

Лабораторная работа 1. Лексический анализатор. Освоение работы с существующими лексическими анализаторами. Разработка собственного лексического анализатора подмножества языка

п р о г р а м м и р о в а н и я, для чего определить лексические правила и выполнить перевод потока символов в поток лексем (токенов).

Лексический анализатор читает поток символов, составляющих исходную программу, и группирует эти символы в значащие последовательности, называемые лексемами.

На вход программы подается текстовый файл (напр. с именем INPUT.TXT), содержащий строки символов анализируемой программы. Например, строка может задавать переменной значения арифметического выражения в виде

ПЕРЕМЕННАЯ = ВЫРАЖЕНИЕ.

Выражение может включать:

Знаки сложения и умножения («+» и «*»);

Круглые скобки («(» и «)»);

Константы (например, 5; 3.8; 1e+18, 8.41E-10);

Имена переменных.

Имя переменной – это последовательность букв и цифр, начинающаяся с буквы.

Разбор выражения $COST = (PRICE + TAX) * 0.98$.

Проанализируем выражение:

$COST$, $PRICE$ и TAX –

лексемы-идентификаторы;

0.98 – лексема-константа; =, +, * – просто

лексемы.

Пусть все константы и идентификаторы можно отображать в лексемы типа <идентификатор> (<ИД>). Тогда выходом лексического анализатора будет последовательность лексем $\langle ИД_1 \rangle = (\langle ИД_2 \rangle + \langle ИД_3 \rangle) * \langle ИД_4 \rangle$.

Вторая часть компоненты лексемы (указатель, т.е. номер лексемы в таблице имен) – показана в виде индексов. Символы «=», «+» и «*» трактуются как лексемы, тип которых представляется ими самими. Они не имеют связанных с ними данных и, следовательно, не имеют указателей.

РАБОТА С ТАБЛИЦЕЙ ИМЕН

После того, как в результате лексического анализа лексемы распознаны, информация о некоторых из них собирается и записывается в таблицу имен.

Для нашего примера *COST*, *PRICE* и *TAX* – переменные с плавающей точкой. Рассмотрим вариант такой таблицы. В ней перечислены все идентификаторы вместе с относящейся к ним информацией (табл. 1).

Табл. 1 – Таблица имен

Номер элемента	Идентификатор	Информация
1	<i>COST</i>	Переменная с плавающей точкой
2	<i>PRICE</i>	Переменная с плавающей точкой
3	<i>TAX</i>	Переменная с плавающей точкой
4	0.98	Константа с плавающей точкой

Если позднее во входной цепочке попадается идентификатор, надо справиться в этой таблице, не появлялся ли он ранее. Если да, то лексема, соответствующая новому вхождению этого идентификатора, будет той же, что и у предыдущего вхождения.

Возможно создание таблиц констант, ключевых слов, разделителей, логических/математических операторов

При определении неверной последовательности символов, необходимо обнаружить эту ошибку и выдать сообщение о ней.

Лабораторная работа 2. Синтаксический анализатор.

В ходе синтаксического анализа исходный текст программы проверяется на соответствие синтаксическим нормам языка с построением дерева разбора (синтаксическое дерево), которое отражает синтаксическую структуру входной последовательности и удобно для дальнейшего использования, а также в случае несоответствия – позволяет вывести сообщения об ошибках.

Как правило, результатом синтаксического анализа является синтаксическое строение предложения, представленное либо в виде дерева зависимостей, либо в виде дерева составляющих, либо в виде некоторого сочетания первого и второго способов представления.

Таким образом на основе анализа выражений, состоящих из литералов, операторов и круглых скобок выполняется группирование токенов исходной программы в грамматические фразы, используемые для синтеза вывода.

Представление грамматических фраз исходной программы выполнить в виде дерева. Реализовать синтаксический анализатор с использованием одного из табличных методов (LL-, LR-метод, метод предшествования и пр.).

Разбор выражения $COST = (PRICE + TAX) * 0.98$.

Выходом анализатора служит дерево, которое представляет синтаксическую структуру, присущую исходной программе.

$$\langle ID_1 \rangle = (\langle ID_2 \rangle + \langle ID_3 \rangle) * \langle ID_4 \rangle.$$

По этой цепочке необходимо выполнить следующие действия:

- 1) $\langle ID_3 \rangle$ прибавить к $\langle ID_2 \rangle$;
- 2) результат (1) умножить на $\langle ID_4 \rangle$;
- 3) результат (2) поместить в ячейку, резервированную для $\langle ID_1 \rangle$. Этой последовательности соответствует дерево, изображенное на рис. 2.1.

