

Абстрактный и конкретный синтаксис

Коновалов А.В.

29 марта 2023 г.

Абстрактный синтаксис

Абстрактный синтаксис (abstract syntax) — упрощённая грамматика языка, в которой отсутствует информация, гарантирующая построение уникальных деревьев вывода.

В отличие от конкретного синтаксиса, абстрактный синтаксис не предназначен для ни для вывода цепочек символов, принадлежащих языку, ни для синтаксического анализа этих цепочек.

Абстрактный синтаксис предназначен для описания *абстрактных синтаксических деревьев*.

В сущности, грамматику абстрактного синтаксиса можно считать разновидностью описания алгебраических типов данных.

Абстрактное синтаксическое дерево (1)

Абстрактное синтаксическое дерево или **дерево абстрактного синтаксиса** (abstract syntax tree) — структура данных, отражающая содержательную информацию о структуре программы. Узлами абстрактного синтаксического дерева служат конструкции языка программирования: внутренние узлы соответствуют операторам, их потомки — операндам.

В общем случае, любая конструкция языка программирования может быть обработана путём создания оператора для данной конструкции и рассмотрения семантически значащих компонентов конструкции в качестве операндов этого оператора.

Абстрактное синтаксическое дерево (2)

В отличие от конкретного синтаксического дерева (дерева вывода или дерева разбора) абстрактное синтаксическое дерево не включает вспомогательные языковые средства: скобки для группировки, разделители вроде запятых или точек с запятой и т.д.

Также в абстрактном синтаксическом дереве отсутствуют вспомогательные нетерминалы, используемые для указания приоритета.

Далее прилагательное *абстрактное* мы будем опускать. Конкретные синтаксические деревья мы будем называть деревьями вывода или деревьями разбора.

Примеры описаний абстрактного синтаксиса (1)

2. Language

In this section, we describe the language which will be used throughout this paper.

DEFINITION 2.1 (Language). The language for which the described transformations are to be performed is a simple higher-order functional language as shown in Figure 3.

$prog ::= e_0$	where	$f_1 = e_1 \dots f_n = e_n$	Program
e	$::=$	v	Variable
		$ c \ e_1 \dots e_n$	Constructor
		$ f$	Function Call
		$ \lambda v. e$	λ -Abstraction
		$ e_0 \ e_1$	Application
		$ \text{case } e_0 \text{ of } p_1 : e_1 \mid \dots \mid p_k : e_k$	Case Expression
p	$::=$	$c \ v_1 \dots v_n$	Pattern

Figure 3. Language Grammar

Скриншот из статьи *Hamilton, Geoff W. "Distillation: extracting the essence of programs." ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation (2007).*

https://doi.org/10.1007/978-3-642-11486-1_13

Примеры описаний абстрактного синтаксиса (2)

3.1 Syntax

Definition 3. Assume denumerable disjoint sets of symbols for constructors \mathcal{C} , functions \mathcal{F} , pattern functions \mathcal{G} (ranged over by c , f , and g , respectively), and variables \mathcal{X} (ranged over by x, y). Then the set of programs \mathcal{Q} , definitions \mathcal{D} , terms \mathcal{T} , and patterns \mathcal{P} (ranged over by q , d , t , and p , respectively) are defined by the abstract syntax grammar

(program)	$\mathcal{Q} \ni q ::= d^m$	(definitions)
(definition)	$\mathcal{D} \ni d ::= f\ x^n = t \mid (g\ p.\ x^n = t.)^m$	(function/matcher)
(pattern)	$\mathcal{P} \ni p ::= c\ x^n$	(flat pattern)
(term)	$\mathcal{T} \ni t ::= x \mid c\ t^n \mid f\ t^n \mid g\ t^m$	(variable/construction/application/match)
(value)	$\mathcal{V} \ni v ::= c\ v^n$	

where $n \geq 0$ and $m > 0$. We require that

Скриншот из статьи *Secher, J.P. (2001). Driving in the Jungle. In: Danvy, O., Filinski, A. (eds) Programs as Data Objects. PADO 2001. Lecture Notes in Computer Science, vol 2053. Springer, Berlin, Heidelberg. https://doi.org/10.1007/3-540-44978-7_12*

Примеры описаний абстрактного синтаксиса (3)

Как видно из приведённых примеров, описание абстрактного синтаксиса зачастую неформально и сочетает в себе черты математической нотации.

