light1, \sim Light2\},$

$\{Up(t2), Relay, \sim Light2, \sim Up(t3)\})$ по (1);

S1 = $\{\sim Up(t1), Up(t2), \sim Up(t3), Up(t4), Relay, \sim Light1, \sim Light2\}$ после $[\alpha(t2)]$;

Далее предположим, что мы получили новое восприятие:

B1 = $\{Up(t2), Light1, \sim Light2\}$ в **S1**;

Из него мы можем сделать заключение, что в то время, которое соответствует **S1**, имеют место также $\sim Up(t3)$, $\sim Up(t4)$ и $Relay$. Но значение $Up(t1)$ не выводится из восприятия и ограничений состояний. В связи с этим возможно построить два возможных мира с разными ситуациями:

w'1: **S'1** = $\{Up(t1), Up(t2), \sim Up(t3), \sim Up(t4), Relay, Light1, \sim Light2\}$ после $[\]$;

$w'2: S'1 = \{\sim Up(t1), Up(t2), \sim Up(t3), \sim Up(t4), Relay, Light1, \sim Light2\}$ после [];³¹

Итак, мы видим, что ситуации в двух случаях получились разные. Отличается даже их количество для разных способов построения. Однако есть случай, когда они совпадают — $w1$ и $w'1$. Это то самое пересечение двух построенных разными способами систем возможных миров (в данном случае представленное единичным множеством), которое позволяет сохранить определенную степень доверия к нашей модели сценария.

Самым важным в этой ситуации является то, что степень неизвестности при несовпадении возможных миров снижается. В приведенном примере неизвестность в принципе была исключена: $S1$ оказалась определена однозначно соотношением ситуаций, возможных с разных точек зрения. Но даже если в результате такого сопоставления мы получим не единичное множество миров, возможных с обеих точек зрения, оно все же будет меньше изначальных. В случае же совпадения систем возможных миров степень информированности агента не меняется.

Отдельно стоит подчеркнуть, что в случаях, когда отбрасываются возможные миры с ситуациями, выведенными из действия, информированность возрастает не только относительно настоящего, но и относительно прошлого. Указанные миры имеют историю (в отличие от возможных миров восприятий), и снижение их количества снижает также степень неопределенности относительно прошлого. В нашем примере выяснилось, что изначально имела место ситуация $\{\sim Up(t1), \sim Up(t2), Up(t3), \sim Up(t4), \sim Relay, \sim Light1, Light2\}$ ($S0$ в мире $w1$). Выше мы уже замечали, что предлагаем систему, позволяющую делать выводы «о прошлом». Сейчас мы подробнее рассмотрим что для этого требуется. Заметим, что снова перед нами встает вопрос, близкий к тематике планирования: какое действие нужно выбрать, чтобы снизить неопределенность?

³¹ Мы оставляем здесь квадратные скобки с последовательностью действий пустыми для того, чтобы подчеркнуть различие между ситуациями, выведенными из действий, и ситуациями, выведенными из восприятий.

Восприятия в последующих (не изначальном) состояниях позволяют производить проверку наших предположений. Это не прямая проверка данных — мы предполагаем, что возможности наших сенсоров с течением времени (развитием сценария) не меняются. То есть информация, которая не была доступна им в изначальном состоянии не становится в последующих доступной напрямую. Однако мы можем предпринять какие-либо действия, которые позволят нам сделать вывод о предыдущем состоянии среды. Это действие должно приводить к одному доступному проверке состоянию-наследнику в одном возможном мире, и к другому — в другом. Проверая данные в наследующем состоянии мы тем самым определяем, в каком возможном мире мы оказались.

Описанная ситуация возможна тогда, когда в изначальном состоянии мы не имеем точной информации о литералах, составляющих условия какого-либо действия или причинных связей, эффект от которых — доступный нашему восприятию литерал. Тогда в случае, если действие успешно (или закон причинной связи сработал), мы можем заключить, что условия были соблюдены. Если же нет — условия соблюдены не были. Причем заметим, что в случае причинных связей может быть две причины того, что она не была актуализирована: либо не выполнено условие, либо подвергаемый воздействию литерал находится не в том состоянии, которое указывается в законе первым (изначальным). Для действий такого не случается, т. к. мы предполагаем, что тот литерал, который доступен манипулированию, дан нам в восприятии.

Возможность проверки присутствует не всегда. Она отсутствует, например, в ситуациях, где мы не можем воздействовать на доступные наблюдению литералы, или воздействие это не зависит от неизвестных нам изначально литералов. Тогда никакие возможные миры не будут вычеркнуты, и просто продолжат развиваться параллельно. Также нет возможности проверки прошлого в рассмотренном ранее варианте срастающихся возможных миров: там никакие наблюдения не позволят определить прошлое, т. к. история миров

«обнуляется». Стратегия «обнуления» противоположна по своей сути стратегии «критического эксперимента», которая основывается на проверке восприятием. В первом случае, действуя, мы стремимся к снижению различий между ожидаемыми состояниями-наследниками. Во втором же, наоборот, сохраняем различия, и стремимся сделать их явными, доступными нашему восприятию. В состояниях, построенных гипотетически (перебором непротиворечивых и удовлетворяющих первичному восприятию вариантов), различия не явны — они в той части состояния, которая не доступна нашему восприятию. Суть стратегии «критического эксперимента» в том, что мы подбираем такое действие, которое сделает явным это различие. Оно должно влиять (прямо, или через посредничество причинных связей) на то, что доступно восприятию; и результат (самого действия, или применения причинных связей) должен зависеть от значения тех литералов, значение которых нам не известно точно.

