

3. 语义分析

本章实验为**实验二**，任务是在词法分析和语法分析程序的基础上编写一个程序，对C—源代码进行语义分析和类型检查，并打印分析结果。与实验一不同的是，实验二不再借助已有的工具，所有的任务都必须手写代码来完成。另外，虽然语义分析在整个编译器的实现中并不是难度最大的任务，但却是最细致、琐碎的任务。因此需要用心地设计诸如符号表、变量类型等数据结构的实现细节，从而正确、高效地实现语义分析的各种功能。

需要注意的是，由于在后面的实验中还会用到本次实验已经写好的代码，因此保持一个良好的代码风格、系统地设计代码结构和各模块之间的接口对于整个实验来讲相当重要。

3.1 实验内容

3.1.1 实验要求

在本次实验中，我们对C—语言做如下假设，你可以认为这些就是C—语言的特性（注意，假设3、4、5可能因后面的不同选做要求而有所改变）：

- 1) **假设1**：整型（int）变量不能与浮点型（float）变量相互赋值或者相互运算。
- 2) **假设2**：仅有int型变量才能进行逻辑运算或者作为if和while语句的条件；仅有int型和float型变量才能参与算术运算。
- 3) **假设3**：任何函数只进行一次定义，无法进行函数声明。
- 4) **假设4**：所有变量（包括函数的形参）的作用域都是全局的，即程序中所有变量均不能重名。
- 5) **假设5**：结构体间的类型等价机制采用**名等价（Name Equivalence）**的方式。
- 6) **假设6**：函数无法进行嵌套定义。
- 7) **假设7**：结构体中的域不与变量重名，并且不同结构体中的域互不重名。

以上假设1至7也可视为要求，违反即会导致各种语义错误，不过我们只对后面讨论的17种错误类型进行考察。此外，你可以安全地假设输入文件中不包含注释、八进制数、十六进制数、以及指数形式的浮点数，也不包含任何词法或语法错误（除了特别说明的针对选做要求的测试）。

你的程序需要对输入文件进行语义分析（输入文件中可能包含函数、结构体、一维和高维数组）并检查如下类型的错误：

- 1) **错误类型1**：变量在使用时未经定义。

- 2) **错误类型2**: 函数在调用时未经定义。
- 3) **错误类型3**: 变量出现重复定义, 或变量与前面定义过的结构体名字重复。
- 4) **错误类型4**: 函数出现重复定义 (即同样的函数名出现了不止一次定义)。
- 5) **错误类型5**: 赋值号两边的表达式类型不匹配。
- 6) **错误类型6**: 赋值号左边出现一个只有右值的表达式。
- 7) **错误类型7**: 操作数类型不匹配或操作数类型与操作符不匹配 (例如整型变量与数组变量相加减, 或数组 (或结构体) 变量与数组 (或结构体) 变量相加减)。
- 8) **错误类型8**: `return`语句的返回类型与函数定义的返回类型不匹配。
- 9) **错误类型9**: 函数调用时实参与形参的数目或类型不匹配。
- 10) **错误类型10**: 对非数组型变量使用 “[...]” (数组访问) 操作符。
- 11) **错误类型11**: 对普通变量使用 “(…)” 或 “()” (函数调用) 操作符。
- 12) **错误类型12**: 数组访问操作符 “[...]” 中出现非整数 (例如 `a[1.5]`)。
- 13) **错误类型13**: 对非结构体型变量使用 “.” 操作符。
- 14) **错误类型14**: 访问结构体中未定义过的域。
- 15) **错误类型15**: 结构体中域名重复定义 (指同一结构体中), 或在定义时对域进行初始化 (例如 `struct A { int a = 0; }`)。
- 16) **错误类型16**: 结构体的名字与前面定义过的结构体或变量的名字重复。
- 17) **错误类型17**: 直接使用未定义过的结构体来定义变量。

其中, 要注意三点: 一是关于数组类型的等价机制, 同C语言一样, 只要数组的基类型和维数相同我们即认为类型是匹配的, 例如 `int a[10][2]` 和 `int b[5][3]` 即属于同一类型; 二是我们允许类型等价的结构体变量之间的直接赋值 (见后面的测试样例), 这时的语义是, 对应的域相应赋值 (数组域也如此, 按相对地址赋值直至所有数组元素赋值完毕或目标数组域已经填满); 三是对于结构体类型等价的判定, 每个匿名的结构体类型我们认为均具有一个独有的隐藏名字, 以此进行名等价判定。