Абстрактный синтаксис часто подразумевает повторение некоторых дочерних конструкций несколько раз, что часто выражается через многоточие или (как в примере из статьи Secher'a) — через степень.

В абстрактном синтаксическом дереве (в его программном воплощении) повторению однородных дочерних узлов могут соответствовать массивы/списки соответствующего типа.

Абстрактный синтаксис и синтаксическое дерево

Абстрактный синтаксис определяет типы узлов синтаксических деревьев и допустимые типы потомков для каждого из узлов.

Пример. Грамматика абстрактного синтаксиса с двумя типами узлов для выражений — арифметических и логических:

```
ArithExpr → NUMBER  
           | ArithExpr BinOp ArithExpr  
           | - Expr  
           | BoolExpr ? ArithExpr : ArithExpr
```

```
BoolExpr → TRUE | FALSE  
          | ArithExpr RelOp ArithExpr  
          | BoolExpr BoolOp BoolExpr  
          | NOT BoolExpr
```

```
BinOp → + | - | * | /
```

```
RelOp → = | ≠ | > | < | ≤ | ≥
```

```
BoolOp → AND | OR
```


Представление синтаксического дерева в программе (1)

В объектно-ориентированных языках нетерминалам абстрактного синтаксиса соответствуют абстрактные базовые классы (или интерфейсы), продукциям — конкретные классы-наследники. Атрибутами конкретных классов являются семантически значимые компоненты продукций.

В примере выше нам потребуются абстрактные классы `ArithExpr` и `BoolExpr` для представления выражений, типы данных для `BinOp`, `RelOp` и `BoolOp` можно представить при помощи перечислений.

Для продукции для тернарного оператора `ArithExpr` → `BoolExpr` ? `ArithExpr` : `ArithExpr` будет написан класс, полями которого будет одно логическое и два арифметических выражения. Знаки ? и : суть разделители, используемые для удобства записи.

Представление синтаксического дерева в программе (2)

В статически типизированных функциональных языках (Haskell, Standard ML и др.) для описания абстрактного синтаксиса разумно использовать алгебраические типы данных — нетерминалы соответствуют типам, продукции — конструкторам типов.

В объектно-функциональном языке Scala синтаксическое дерево естественно выражается при помощи case-классов.

Переход от абстрактного синтаксиса к конкретному

1. Абстрактный синтаксис описывает синтаксическое дерево — его нужно преобразовать в грамматику, способную порождать цепочки символов данного языка, пусть даже и неоднозначную.
2. Написать вспомогательные нетерминалы + правила, описывающие итерацию при помощи рекурсии.
3. Написать вспомогательные нетерминалы + правила для выражения приоритета и ассоциативности.

Приоритет и ассоциативность (1)

Рассмотрим ситуацию, когда в выражении две операции op_1 и op_2 имеют общий операнд b :

$a \ op_1 \ b \ op_2 \ c$

Тогда, если в языке операции op_1 и op_2 имеют разный **приоритет**, «спорный» аргумент b достанется той из них, чей приоритет выше.

В случае операций с равным приоритетом решение принимается на основе соответствующей их уровню приоритета **ассоциативности**: b принадлежит op_1 в случае **левой** ассоциативности, op_2 — в случае **правой**.