Проще всего проверять условия выполнимости действия. Это делается в один шаг: мы выполняем действие, и проверяем — привело ли оно к предполагаемому результату. Если нет — значит какое-то из условий не выполнено; если да — все условия выполнены. Такой пример мы приводили в начале главы, где из того, завелся двигатель или нет после поворота ключа зажигания, делали вывод о том, есть ли в выхлопной трубе картофелялина.

Отметим отличие связи условий действия и его результата от простой имплицативной связи: если записать её так: $У \ \& \ Д \rightarrow Р$ (где $У$ означает высказывание «условия соблюдены»; $Д$ - «действие выполнено»; $Р$ - «ожидаемый результат достигнут»), то из $Д \ \& \ Р$ мы не сможем вывести значение $У$. Только из $Д \ \& \ \sim Р$ мы можем получить достоверную информацию о значении $У$, что $У =$ «ложь». Это отличие логического следования от следования в случае причинных связей (и действий как их подвида, где причиной выступает агент) и порождает проблему ветвления. Именно потому мы не могли обойтись в этой работе без удовлетворительного решения проблемы ветвления — нам нужна была подходящая база для обеспечения корректных выводов в

прошлое в условиях неполной информации.

Вывод о прошлом посредством законов причинных связей является более сложным, и подвержен определенным ограничениям. Это связано с тем, что мы можем не иметь информации о том, в каком состоянии находится изменяемый литерал, и даже литерал, инициирующий действие закона причинности. Рассмотрим следующий пример:

(0) $\alpha(A)$ меняет $\sim A$ на A при условии C ;

(1) A вызывает смену $\sim B$ на B при условии D ;

$\mathbf{B0} = \{\sim A, \sim B\}$;

Мы можем сделать вывод о том, что с ограничениями состояний точно совместимо:

(a) $A \& D \rightarrow B$;

Это ограничение состояний, может быть, не будут полными. Но так как мы сейчас пытаемся рассмотреть наиболее общий случай, нам этого будет достаточно. Они дают нам исключить из рассмотрения как противоречивые все ситуации, где встречаются $\{A, \sim B, D\}$.

Возможных миров получается 4, в соответствии с возможными непротиворечивыми распределениями истинностных значений недоступных восприятию литералов:

$\mathbf{w1: S0} = \{\sim A, \sim B, C, D\}$ после $[\]$;

$\mathbf{w2: S0} = \{\sim A, \sim B, C, \sim D\}$ после $[\]$;

$\mathbf{w3: S0} = \{\sim A, \sim B, \sim C, D\}$ после $[\]$;

$\mathbf{w4: S0} = \{\sim A, \sim B, \sim C, \sim D\}$ после $[\]$;

И, соответственно:

$\mathbf{w1: S1} = \{A, B, C, D\}$ после $[\alpha(A)]$; 1; 2

$\mathbf{w2: S1} = \{A, \sim B, C, \sim D\}$ после $[\alpha(A)]$; 1

$\mathbf{w3: S1} = \{\sim A, \sim B, \sim C, D\}$ после $[\alpha(A)]$;

$\mathbf{w4: S1} = \{\sim A, \sim B, \sim C, \sim D\}$ после $[\alpha(A)]$;

Что будет, если сравнить эти гипотетические ситуации с восприятием в

S1? Если $\mathbf{B1} = \{A, B\}$, то остается только один возможный мир — $\mathbf{w1}$. Это тот возможный мир, который испытал действие всех имеющихся законов причинных связей, и тем самым может подтвердить, что все условия выполнены. То есть, если вывод ситуации заканчивается изменением восприятия, мы можем сделать вывод, что все условия, необходимые для такой ситуации, были выполнены. Это еще раз показывает важность возможности вывода от подтверждения, а не от опровержения предположения. Особенно на фоне того, как ведет себя вывод от опровержения. Так, $\mathbf{B1} = \{\sim A, \sim B\}$ не дает нам определить мир, в котором мы находимся, однозначно. Этому восприятию соответствуют два мира: $\mathbf{w3}$ и $\mathbf{w4}$. Причем заметим, что информацию, которую мы получаем в этом случае, мы могли бы получить и без закона причинной связи между A и B. Для вывода здесь нам достаточно закона действия.

Интересной ситуацией здесь является случай $\mathbf{B1} = \{A, \sim B\}$. Здесь отрицание следствия закона причинности все же дает нам точную картину. Но причиной этого является простота примера. Если мы добавим еще один шаг причинно-следственных связей между прямым результатом действия и наблюдаемым результатом, то ситуация изменится:

- (0) $\alpha(A)$ меняет $\sim A$ на A при условии D ;
 - (1) A вызывает смену $\sim B$ на B при условии E ;
 - (2) B вызывает смену $\sim C$ на C при условии F ;
- $\mathbf{B0} = \{\sim A, \sim C\}$;

Ограничения дают нам исключить из рассмотрения ситуации, где встречаются $\{A, \sim B, E\}$ и $\{B, \sim C, F\}$:

- (a) $A \& E \rightarrow B$;
- (b) $B \& F \rightarrow C$;

