

УДК 004.8 (575.2)(04)

ОНТОЛОГИИ И АЛГЕБРАИЧЕСКИЕ СПЕЦИФИКАЦИИ

В.А. Лапшин

Приводятся различные аспекты рассмотрения онтологий как алгебраических спецификаций. Алгебраический подход для формирования спецификаций разработан достаточно хорошо, поэтому представляется полезным использовать накопленный опыт для построения онтологий.

Ключевые слова: онтологии; алгебраические спецификации; открытые языки представления спецификаций; CASL; institution.

Введение. Большинство популярных инструментов построения онтологий используют в качестве основного формализма моделирования те или иные логики исчисления предикатов [3]. Так, в языке OWL [3, 6], который на данный момент является “de facto” стандартным языком описания содержания онтологий, используется формализм логики предикатов первого порядка с элементами языка второго порядка (различными встроенными в язык ограничениями на вид предикатов онтологии). В данной работе описывается другой подход, на основе которого можно формировать онтологии – представление онтологий в виде алгебраической спецификации [9, 11].

Алгебраическая спецификация представляется в виде теории, высказывания которой формируются на основе множества имен типов и имен операций, называемый сигнатурой. Операции выражают в спецификации как отношения между объектами, так и средства построения новых типов из уже существующих. Смысл операций задается соотношениями между ними, выраженными в виде равенств.

Для образования выражений на основе имен операций сигнатуры обычно вводится множество переменных, которые помечены типами этой сигнатуры. На основе имен операций и заданного множества переменных можно сформировать множество синтаксически правильных выражений этой сигнатуры, или термов. Некоторые термы могут быть преобразованы друг в друга с помощью подстановок, правила которых задаются соотношениями – равенствами данной сигнатуры. Такие термы можно считать семантически эквивалентными. Если вместо переменных использовать константы (т.е. имена опера-

ций без параметров), то множество всех термов данной сигнатуры разбивается на непересекающиеся классы эквивалентности. Полученное множество классов эквивалентности является внутренним (машинным) представлением спецификации и может быть использовано для вычислений.

Каждый класс эквивалентности в описанной выше конструкции может содержать бесконечное множество термов. Оперировать такими множествами на машинном уровне невозможно, поэтому для каждого класса выбирается его представитель или дескриптор. Пусть, например, операции сигнатуры задают арифметические действия:

“сложение” $+:INT,INT \rightarrow INT$
и “вычитание” $-:INT,INT \rightarrow INT$.

Тогда константа $5: \rightarrow INT$ будет дескриптором класса, содержащего термы $+(2,3)$, $-(+ (2,7),4)$ и т.п. Каждый терм здесь представляет арифметическое выражение, для вычисления которого надо найти дескриптор класса, которому принадлежит данный терм.

Имена операций в сигнатуре синтаксически можно рассматривать как правила контекстно-свободной порождающей грамматики Хомского [4]. Например, операцию “сложения” $+:INT,INT \rightarrow INT$ можно задать как $INT \rightarrow +(INT,INT)$. Это правило контекстно-свободной грамматики можно записать в инфиксном виде. В этом случае получим правило $INT \rightarrow INT+INT$. Здесь символ “+” выступает в качестве терминала контекстно-свободной грамматики, а имя типа INT – как нетерминальный символ.

Абстрактные алгебраические типы. В алгебраическом подходе для представления он-

тологий используются так называемые много-
сортные алгебры [1, 8]. Для определения много-
сортной алгебры сначала задается ее сигнатура
как пара $\Sigma=(S,OP)$, состоящая из множества S
имен типов (сортов) и множества имен опера-
ций OP .

Обычно на множестве S не задается ника-
ких отношений, что, в общем, не отражает того,
что типы в онтологии в подавляющем большин-
стве случаев образуют иерархию. Поэтому часто
на множестве S определяется отношение час-
тичного порядка так, чтобы это множество об-
разовывало *решетку*. Решеткой называется ал-
гебраическая структура, заданная на частично-
упорядоченном множестве таким образом, что
выполняются следующие условия:

Для любых двух элементов $s_1, s_2 \in S$ имеет-
ся их непосредственный родительский элемент
 $s_p \in S$, представляющий собой наименьшую
верхнюю грань этих элементов.