Приоритет и ассоциативность (2)

Левой ассоциативностью обладает, например, операция вычитания:

$$10 - 5 - 2 = (10 - 5) - 2 = 3 \neq 10 - (5 - 2) = 7$$

Правой ассоциативностью обладает, например, операция присваивания в Си ($a = b = c$ эквивалентно $a = (b = c)$) или операция возведения в степень в Python:

```
>>> 2 ** 3 ** 2
512
```

Если операцию нельзя объединять в цепочку вида $a \text{ op } b \text{ op } c$, то говорят, что операция **неассоциативна**. Например, в Паскале неассоциативны операции отношения, т.е. выражение вида $a > b > c$ синтаксически некорректно.

Пример перехода к конкретному синтаксису (1)

Пусть мы имеем следующий абстрактный синтаксис:

$$\begin{aligned} Expr \rightarrow & Expr + Expr \mid Expr - Expr \mid \\ & \mid Expr * Expr \mid Expr / Expr \\ & \mid Expr^{\wedge} Expr \\ & \mid v \\ & \mid f(Expr_1, \dots, Expr_N) \end{aligned}$$

Здесь v и f — имена переменных и функций, идентификаторы, \wedge — операция возведения в степень.

В спецификации языка говорится, что вызов функции должен принимать как минимум один параметр (т.е. $N > 0$) и для указания порядка выполнения операций можно использовать круглые скобки.

Пример перехода к конкретному синтаксису (2)

Запишем (неоднозначную) грамматику, позволяющую выводить строки символов рассматриваемого языка. Повторение параметров пока напишем условно:

$$\begin{aligned} \text{Expr} \rightarrow & \text{Expr} + \text{Expr} \mid \text{Expr} - \text{Expr} \\ & \mid \text{Expr} * \text{Expr} \mid \text{Expr} / \text{Expr} \\ & \mid \text{Expr} ^ \text{Expr} \\ & \mid \text{IDENT} \\ & \mid \text{IDENT}(\text{Expr}, \dots, \text{Expr}) \\ & \mid (\text{Expr}) \end{aligned}$$

В грамматику мы добавили круглые скобки, которые можно использовать для указания порядка действий.

Пример перехода к конкретному синтаксису (3)

Один или более параметров функции изобразим при помощи вспомогательного нетерминала *Args* с рекурсивным правилом:

```
Expr → Expr + Expr | Expr - Expr  
      | Expr * Expr | Expr / Expr  
      | Expr ^ Expr  
      | IDENT  
      | IDENT(Args)  
      | (Expr)  
Args → Expr | Args, Expr
```


Пример перехода к конкретному синтаксису (4) — приоритет

Разберёмся с приоритетами операций, игнорируя пока ассоциативность.

В рассматриваемом языке наивысшим приоритетом обладает возведение в степень, затем следуют умножение и деление, затем сложение и вычитание.

Для простоты разберёмся сначала со сложением и вычитанием. Операндами $+$ и $-$ у нас могут быть любые выражения, операндами операций $*$, $/$ и $^$ выражения вида $a + b$ и $a - b$ быть не могут (если они, конечно, не заключены в скобки).

Разделим нетерминал Expr на два: Expr_1 , который может содержать любые операции, и Expr_2 , который операции $+$ и $-$ содержать не может.

Пример перехода к конкретному синтаксису (5) — приоритет

Любое выражение Expr_2 будет выражением Expr_1 . Это можно выразить правилом $\text{Expr}_1 \rightarrow \text{Expr}_2$.

Операндами операций $+$ и $-$ могут быть выражения, содержащие любые знаки операций, т.е. Expr_1 . Это можно выразить правилами

$$\text{Expr}_1 \rightarrow \text{Expr}_1 + \text{Expr}_1 \mid \text{Expr}_1 - \text{Expr}_1$$

Выражения Expr_1 операндами операций $*$, $/$ и $^$ быть не могут. Поэтому для нетерминала Expr_2 просто повторяем исходные productions для Expr , заменяя в них Expr на Expr_2 (кроме выражения в скобках — там должно остаться Expr_1).