Возможных миров получается 12:

- $\mathbf{w1: S0} = \{\sim A, \sim B, \sim C, D, E, F\}$ после $[\]$;
- $\mathbf{w2: S0} = \{\sim A, B, \sim C, D, E, \sim F\}$ после $[\]$;
- $\mathbf{w3: S0} = \{\sim A, \sim B, \sim C, D, E, \sim F\}$ после $[\]$;

- w4:** $S_0 = \{\sim A, \sim B, \sim C, \sim D, E, F\}$ после [];
- w5:** $S_0 = \{\sim A, B, \sim C, \sim D, E, \sim F\}$ после [];
- w6:** $S_0 = \{\sim A, \sim B, \sim C, \sim D, E, \sim F\}$ после [];
- w7:** $S_0 = \{\sim A, \sim B, \sim C, D, \sim E, F\}$ после [];
- w8:** $S_0 = \{\sim A, B, \sim C, D, \sim E, \sim F\}$ после [];
- w9:** $S_0 = \{\sim A, \sim B, \sim C, D, \sim E, \sim F\}$ после [];
- w10:** $S_0 = \{\sim A, \sim B, \sim C, \sim D, \sim E, F\}$ после [];
- w11:** $S_0 = \{\sim A, B, \sim C, \sim D, \sim E, \sim F\}$ после [];
- w12:** $S_0 = \{\sim A, \sim B, \sim C, \sim D, \sim E, \sim F\}$ после [];

Откуда

- w1:** $S_1 = \{A, B, C, D, E, F\}$ после [$\alpha(A)$]; 123
- w2:** $S_1 = \{A, B, \sim C, D, E, \sim F\}$ после [$\alpha(A)$]; 1;
- w3:** $S_1 = \{A, B, \sim C, D, E, \sim F\}$ после [$\alpha(A)$]; 12
- w4:** $S_1 = \{\sim A, \sim B, \sim C, \sim D, E, F\}$ после [$\alpha(A)$];
- w5:** $S_1 = \{\sim A, B, \sim C, \sim D, E, \sim F\}$ после [$\alpha(A)$];
- w6:** $S_1 = \{\sim A, \sim B, \sim C, \sim D, E, \sim F\}$ после [$\alpha(A)$];
- w7:** $S_1 = \{A, \sim B, \sim C, D, \sim E, F\}$ после [$\alpha(A)$]; 1
- w8:** $S_1 = \{A, B, \sim C, D, \sim E, \sim F\}$ после [$\alpha(A)$]; 1
- w9:** $S_1 = \{A, \sim B, \sim C, D, \sim E, \sim F\}$ после [$\alpha(A)$]; 1
- w10:** $S_1 = \{\sim A, \sim B, \sim C, \sim D, \sim E, F\}$ после [$\alpha(A)$];
- w11:** $S_1 = \{\sim A, B, \sim C, \sim D, \sim E, \sim F\}$ после [$\alpha(A)$];
- w12:** $S_1 = \{\sim A, \sim B, \sim C, \sim D, \sim E, \sim F\}$ после [$\alpha(A)$];

Здесь восприятие $\mathbf{B}_0 = \{A, \sim C\}$, показывающее, что действие было выполнено, но не привело к видимому изменению среды через причинные связи, не дает сделать однозначного вывода о том, почему так случилось. Ему соответствуют миры **w2**, **w3**, **w7**, **w8**, **w9**. В **w7** и **w9** не выполнено условие E, для смены $\sim B$ на B. В **w8**, хоть это условие и не выполнено, B является истинным, и потому применение закона не возможно. Однако, так как B здесь было истинно изначально, этот факт не может послужить поводом для

применения закона причинной связи (2). Потому же неприменимы законы действий (1) и (2) в мире w_2 , хоть там и истинно условие первой причинной связи. Мир w_2 - самый близкий к состоянию определенности, однако знать наверняка, находимся мы в нем, или нет, мы не можем.

Итак, данные о результате действия делят количество возможных миров пополам. А вот данные о результатах применения законов причинности позволяют нам либо узнать наверняка, либо остаться в множестве остальных $(n/2)-1$ исходов (где n – кол-во возможных миров в S_0). Обратим внимание, что при добавлении шагов между прямым результатом действия и доступным для проверки непрямым, количество возможных миров будет увеличиваться, и новое восприятие также либо определит состояние точно, либо оставит нас среди $(n/2)$ или $(n/2)-1$ (в зависимости от успешности действия) возможных миров.

Обратим внимание также на то, что в данном примере произошло слияние двух возможных миров — w_2 и w_3 . Они идентичны друг другу в S_1 .

Полученные результаты — это худший вариант, так как здесь не строгие ограничения состояний, восприятию доступно минимальное количество литералов, рассмотрено только одно действие. С другой стороны, здесь условия являются единичными литералами, а если они состоят из конъюнкции хотябы двух литералов, то информация о том, что не выполнено именно условие, оставляла бы нам три возможных мира: где ложны оба, где ложен только первый, или где ложен только второй. Это касается и выводов по законам действий.