Для любых двух элементов $s_1, s_2 \in S$ имеет-
ся их непосредственный дочерний элемент $s_c \in S$,
представляющий собой наибольшую нижнюю
грань этих элементов.

Обычно используют так называемые пол-
ные решетки, т.е. такие решетки, у которых
существуют наибольший и наименьший эле-
менты для любого множества элементов. В
онтологии наименьший тип (класс) обычно на-
зывают “ничто”, а наибольший носит название
“объект” (thing). Отношение иерархии форми-
рует дерево, если для каждого типа разрешается
иметь не более одного родителя (в этом слу-
чае наибольшая нижняя грань любых двух
элементов является наименьшим элементом), и
образует направленный ациклический граф,
если разрешено множественное наследование
типов.

Имена операций имеют вид $\sigma:s_1 \times s_2 \times \dots \times s_n \rightarrow s$, т.е. вместе с именем операции указываются типы ее аргументов и тип результата.

Σ -алгебра для сигнатуры $\Sigma=(S,OP)$ – это семейство A множеств A_s , заданных для каж-
дого сорта $s \in S$ и называемых носителями
алгебры A , а также семейство отображений
 $\alpha\sigma:A_{s_1} \times A_{s_2} \times \dots \times A_{s_n} \rightarrow A_s$, заданных для каж-
дого имени операции σ в сигнатуре $\Sigma=(S,OP)$.
Иначе говоря, на семействе множеств, пред-
ставляющих в Σ -алгебре типы, имена которых
заданы в множестве S , задаются отображения,
имена и типы которых заданы в множестве Σ .
Каждая такая Σ -алгебра представляет собой мо-
дель вычисления, заданную на именах типов и
операций, перечисленных в сигнатуре.

Пусть X обозначает некоторое множество
(множество переменных), элементы которого ти-
пизированы типами из множества S . Тогда термом
сигнатуры называется выражение вида $\sigma(T_1, T_1, \dots,$
 $T_n)$, где $\sigma \in \Sigma$ – это имя операции, а T_1, T_1, \dots, T_n
– это либо переменные соответствующих типов,
либо термы, сформированные из имен операций
данной сигнатуры тем же способом. Например,
пусть $+:INT \times INT \rightarrow INT$ – это операция из сигна-
туры $\Sigma=(S,OP)$ и $a, b, c \in X$ – переменные. Тогда
выражение $+(a,b)$ образует терм. Вместо параме-
тров можно поставить имя операции “+”, так как
тип ее результата совпадает с типами параметров.
Например, выражение $+(+(a,b),c)$ также образует
терм. Множество термов, образованных на осно-
ве операций данной Σ -алгебры с переменными X ,
обозначается как $T_\Sigma(X)$. В качестве переменных
можно использовать константы, образуемые из
имен операций с местностью ноль, т.е. операций
вида $\sigma: \rightarrow s$ с пустым списком параметров. Мно-
жество всевозможных термов, образованных из
констант и имен операций для данной сигнатуры,
обозначается T_Σ .

Множество термов $T_\Sigma(X)$ представляет со-
бой Σ -алгебру, т.е. имеет алгебраическую струк-
туру. Действительно, каждому имени типа $s \in S$
в сигнатуре $\Sigma=(S,OP)$ можно сопоставить
множество термов, чья сигнатура имени опе-
рации имеет тип s в качестве результата. На-
пример, типу INT из примера выше будут со-
ответствовать термы вида $+(a,b)$, $+(+(a,b),c)$,
 $+(+(a,b),+(c,b))$ и т.д. Семантика операций также
прозрачна: каждому имени операции $\sigma:s_1 \times s_2 \times \dots \times s_n \rightarrow s$
и каждой n -ке термов t_1, t_1, \dots, t_n типов $s_1,$
 s_1, \dots, s_n , соответственно, дается в качестве зна-
чения терм типа s вида $\sigma(t_1, t_1, \dots, t_n)$. Например,
для операции $*:INT \times INT \rightarrow INT$ термам $+(a,b)$
и $+(+(a,b),+(b,c))$ дается в соответствие терм
 $*(+(a,b),+(+(a,b),+(c,b)))$.