除此之外, 你的程序可以选择完成以下部分或全部的要求:

- 1) **要求2.1**: 修改前面的C—语言假设3, 使其变为 “函数除了在定义之外还可以进行声明”。函数的定义仍然不可以重复出现, 但函数的声明在相互一致的情况下可以重复出现。任一函数无论声明与否, 其定义必须在源文件中出现。在新的假设3下, 你的程序还需要检查两类新的错误和增加新的产生式:

- a) **错误类型18**: 函数进行了声明, 但没有被定义。
- b) **错误类型19**: 函数的多次声明互相冲突(即函数名一致, 但返回类型、形参数量或者形参类型不一致), 或者声明与定义之间互相冲突。
- c) 由于C—语言文法中并没有与函数声明相关的产生式, 因此你需要先对该文法进行适当修改。对于函数声明来说, 我们并不要求支持像“`int foo(int, float)`”这样省略参数名的函数声明。在修改的时候要留意, 你的改动应该以不影响其它错误类型的检查为原则。

2) **要求2.2**: 修改前面的C—语言假设4, 使其变为“变量的定义受可嵌套作用域的影响, 外层语句块中定义的变量可在内层语句块中重复定义(但此时在内层语句块中就无法访问到外层语句块的同名变量), 内层语句块中定义的变量到了外层语句块中就会消亡, 不同函数体内定义的局部变量可以相互重名”。在新的假设4下, 完成错误类型1至17的检查。

3) **要求2.3**: 修改前面的C—语言假设5, 将结构体间的类型等价机制由名等价改为**结构等价 (Structural Equivalence)**。例如, 虽然名称不同, 但两个结构体类型`struct a { int x; float y; }`和`struct b { int y; float z; }`仍然是等价的类型。注意, 在结构等价时不要将数组展开来判断, 例如`struct A { int a; struct { float f; int i; } b[10]; }`和`struct B { struct { int i; float f; } b[10]; int b; }`是不等价的。在新的假设5下, 完成错误类型1至17的检查。

3.1.2 输入格式

你的程序的输入是一个包含C—源代码的文本文件, 该源代码中可能会有语义错误。你的程序需要能够接收一个输入文件名作为参数。例如, 假设你的程序名为`cc`、输入文件名为`test1`、程序和输入文件都位于当前目录下, 那么在Linux命令行下运行`./cc test1`即可获得以`test1`作为输入文件的输出结果。

3.1.3 输出格式

实验二要求通过标准输出打印程序的运行结果。对于那些没有语义错误的输入文件, 你的程序不需要输出任何内容。对于那些存在语义错误的输入文件, 你的程序应当输出相应的错误信息, 这些信息包括错误类型、出错的行号以及说明文字, 其格式为:

```
Error type [错误类型] at Line [行号]: [说明文字].
```

说明文字的内容没有具体要求, 但是错误类型和出错的行号一定要正确, 因为这是判断输出的错误提示信息是否正确的唯一标准。请严格遵守实验要求中给定的错误分类, 否则将影响

你的实验评分。

输入文件中可能包含一个或者多个错误（但每行最多只有一个错误），你的程序需要将它们全部检查出来。当然，有些时候输入文件中的一个错误会产生连锁反应，导致别的地方出现多个错误（例如，一个未定义的变量在使用时由于无法确定其类型，会使所有包含该变量的表达式产生类型错误），我们只会去考察你的程序是否报告了较本质那个的错误（如果难以确定哪个错误更本质一些，建议你报告所有发现的错误）。但是，如果源程序里有错而你的程序没有报错或报告的错误类型不对，又或者源程序里没有错但你的程序却报错，都会影响你的实验评分。

3.1.4 测试环境

你的程序将在如下环境中被编译并运行（同实验一）：

- 1) GNU Linux Release: Ubuntu 12.04, kernel version 3.2.0-29;
- 2) GCC version 4.6.3;
- 3) GNU Flex version 2.5.35;
- 4) GNU Bison version 2.5。

一般而言，只要避免使用过于冷门的特性，使用其它版本的Linux或者GCC等，也基本上不会出现兼容性方面的问题。注意，实验二的检查过程中不会去安装或尝试引用各类方便编程的函数库（如glib等），因此请不要在你的程序中使用它们。