Но, в общем случае, мы все же можем говорить о том, что количество информации о среде после проверки увеличивается, и если выбирать подходящие действия, то проверка может принести полное знание о среде. Лучшими действиями будут те, которые влекут точное определение значения того или иного литерала. То есть это либо действия, результаты которых могут указать значение условий их выполнения; либо действия, для не прямых

эффектов которых определены два литерала из трех возможных неизвестных — условия выполнения, наличие эффекта-переключателя, и изначальное состояние изменяемого литерала. Если мы подберем такое действие, позволяющее посредством проверки точно сделать вывод о состоянии множества (возможно единичного) литералов, то это и будет «критический эксперимент».

В заключении рассуждения о связи ситуаций и восприятий, определим более точно их взаимосвязь. *Гипотетическими ситуациями* мы назовем все ситуации, полученные просто перебором всех согласующихся с ограничениями ситуаций, а после выполнения действий — еще с законами данных действий и причинных связей, множеств литералов, представляющих собой полные описания среды. Восприятия остаются множеством тех фактических данных о среде, которые агент получает посредством наблюдения. Расширенным восприятием мы назовем то множество литералов, которое получается подстановкой данных восприятия в ограничения ситуаций и последующим выводом. *Проверка* — процесс сопоставления данных восприятия и гипотетических ситуаций, в котором не соответствующие восприятию ситуации отбрасываются. Итогом проверки станет множество (возможно единичное) *фактических ситуаций*. Таким термином мы называем те гипотетические ситуации, которые согласуются с фактами, данными сенсоров. Отметим, что их все же может быть несколько, и каждая из них потому несет некоторый характер необязательности. Подлежащую реальную ситуацию мы будем называть *действительной* (по аналогии с действительным миром, выделенным по какому-то принципу среди возможных).

Теперь опишем получившуюся модель более строго. Как и ранее, домен будет состоять из общих, не зависящих от ситуации знаний. $D = (O, F, A, AL, C, I)$, где O - множество индивидов, F - множество функций, A - множество действий, AL - множество законов действий, C - множество ограничений состояний и I - множество упорядоченных пар, задающих информацию о

влияниях. Из C и I могут быть выведены законы причинных связей, обозначаемые CL . Изменения касаются в основном сценария.

Состоянием теперь мы будем называть некоторый временной срез, в котором могут образоваться несколько ситуаций, задающих несколько возможных миров. Множество состояний мы будем обозначать $ST = \{s1, \dots, sn\}$. С состоянием ассоциировано некоторое восприятие $\mathbf{B}i = \{Fi, \dots\}$ в $s1$, которое сообщает нам фактическую информацию (возможно не полную) о среде. Расширенным восприятием мы назовем результат подстановки данных восприятия в ограничения состояний и последующего вывода, и обозначим его $\mathbf{B}'i$. Заметим, что теперь мы будем ассоциировать восприятие именно с состоянием, которое четко задает отрезок времени. Также состоянию сопоставлено множество гипотетических ситуаций, представляющих собой все возможные полные непротиворечивые и приемлемые с точки зрения ограничений ситуации, а после изначального состояния — еще и согласующиеся с законами действий и причинных связей. Гипотетические ситуации мы будем обозначать $\mathbf{S}'i = \{Fi, \dots\}$ после $[\alpha]$, множество гипотетических ситуаций - $\mathbf{S}' = \{\mathbf{S}'1, \dots, \mathbf{S}'n\}$. Фактические ситуации — это те ситуации, которые прошли проверку на согласование с восприятием. Они, как и ранее, обозначаются $\mathbf{S}i = \{Fi, \dots\}$ после $[\alpha]$, их множество - $\mathbf{S} = \{\mathbf{S}1, \dots, \mathbf{S}'n\}$. Ситуации дополнены указанием на то, какая последовательность действий к ним привела.

Каждая ситуация задает возможный мир, в котором развивается посредством применения действий. Возможные миры развиваются параллельно, но могут случаться ситуации, когда они пересекаются — когда действия в различных ситуациях приводят к той-же ситуации-наследнику. Новое расщепление миров невозможно, т. к. все согласующиеся с массивом данных восприятий и ранее произведенных действий ситуации уже составлены путем перебора, и новых вариаций быть не может.

Отбрасывая возможные миры по результатам проверки, мы тем самым

снижаем неопределенность не только в настоящем, но и в прошлом. Это связано с тем, что возможные миры в нашей модели могут ветвиться назад во времени, но не могут ветвиться вперед.

Основываясь на этой модели попытаемся предложить способ, позволяющий совершать умозаключения по умолчанию.

3.2. Состояния по умолчанию: нормы предположений.

В этом заключительном разделе мы рассмотрим собственно то, что составляет решение проблемы квалификации, и рассмотрим некоторые возможные дальнейшие расширения построенной модели.

Проблема квалификации — это проблема нахождения такого формализма, который бы позволил указать причины возможных неудач действия, но позволил бы действовать, не располагая всей полнотой информации относительно этих условий. Предложенная нами модель возможных миров позволяет это сделать введением модальных операторов.

Информация о том, какие условия необходимы для реализации действия, задается нами в самом законе действия. Нам не требуется дополнительное множество ограничений. Кроме того, схожесть по форме этих условий с условиями исполнения законов причинных связей позволяет нам рассуждать об успешности неявно заданных действий — таких, где в качестве результата указан не прямой эффект (действие «включить свет», а не «поменять положение выключателя»). Все что нам необходимо сделать, это задать условия умолчания, то есть некоторых привилегированных ситуаций, которые мы будем считать наиболее вероятными.