Пример перехода к конкретному синтаксису (6) — приоритет

Получим следующую грамматику:

$$\text{Expr}_1 \rightarrow \text{Expr}_1 + \text{Expr}_1 \mid \text{Expr}_1 - \text{Expr}_1 \mid \text{Expr}_2$$
$$\text{Expr}_2 \rightarrow \text{Expr}_2 * \text{Expr}_2 \mid \text{Expr}_2 / \text{Expr}_2$$
$$\mid \text{Expr}_2 ^ \text{Expr}_2$$
$$\mid \text{IDENT}$$
$$\mid \text{IDENT}(\text{Args})$$
$$\mid (\text{Expr}_1)$$
$$\text{Args} \rightarrow \text{Expr}_1 \mid \text{Args}, \text{Expr}_1$$

В полученной грамматике выражения вида $a + b$ и $a - b$ не могут быть операндами операций $*$, $/$ и $^$, следовательно, в выражениях вида $a + b * c$ приоритет $*$ будет выше, чем $+$ (и аналогично для других операций).

Пример перехода к конкретному синтаксису (7) — приоритет

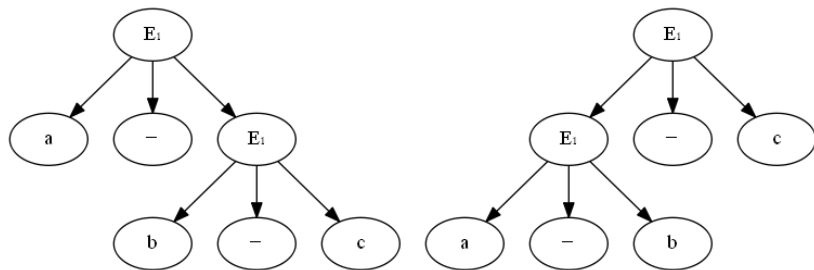
Повторяя те же рассуждения к операции $^$ с наивысшим приоритетом (её операндами не могут быть выражения вида $a * b$ или a / b), получим грамматику

$$\text{Expr}_1 \rightarrow \text{Expr}_1 + \text{Expr}_1 \mid \text{Expr}_1 - \text{Expr}_1 \mid \text{Expr}_2$$
$$\text{Expr}_2 \rightarrow \text{Expr}_2 * \text{Expr}_2 \mid \text{Expr}_2 / \text{Expr}_2 \mid \text{Expr}_3$$
$$\text{Expr}_3 \rightarrow \text{Expr}_3 ^ \text{Expr}_3$$
$$\mid \text{IDENT}$$
$$\mid \text{IDENT}(\text{Args})$$
$$\mid (\text{Expr}_1)$$
$$\text{Args} \rightarrow \text{Expr}_1 \mid \text{Args}, \text{Expr}_1$$

Таким образом, если операции в грамматике имеют различный приоритет, для каждого уровня приоритета нужно ввести вспомогательный нетерминал.

Пример перехода к конкретному синтаксису (8) — ассоциативность

Полученная грамматика по-прежнему неоднозначная. Для выражений вида IDENT - IDENT - IDENT можно построить два дерева разбора



Очевидно, правое дерево будет правильным, а левое — нет.

Пример перехода к конкретному синтаксису (9) — ассоциативность

Т.к. ассоциативность операции $-$ — левая, то правым операндом не может быть выражение вида $b - c$.

А это означает, что правым операндом операций $+$ и $-$ (у них одинаковый приоритет и ассоциативность) не может быть выражение, содержащее $+$ или $-$, т.е. правым операндом должно быть Expr_2 .