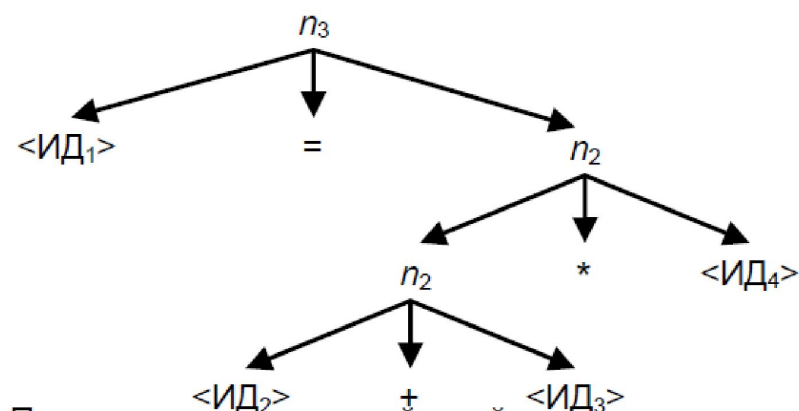


Рис. 2 – Последовательность действий при вычислении выражения

Т.е. мы имеем последовательность шагов в виде помеченного дерева.

Внутренние вершины представляют те действия, которые можно выполнять. Прямые потомки каждой вершины либо представляют аргументы, к которым нужно применять действие (если соответствующая вершина помечена идентификатором или является внутренней), либо помогают определить, каким должно быть это действие, в частности, знаки «+», «*» и «=». Скобки отсутствуют, т.к. они только определяют порядок действий.

LL- и LR- методы позволяют обнаружить ошибки на самых ранних стадиях, т.е. когда разбор потока токенов от лексического анализатора в соответствии с грамматикой языка становится невозможен.

Можно использовать нисходящий (англ. top-down parser) со стартового символа, до получения требуемой последовательности токенов. Для этих целей применим метод рекурсивного спуска либо LL-анализатор. Или использовать восходящий (англ. bottom-up parser) - продукции восстанавливаются из правых частей, начиная с токенов и кончая стартовым символом - LR-анализатор и проч.

Лабораторная работа 3. Семантический анализатор.

В процессе семантического анализа проверяется наличие семантических ошибок в исходной программе и накапливается информация о типах для следующей стадии – генерации кода. При семантическом анализе используются иерархические структуры, полученные во время синтаксического анализа для идентификации операторов и операндов выражений и инструкций.

Важным аспектом семантического анализа является проверка типов, когда компилятор проверяет, что каждый оператор имеет операнды допустимого спецификациями языка типа. Например, определение многих языков программирования требует, чтобы при использовании действительного числа в качестве индекса массива генерировалось сообщение об ошибке. В то же время спецификация языка может позволить определенное насильственное преобразование типов, например, когда бинарный арифметический оператор применяется к операндам целого и действительного типов. В этом случае компилятору может потребоваться преобразование целого числа в действительное.

В большинстве языков программирования имеет место неявное изменение типов (иногда называемое приведением типов (coercion)). Реже встречаются языки, подобные Ada, в которых большинство изменений типов должно быть явным.

В языках со статическими типами, например C, все типы известны во время компиляции, и это относится к типам выражений, идентификаторам и литералам. При этом неважно, насколько сложным является выражение: его тип может определяться во время компиляции за определенное количество шагов, исходя из типов его составляющих. Фактически, это позволяет производить контроль типов во время компиляции и находить заранее (в процессе компиляции, а не во время выполнения программы!) многие программные ошибки.

Лабораторная работа 4. Интерпретатор.

На основе результатов анализа лабораторных работ 1-3 выполнить интерпретацию программы.

Перечень языков программирования для вариантов работ:

Pascal (Extended Pascal, Object Pascal, PASCAL-XSC, Standard Pascal) - Delphi

C (C#/C++/Ruby C/JRuby/ANSI C).

Java/JavaScript

Python/Jython

Oracle SQL/ PL

F#

Lisp

Perl

Smalltalk

Prolog (диалекты: Edinburgh Prolog, ISO Prolog, Strawberry Prolog)