Построенная нами модель с возможными мирами позволяет ввести модальности, близкие к алетическим: мы можем говорить о некоторых необходимых в данной ситуации положениях дел, а можем — о возможных. Можем говорить, что такое-то действие с необходимостью повлечет такой-то результат (если во всех возможных мирах соблюдены все условия). Но наша

задача на данный момент заключается в другом: отложить построение полной картины всех возможностей, и выбрать какую-то одну (однако не забывая о том, что она лишь возможна, а не достоверно необходима). Достичь этого можно посредством введения модальностей, близких к деонтическим — которые скажут нам как *должно* себя вести.

Итак нам нужно выбрать какой-либо привилегированный возможный мир. Мы сможем сделать это на основе оценок. Оценивая некоторые литералы как наиболее вероятные, как те, что «обычно» таковы, мы сможем задать норму, обязывающую нас выбрать такой возможный мир, где они выполняются. Однако сохраняя информацию о том, что это было только предположением, мы сможем из опровержения этого предположения сделать вывод о его ошибочности.

В упомянутом выше сценарии, в котором наша задача — завести двигатель, мы можем сделать таким наиболее вероятным литерал, указывающий на отсутствие в выхлопной трубе посторонних предметов. Тогда мы *должны* будем предположить, что результатом поворота ключа в замке зажигания станет работающий двигатель.

Рассмотрим более комплексную формулировку этого домена. Индивид будет только один — {key}, обозначающий ключ зажигания. Функции — {Clearpipe, Engineruns, Bataryloaded, Onposition(x)}, означающие, соответственно, что труба чиста, двигатель заведен, аккумулятор заряжен, x находится в верном положении (имеется в виду то положение ключа зажигания, которое соответствует работе двигателя). Действие {_turn(x)} с законом

(0) _turn(x) меняет \sim Onposition(x) на Onposition(x) при условии \top
будет обозначать поворот ключа зажигания.

Ограничения состояний будут включать формулу

Engineruns \equiv Clearpipe & Bataryloaded & Onposition(key),

которая в соединении с информацией о влиянии (Onposition(key), Engineruns) дает следующие причинные связи:

(1) $\text{Onposition}(\text{key})$ вызывает смену $\sim\text{Engineruns}$ на Engineruns при условии $\{\text{Clearpipe}, \text{Bataryloaded}\}$;

(2) $\sim\text{Onposition}(\text{key})$ вызывает смену Engineruns на $\sim\text{Engineruns}$ при условии \top .

$\mathbf{B0} = \{\sim\text{Onposition}(\text{key}), \sim\text{Engineruns}\}$ в $s0$.

Далее, мы добавим информацию о том, что Clearpipe *обычно истинно*. Для информации об этом мы введем множество оценок \mathbf{P} . В данном случае $\mathbf{P} = \{\text{Clearpipe}\}$. Этим мы выражаем то, что очень редко бывает ситуация, когда выхлопная труба забита.

Так, в состоянии $s0$ мы имеем

$\mathbf{w1: S0} = \{\sim\text{Onposition}(\text{key}), \sim\text{Engineruns}, \text{Clearpipe}, \text{Bataryloaded}\}$.

$\mathbf{w2: S0} = \{\sim\text{Onposition}(\text{key}), \sim\text{Engineruns}, \text{Clearpipe}, \sim\text{Bataryloaded}\}$.

$\mathbf{w3: S0} = \{\sim\text{Onposition}(\text{key}), \sim\text{Engineruns}, \sim\text{Clearpipe}, \text{Bataryloaded}\}$.

$\mathbf{w4: S0} = \{\sim\text{Onposition}(\text{key}), \sim\text{Engineruns}, \sim\text{Clearpipe}, \sim\text{Bataryloaded}\}$.

Сейчас, вместо того, чтобы просматривать возможное дальнейшее развитие ситуации во всех мирах, мы составим некоторое предположение. То есть просто выберем один из возможных миров. Наличие \mathbf{P} помогает нам сделать выбор. Мы *должны* остановить его на таком мире, где присутствует $\{\text{Clearpipe}\}$. В общем виде нам *следует выбирать* такое S^i (верхний индекс здесь означает возможный мир), для которого выполняется

$$\forall S^j (|S^j \cap \mathbf{P}| \leq |S^i \cap \mathbf{P}|)$$

То есть мощность пересечения его с \mathbf{P} больше, чем у всех остальных, или равна им. Тем самым мы обеспечиваем максимальное совпадение этих двух множеств.

При этом нам не обязательно протраивать все возможные миры. Мы можем просто набирать ситуацию сначала из восприятий, затем из следствий, получаемый подстановкой результатов наблюдений в ограничения состояния,

затем из множества P , затем случайным образом (например — присвоить всем литералам значение «истина»). Каждый шаг делается только в случае, если ситуация еще не полна, то есть если значения есть еще не у всех литералов.

Однако нам нужно как-то выразить, что взятый возможный мир — только один из возможных. С этой целью мы вводим оператор Δ , который будет ставиться перед теми литералами, чье значение только предположительно. Так мы получаем например такую ситуацию:

$\{\sim\text{Onposition}(\text{key}), \sim\text{Engineruns}, \Delta\text{Clearpipe}, \Delta\text{Bataryloaded}\}$.