Для наложения ограничений на операции
 Σ -алгебры используются соотношения в виде
равенств. Точнее, равенство – это тройка вида
 (X, T_1, T_2) , где T_1 и T_2 – термы, образованные
из операций Σ -алгебры и переменных из X . На-
пример, $(\{a\}, +(a,b), +(b,a))$ выражает свойство
коммутативности операции $+$. Часто равенства
записывают просто в виде $T_1 = T_2$, а множество
переменных определяется из контекста. Напри-
мер, равенство $(\{a\}, +(a,b), +(b,a))$ можно запи-
сать просто $+(a,b) = +(b,a)$.

Множество термов T_Σ , где в качестве пе-
ременных используются константы (т.е. сим-
волы операций, список аргументов которых
пуст. Например, константа 5 определяется как

операция $5: \rightarrow \text{INT}$), образующие так называемую алгебру замкнутых термов. Множество T_Σ можно факторизовать по отношению эквивалентности, задаваемому соотношениями равенства. Полученное множество также образует алгебру и называется инициальной Σ -алгеброй. Классы эквивалентности этой алгебры – это бесконечные множества термов. На практике используется конечное приближение классов эквивалентности инициальной Σ -алгебры, в котором каждый класс содержит только конечное множество термов. Такое приближение называется *конечной аппроксимацией* инициальной Σ -алгебры.

Каждому классу эквивалентности инициальной алгебры обычно дается в соответствие его представитель, или *дескриптор*. Обычно дескриптор представляет собой самый короткий терм данного класса эквивалентности среди построенных или специальное выражение, введенное пользователем. Вычислить терм означает найти дескриптор класса эквивалентности, которому этот терм принадлежит. Вычисление можно производить переписыванием термов, если рассматривать соотношения эквивалентности как задание правил переписывания.

Вычисление термина для некоторых операций можно также реализовать специальным образом. Например, для операции $*: \text{INT} \times \text{INT} \rightarrow \text{INT}$ можно ввести семантику простым процессорным вычислением произведения двух чисел. Вычисление производится по древовидному представлению термина путем обхода в глубину слева направо. Таким образом, для термина $*(+(3,2),4)$ сначала будут вычислены термы $*(+(3,2)$ и 4 , а затем задействована внешняя процедура, которая вычислит результат сложения 5 и 4 , и этот результат будет возвращен в качестве дескриптора термина $*(+(3,2),4)$.

Алгебраический подход к спецификации достаточно хорошо изучен. Например, язык алгебраических спецификаций CASL [6] разработан специально созданной для этой цели группой CoFI [5]. Этот язык используется преимущественно для моделирования программных модулей, но известны работы, посвященные использованию языка CASL для моделирования онтологий [11].

Описание синтаксиса языка онтологии шаблонами. Термы представляющей онтологию конечной аппроксимации инициальной Σ -алгебры формируются на основе операций сигнатуры $\Sigma=(S,OP)$. Операции сигнатуры записываются в префиксном виде, но можно ввести инфикс-

ную запись операций. Например, терм операции $+: \text{INT} \times \text{INT} \rightarrow \text{INT}$ вида $+(3,5)$ можно записать в виде $3+5$. Каждая такая инфиксная запись представляет собой правило контекстно-свободной порождающей грамматики [4], в левой части которого стоит символ сорта результата операции. Инфиксную запись операции, в которой кроме параметров могут быть строки символов, будем называть *шаблоном*.

Шаблоны позволяют расширять синтаксис языка представления онтологии. Синтаксически шаблон представляет собой строку, в которой могут быть “дырки”, которые имеют типы. Также шаблон имеет тип результата. Таким образом, каждый шаблон – это операция в сигнатуре Σ -алгебры, представляющей данную онтологию. Добавление синтаксических шаблонов расширяет множество имен операций Σ -алгебры онтологии и позволяет расширить синтаксис языка представления онтологии.

Шаблоны рассматриваются как правила контекстно-свободной грамматики, в которых имена типов выступают как нетерминальные символы, а строки между “дырками” могут рассматриваться как терминалы. Например, инфиксная запись операции $+: \text{INT} \times \text{INT} \rightarrow \text{INT}$ – это правило $\text{INT} \rightarrow \text{INT} \text{ ' + ' INT}$, где INT – это нетерминальный символ, а ' + ' представляет собой терминал грамматики.