Пример перехода к конкретному синтаксису (10) — ассоциативность

Получаем грамматику с правильной ассоциативностью для операций + и -:

$$\text{Expr}_1 \rightarrow \text{Expr}_1 + \text{Expr}_2 \mid \text{Expr}_1 - \text{Expr}_2 \mid \text{Expr}_2$$
$$\text{Expr}_2 \rightarrow \text{Expr}_2 * \text{Expr}_2 \mid \text{Expr}_2 / \text{Expr}_2 \mid \text{Expr}_3$$
$$\text{Expr}_3 \rightarrow \text{Expr}_3 ^ \text{Expr}_3$$
$$\mid \text{IDENT}$$
$$\mid \text{IDENT}(\text{Args})$$
$$\mid (\text{Expr}_1)$$
$$\text{Args} \rightarrow \text{Expr}_1 \mid \text{Args}, \text{Expr}_1$$

Пример перехода к конкретному синтаксису (11) — ассоциативность

Аналогично исправляем ассоциативность для умножения и деления:

$$\text{Expr}_1 \rightarrow \text{Expr}_1 + \text{Expr}_2 \mid \text{Expr}_1 - \text{Expr}_2 \mid \text{Expr}_2$$
$$\text{Expr}_2 \rightarrow \text{Expr}_2 * \text{Expr}_3 \mid \text{Expr}_2 / \text{Expr}_3 \mid \text{Expr}_3$$
$$\text{Expr}_3 \rightarrow \text{Expr}_3 ^ \text{Expr}_3$$
$$\mid \text{IDENT}$$
$$\mid \text{IDENT}(\text{Args})$$
$$\mid (\text{Expr}_1)$$
$$\text{Args} \rightarrow \text{Expr}_1 \mid \text{Args}, \text{Expr}_1$$

Пример перехода к конкретному синтаксису (12) — ассоциативность

Ассоциативность операции возведения в степень $^$ — правая. Из тех же соображений получаем, что её левым операндом не может быть выражение вида $a \wedge b$. Для обеспечения этого вводим ещё один нетерминал Expr_4 :

$$\text{Expr}_1 \rightarrow \text{Expr}_1 + \text{Expr}_2 \mid \text{Expr}_1 - \text{Expr}_2 \mid \text{Expr}_2$$
$$\text{Expr}_2 \rightarrow \text{Expr}_2 * \text{Expr}_3 \mid \text{Expr}_2 / \text{Expr}_3 \mid \text{Expr}_3$$
$$\text{Expr}_3 \rightarrow \text{Expr}_4 \wedge \text{Expr}_3 \mid \text{Expr}_4$$
$$\text{Expr}_4 \rightarrow \text{IDENT} \mid \text{IDENT}(\text{Args}) \mid (\text{Expr}_1)$$
$$\text{Args} \rightarrow \text{Expr}_1 \mid \text{Args}, \text{Expr}_1$$

Пример перехода к конкретному синтаксису (13)

Давайте теперь дадим более осмысленные имена нетерминалам (term — слагаемое, factor — сомножитель, base — основание степени):

$$\text{Expr} \rightarrow \text{Expr} + \text{Term} \mid \text{Expr} - \text{Term} \mid \text{Term}$$
$$\text{Term} \rightarrow \text{Term} * \text{Factor} \mid \text{Term} / \text{Factor} \mid \text{Factor}$$
$$\text{Factor} \rightarrow \text{Base} ^ \text{Factor} \mid \text{Base}$$
$$\text{Base} \rightarrow \text{IDENT} \mid \text{IDENT}(\text{Args}) \mid (\text{Expr})$$
$$\text{Args} \rightarrow \text{Expr} \mid \text{Args}, \text{Expr}$$

Можем сделать **вывод**, что левая ассоциативность соответствует левой рекурсии в правилах, правая — правой.

Ещё заметим, что в этой грамматике мы можем считать, что запятая имеет приоритет ещё ниже, чем сложение.

Переход от абстрактного синтаксиса к конкретному

Выводы:

- ▶ Повторение выражается рекурсией.
- ▶ Для каждого уровня приоритета вводится вспомогательный нетерминал.
- ▶ Левая ассоциативность соответствует в правилах левой рекурсии, правая — правой.