Далее мы зададим правила, по которым этот оператор будет проноситься сквозь законы действий и причинных связей. И они следуют из предшествовавшего рассмотрению выводов в прошлое:

A меняет $\sim B$ на ΔB при условии ΔC ;

ΔA вызывает смену $\sim B$ на ΔB при условии C ;

A вызывает смену $\Delta \sim B$ на ΔB при условии C ;

A вызывает смену $\sim B$ на ΔB при условии ΔC ;

ΔA вызывает смену $\Delta \sim B$ на ΔB при условии C ;

ΔA вызывает смену $\sim B$ на ΔB при условии ΔC ;

ΔA вызывает смену $\Delta \sim B$ на ΔB при условии ΔC ;

Коротко можно выразить это, сказав, что если хоть что-то в законе предположительно, то предположителен и результат.

Для того, чтобы сделать обратные выводы, нам потребуется еще одна вспомогательная конструкция. В предварительных ситуациях-наследниках, кроме информации об эффектах, мы будем записывать информацию о зависимостях (Dep). Мы будем указывать, какой предположительный вывод мы сделали, опираясь на какие посылки. Это будет еще одно множество, состоящее из записей вида: $V[A, C]$, что означает « V — предположительный вывод на основе предположительных посылок A и C ». Эта информация перейдет в гипотетическую ситуацию, и там должна сравниваться с расширенным восприятием. Для обратного вывода мы установим следующие *нормы*,

регулирующие выбор предположений:

- Если $(B[A_1, \dots, A_n] \in S'i) \ \& \ (B \in B'i)$ то следует заменить $\Delta A_1, \dots, \Delta A_n$ на A_1, \dots, A_n ;
- Если $(B[A] \in S'i) \ \& \ (\sim B \in B'i)$ то следует заменить ΔA на $\sim A$;
- Если $(B[A_1, \dots, A_n] \in S'i) \ \& \ (\sim B \in B'i)$ то:
 - для всех A_i если $\sim(A_i \in P)$ то следует заменить ΔA_i на $\Delta \sim A_i$;
 - если для всех A_i верно $(A_i \in P)$ то следует заменить ΔA_1 на $\Delta \sim A_1$.

То есть если мы имеем несколько предположений, и их результат оправдался — предположения верны. Если наш результат зависел от одного предположения, и он не оправдался — значит предположение неверно. Самая сложная ситуация - если несколько предположений, но результат не оправдался. В таком случае мы для всех предположений, не входящих в P , устанавливаем обратное предположение, что условия не были выполнены (заменяем предположения истины на предположения ложности). Если же таких предположений, которые не входят в P нет, то мы очень осторожно, по одному, начинаем предполагать их ложность.

Теперь мы готовы вычислить результат действия $_turn(key)$:

$S_0 = \{\sim Onposition(key), \sim Engineruns, \Delta Clearpipe, \Delta Bataryloaded\}$.

$(\{Onposition(key), \sim Engineruns, \Delta Clearpipe, \Delta Bataryloaded\}, \{Onposition(key)\})$ по (0);

$(\{Onposition(key), \Delta Engineruns, \Delta Clearpipe, \Delta Bataryloaded\}, \{Onposition(key), \Delta Engineruns\}, \{Engineruns[Clearpipe, Bataryloaded]\})$ по (1);

$S'1 = (\{Onposition(key), \Delta Engineruns, \Delta Clearpipe, \Delta Bataryloaded\}, \{Engineruns[Clearpipe, Bataryloaded]\})$ после $[_turn(key)]$

Предположим, $B1 = \{Onposition(key), \sim Engineruns\}$ в $s1$. Тогда, согласно последней норме, мы предполагаем $\Delta \sim Bataryloaded$.

$S1 = \{Onposition(key), \sim Engineruns, \Delta Clearpipe, \Delta \sim Bataryloaded\}$.

И далее, если у нас в распоряжении есть еще какие-либо действия для проверки, мы будем выполнять их.

Таковая общая идея нормативного подхода к рассуждениям в мире с неполной информацией: не имея достаточных оснований для точного решения, мы нормативно задаем выбор того или иного состояния. В данной работе идея приводится в неразработанном до конца виде, остается много вопросов и требующих более строгого анализа моментов. Однако можно наметить основные пути развития этой идеи.

Во-первых, можно предпринять попытку внедрения на указанной модели аспектов логики предпочтения³². Она может дать более широкую дисперсию вариантов значимости литералов друг относительно друга, а благодаря этому выбор предположений станет более рациональным, менее случайным.

Во-вторых, интересно было бы рассмотреть возможности включения в универсум рассуждения конкурентных действий и случайных изменений. С одной стороны это весьма просто — нужно просто разрешить ветвление не только в прошлое, но и в будущее. Тогда предположениями необходимо будет пользоваться не только в отношении того, что не было известно изначально, но и относительно того, что было известно, но могло поменять свое значение.

И в третьих, интересно было бы рассмотреть возможности выведения различных аспектов нашей модели, вроде оценок (а исходя из них — норм предпочтений), ограничений состояний и даже законов действия в автоматическом режиме. Это открыло бы перед теориями действий совершенно новые горизонты.

³² См.: Boutilier C. Toward a Logic for Qualitative Decision Theory Интернет источник. Режим доступа: http://www.cs.toronto.edu/kr/papers/craig_kr94.pdf (дата обращения 15.05.2014 г.)