Для разрешения синтаксических неоднозначностей иногда приходится вводить в сигнатуру “синтаксические” типы, которые необходимы только для того, чтобы управлять процессом синтаксического анализа. Например, можно задать грамматику арифметических выражений следующим образом. Сигнатура $\Sigma=(S,OP)$ состоит из типов INT , PLUS и MUL с операциями $+: \text{MUL} \times \text{PLUS} \rightarrow \text{PLUS}$ и $*: \text{INT} \times \text{MUL} \rightarrow \text{MUL}$. Также сделаем тип PLUS подтипом типа INT , а тип MUL – подтипом типа PLUS . Каждое такое отношение тип-подтип генерирует правило вида тип \rightarrow подтип. Добавим также тип NUM для выделения констант и сделаем его подтипом типа MUL . Таким образом, для данной сигнатуры можно определить следующую грамматику шаблонов:

$\text{PLUS} \rightarrow \text{MUL} \text{ ' + ' PLUS}$

$\text{MUL} \rightarrow \text{NUM} \text{ ' * ' MUL}$

$\text{INT} \rightarrow \text{PLUS}$

$\text{PLUS} \rightarrow \text{MUL}$

$\text{MUL} \rightarrow \text{NUM}$

Предположим теперь, что в сигнатуре $\Sigma=(S,OP)$ заданы константы $1: \rightarrow \text{NUM}$, $2: \rightarrow \text{NUM}$

и $3: \rightarrow \text{NUM}$. Тогда выражение $1+2*3$ будет преобразовано в терм $+(1,*(2,3))$.

Заключение. В настоящей работе представлен метод представления онтологий посредством алгебраического подхода. Онтология рассматривается как алгебраическая спецификация, в которой определяются сигнатура и аксиомы алгебраической теории, представляющей данную онтологию. Для онтологии на основе ее алгебраической теории строится алгебра, получаемая как множество всех правильно построенных выражений в сигнатуре онтологии, профакторизованное по отношению эквивалентности, следующей из аксиом-равенств, введенных в онтологию. Эта алгебра может быть использована в качестве машинного представления онтологии, позволяющего производить формальные манипуляции над онтологией: вести диалог с онтологией, передавать машинное представление в последовательной форме, объединять онтологии друг с другом и т.д.

Данный подход реализован в системе ЭЗОП [2] (<http://www.ezop-project.ru>), представляющей собой Веб-сервер для коллективного построения онтологий. В системе реализовано расширение синтаксиса языка представления знаний с помощью описанного в статье механизма. Онтологий в системе представляются как алгебраические спецификации, на основе которых строится внутреннее представление, как это описано в данной работе. Система предоставляет сервис для осуществления диалога с описываемой онтологией, а также возможность импорта онтологий из других форматов представления онтологии. Дальнейшие перспективы разработки связаны с реализацией расширения языка представления онтологии на лексическом уровне.

Литература

1. *Бениаминов Е.М.* Алгебраические методы в теории баз данных и представлении знаний. М.: Научный мир, 2003.
2. *Бениаминов Е.М., Болдина Д.М.* Система представления и обработки знаний ЭЗОП // Материалы конференции Диалог'20001. Прикладные проблемы, 2001.
3. *Лапшин В.А.* Онтологии в компьютерных системах. М.: Научный мир, 2010.
4. *Лапшин В.А.* Лекции по математической лингвистике. М.: Научный мир, 2010.
5. Страница группы CoFI. <http://www.informatik.uni-bremen.de/cofi/wiki/index.php/CoFI>.
6. Страница описания языка OWL. <http://www.w3.org/TR/owl-features>.
7. *Bidoit M., Mosses P.* CASL User Manual. Introduction to Using the Common Algebraic Specification Language. Series: Lecture Notes in Computer Science, Vol. 2900, 2004.
8. *Goguen J., Malcolm G.* Algebraic semantics of imperative programs. MIT Press, 1996.
9. *Goguen J.* Data, Schema, Ontology, and Logic Integration // Proceedings CombLog'04 Workshop, Carnielli W., Dionisio M., Mateus P. Lisbon, Portugal, July 2004.
10. *Goguen J., Burstall R.* Institutions: abstract model theory for specification and programming // Journal of the ACM. 39. 1. 1992.
11. *Luttich K., Mossakowski T.* Specification of Ontologies in CASL // A. Varzi, L. Vieu eds., Formal Ontology in Information Systems – Proceedings of the Third International Conference (FOIS-2004), Vol. 114, pp. 140–150. IOS Press, Amsterdam, 2004.