ЗАКЛЮЧЕНИЕ

Проблема фреймов, центральная тема данной работы, вынужденно была объяснена только в рамках самой работы. Это связано с тем широким толкованием, которое она приобрела в научных кругах. Нам необходимо выделить тот фрагмент, который мы сделаем предметом своего исследования. Для наших целей мы ввели формализм, отдаленно напоминающий ситуационное исчисление, и выразили проблему в нем. Проблема в таком виде заключалась в следующем: для спецификации того, что в мире меняется, а что остается неизменным после того, как было произведено какое-либо действие, требуется очень большое количество аксиом. Это мы назвали проблемой ветвлений — проблемой устранения ветвления ситуаций, если быть более точными. То есть нас занимал вопрос о том, как определить однозначно и в соответствии с интуицией результат действия в сложном домене, где имеют место различного рода причинные связи между элементами. И при этом сохранить разумным количество требуемых аксиом. Решение мы нашли во введении причинных связей. Мы описали действия как влияющие на очень маленький фрагмент реальности, но могущие вызвать тем самым цепь причинных связей как своих непрямых действий. Это позволило нам получить формализм, хорошо справляющийся с ветвящимися доменами.

Следующим шагом стало включение в рассмотрение таких ситуаций, где нам не доступна полная информация о состоянии среды. Поэтому наши выводы стали лишь гипотетическими (не из-за неудачности формализма, а из-за отсутствия фактических данных). Научиться работать с такой неопределенностью — это и было нашей задачей в случае проблемы квалификации, второй проблемы, которую мы для себя определили.

Эта проблема была основной для нас. Именно в ней мы предложили свой подход, основывающийся на семантике возможных миров. Мы подробно рассмотрели следствия принятия этой семантики, и в конце обозначили общую идею рассмотрения предположений как принимаемых в согласии с некоторой

нормой. Эта идея основана на том, что человек, его обыденный разум, по нашему мнению, ожидает от мира некоторого «нормального поведения». Полная информация о мире нам не доступна, но постепенно мы начинаем привыкать к этому. И в этом заключается «рациональность» - мы вырабатываем себе некоторые нормы взаимодействия с миром, и перестаем тяготиться отсутствием полной информации о нем.

В данной работе рассмотрен только простейший вариант ситуаций с неполной информацией: когда нам доступна вся общая информация о мире (законы причинных связей, законы действий), не доступны только какие-то фактические данные, значения конкретных литералов. Также в данной работе мы не рассматривали вопросов, возникающих с разрешением конкурентных действий и случайных изменений в среде. Мы обозначили общий путь их описания, но не углублялись в проблему.

Рассмотрение этих вопросов может стать продолжением данного исследования, вместе с попытками применения логики предпочтения к выведению норм предположения, и попытками расширения данного подхода до ситуаций, где нам не известны общие данные о мире, его законы.

Научная новизна исследования заключается в попытке применить сначала семантику возможных миров, а затем построенных на этой основе модальных операторов, нормативных по своей природе, определяющих то, как мы должны рассуждать о действиях и их результатах в мире с неполной информацией.

Список литературы:

1. Вригт Г. Х. фон. Логико-философские исследования. М., 1986
2. Ивин А. А. Модальные теории Яна Лукасевича – М.: Институт философии РАН, 2001.
3. Ивин А. А. Парадоксы модальной логики Я. Лукасевича // философские науки. 1980. № 1. С. 75-83.
4. Карпенко А.С. Логика Лукасевича и простые числа – М.: Наука, 2000.
5. Карпенко А. С. Фатализм и случайность будущего: логический анализ – М.: Наука, 1990.
6. Клини С. К. Введение в метаматематику – М.: Издательство иностранной литературы, 1957.
7. Лукасевич Я. Аристотелевская силлогистика с точки зрения современной формальной логики – М.: Издательство иностранной литературы, 1959.
8. Лукасевич Я. О детерминизме [Электронный ресурс]. URL: http://diakonia.narod.ru/lib/fil_i_log.htm (дата обращения: 14.05.2013).
9. Многозначные логики и их применения: Логические исчисления, алгебры и функциональные свойства. Под ред. Финна В. К. Том 1. М.: УРСС, 2008.
10. Многозначные логики и их применения: Логика в системах искусственного интеллекта. Под ред. Финна В. К. Том 2. М.: УРСС, 2008. 240 с.
11. Смирнов В. А. Логические системы с модальными временными операторами // Модальные и временные логики. М.: Институт философии АН СССР, 1979.
12. Смирнова Е. Д. Логическая семантика и философские основания логики – М. Издательство Московского университета, 1986.
13. Baker, A. Nonmonotonic reasoning in the framework of the situation calculus. Artificial Intelligence, 49, 5-24
14. Belnap N., Perloff M., Ming Xu. Facing the future. Agents and choices in our indeterminist world. Oxford University Press, 2001

15. Boutilier C. Toward a Logic for Qualitative Decision Theory Интернет источник. Режим доступа: http://www.cs.toronto.edu/kr/papers/craig_kr94.pdf (дата обращения 15.05.2014 г.)
16. Davis, E. Representations of commonsense knowledge. San Mateo, CA: Morgan Kaufmann. Davis, E. Representations of commonsense knowledge. San Mateo, CA: Morgan Kaufmann.
17. Elkan C. Reasoning about action in first-order logic. In *Proceedings of the Conference of the Canadian Society for Computational Studies of Intelligence (CSCSI)*, pp/ 221-227, Vancouver, Canada, May 1992. Morgan Kaufmann.
18. Fetzner J.H. The frame problem: Artificial Intelligence meets David Hume. In K.M. Ford and P.J. Hayes (Eds.), *Reasoning agents in a dynamic world: The frame problem*, 55-69. Greenwich, CT: JAI Press.
19. Fikes, R. and Nilsson, N. STRIPS: A new approach to the application of theorem proving to problem solving. *Artificial Intelligence*, 3, 189-208.
20. Finger J. J. Exploiting constraints in design synthesis. Department of Computer Science STAN-CS-88-1204, Stanford University.
21. Hanks S. and McDermott D. Nonmonotonic logic and temporal projection. *Artificial Intelligence*, 33(3):379-412.
22. Harel D., Kozen D., Tiuryn J. *Dynamic logic*. MIT Press, 2000 Burgess J. P.
23. Haugh, B. Simple causal minimizations for temporal persistence and projection. *Proceedings of AAAI 1987*, 218-223.
24. Kamermans M. and Schmits T. The History of the Frame Problem Интернет источник. Режим доступа: <http://pomax.nihongoresources.com/The%20History%20of%20the%20Frame%20Problem%20-%20final%20version.pdf> (дата обращения 15.05.2014 г.)
25. Kautz, H. The logic of persistence. *Proceedings of AAAI 1986*, 401-405.
26. Lifschitz, V. Formal theories of action. *Proceedings of IJCAI 1987*, 966-972.
27. Lifschitz V. Frames in the space of situations. *Artificial Intelligence*, 46:365-376, 1990

28. Lifschitz, V. Pointwise circumscription: Preliminary report. Proceedings of AAAI 1986, 406-410.
29. Lifschitz V. The Frame Problem, Then and Now Интернет источник. Режим доступа: <http://www.cs.utexas.edu/users/vl/papers/jmc.pdf> (дата обращения 15.05.2014 г.)
30. Lighthill, J: Artificial Intelligence: A General Survey // Artificial Intelligence: a paper symposium, Science Research Council
31. McArthur R. P. Factuality and modality in the future tense // Nous. 1974. Vol. 8. P. 283-288.
32. McCarthy, J. Applications of circumscription to formalizing common-sense knowledge. Artificial Intelligence, 28, 86-116. (1986).
33. McCarthy, J. Circumscription - a form of non-monotonic reasoning. Artificial Intelligence, 13, 27-39.
34. McCarthy J. Programs with common sense. / M. L. Minsky (Ed.), Semantic information processing. Cambridge, MA: MIT Press, 1968, 403-409.)
35. McCarthy J. and Hayes P. J. Some philosophical problems from the standpoint of Artificial Intelligence Интернет источник. Режим доступа: <http://www-formal.stanford.edu/jmc/mcchay69.pdf> (дата обращения 15.05.2014 г.)
36. McDermott D. A temporal logic for reasoning about processes and plans. Cognitive Science, 6, 101-155.
37. Morgenstern L. The Problem with Solutions to the Frame Problem Интернет источник. Режим доступа: <http://www-formal.stanford.edu/leora/fp.pdf> (дата обращения 15.05.2014 г.)
38. Morgenstern L. and Stein L. A. Why things go wrong: A formal theory of causal reasoning. Proceedings of AAAI 1988, 518-523.
39. Penberthy, J.S. Planning with continuous change. Technical Report # 93-12-01, University of Washington, Dept. of Computer Science and Engineering.
40. Pednault, E. Generalizing nonlinear plans to handle complex goals and actions with context-dependent effects. Proceedings of IJCAI 1991, 240-245.
41. Prior A. N. Time and modality. Oxford. 1957.

42. Prior A. N. Past, present and future. Oxford. 1967.
43. Reiter, R. The frame problem in the situation calculus: A simple solution (sometimes) and a completeness result for goal regression. In V. Lifschitz, (Ed.) Artificial Intelligence and mathematical theory of computation: Papers in honor of John McCarthy, 359-380. San Diego: Academic Press.
44. Rescher N. Truth and necessity in temporal perspective // The philosophy of time (ed. R. Gale). N. Y. 1967.
45. Rescher N. Many-valued logic. N. Y. 1969.
46. Sandewall E. Reasoning about actions and change with ramification. In *Computer Science Today*, volume 1000 of LNCS. Springer, 1995.
47. Schubert, L.K. Monotonic solution of the frame problem in the situation calculus: An efficient method for worlds with fully specified actions. In H.E. Kyburg, R.P. Loui, and G.N. Carlson (Eds.), Knowledge representation and defeasible reasoning, 23-67. Boston: Kluwer Academic Press.
48. Shoham, Y. Reasoning about change: Time and causation from the standpoint of artificial intelligence. Cambridge, MA: MIT Press, 1988
49. Thielscher M. Challenges for action theories. Springer, 2000
50. Shoham Y. and Val del A. Deriving properties of belief update from theories of action (II). In R. Bajcsy, editor, *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 732-737, Chambéry, France, August 1993. Morgan Kaufmann.
51. Winslett M. Reasoning about action using a possible models approach. In *Proceedings of the AAAI National Conference on Artificial Intelligence*, pp. 89-93, Saint Paul, MN, August 1988.