3.1.5 提交要求

实验二要求提交如下内容（同实验一）：

- 1) Flex、Bison以及C语言的可被正确编译运行的源程序。
- 2) 一份PDF格式的实验报告，内容包括：
 - a) 你的程序实现了哪些功能？简要说明如何实现这些功能。清晰的说明有助于助教对你的程序所实现的功能进行合理的测试。
 - b) 你的程序应该如何被编译？可以使用脚本、makefile或逐条输入命令进行编译，请详细说明应该如何编译你的程序。无法顺利编译将导致助教无法对你的程序所实现的功能进行任何测试，从而丢失相应的分数。
 - c) 实验报告的长度不得超过三页！所以实验报告中需要重点描述的是你的程序中的亮点，是你认为最个性化、最具独创性的内容，而相对简单的、任何人都可以做

的内容则可不提或简单地提一下，尤其要避免大段地向报告里贴代码。实验报告中所出现的最小字号不得小于五号字（或英文11号字）。

3.1.6 样例（必做内容）

实验二的样例包括**必做内容样例**与**选做要求样例**两部分，分别对应于实验要求中的必做内容和选做要求。请仔细阅读样例，以加深对实验要求以及输出格式要求的理解。这节列举必做内容样例。

样例1:

输入:

```
1 int main()
2 {
3     int i = 0;
4     j = i + 1;
5 }
```

输出:

样例输入中变量“j”未定义，因此你的程序可以输出如下的错误提示信息:

```
Error type 1 at Line 4: Undefined variable "j".
```

样例2:

输入:

```
1 int main()
2 {
3     int i = 0;
4     inc(i);
5 }
```

输出:

样例输入中函数“inc”未定义，因此你的程序可以输出如下的错误提示信息:

```
Error type 2 at Line 4: Undefined function "inc".
```

样例3:

输入:

```
1 int main()
2 {
3     int i, j;
4     int i;
5 }
```

输出:

样例输入中变量“i”被重复定义，因此你的程序可以输出如下的错误提示信息:

```
Error type 3 at Line 4: Redefined variable "i".
```

样例4:

输入:

```
1  int func(int i)
2  {
3      return i;
4  }
5
6  int func()
7  {
8      return 0;
9  }
10
11 int main()
12 {
13 }
```

输出:

样例输入中函数“func”被重复定义，因此你的程序可以输出如下的错误提示信息:

```
Error type 4 at Line 6: Redefined function "func".
```

样例5:

输入:

```
1  int main()
2  {
3      int i;
4      i = 3.7;
5  }
```

输出:

样例输入中错将一个浮点常数赋值给一个整型变量，因此你的程序可以输出如下的错误提示信息:

```
Error type 5 at Line 4: Type mismatched for assignment.
```

样例6:

输入:

```
1  int main()
2  {
3      int i;
4      10 = i;
5  }
```

输出:

样例输入中整数“10”出现在了赋值号的左边，因此你的程序可以输出如下的错误提示信息:

```
Error type 6 at Line 4: The left-hand side of an assignment must be a variable.
```

样例7:

输入:

```
1 int main()
2 {
3     float j;
4     10 + j;
5 }
```

输出:

样例输入中表达式“10 + j”的两个操作数的类型不匹配，因此你的程序可以输出如下的错误提示信息:

```
Error type 7 at Line 4: Type mismatched for operands.
```

样例8:

输入:

```
1 int main()
2 {
3     float j = 1.7;
4     return j;
5 }
```

输出:

样例输入中“main”函数返回值的类型不正确，因此你的程序可以输出如下的错误提示信息:

```
Error type 8 at Line 4: Type mismatched for return.
```

样例9:

输入:

```
1 int func(int i)
2 {
3     return i;
4 }
5
6 int main()
7 {
8     func(1, 2);
9 }
```

输出:

样例输入中调用函数“func”时实参数目不正确，因此你的程序可以输出如下的错误提示信息:

```
Error type 9 at Line 8: Function "func(int)" is not applicable for arguments
"(int, int)".
```

样例10:

输入:

```
1 int main()
2 {
3     int i;
4     i[0];
5 }
```

输出:

样例输入中变量“i”非数组型变量，因此你的程序可以输出如下的错误提示信息:

```
Error type 10 at Line 4: "i" is not an array.
```

样例11:

输入:

```
1 int main()
2 {
3     int i;
4     i(10);
5 }
```

输出:

样例输入中变量“i”不是函数，因此你的程序可以输出如下的错误提示信息:

```
Error type 11 at Line 4: "i" is not a function.
```

样例12:

输入:

```
1 int main()
2 {
3     int i[10];
4     i[1.5] = 10;
5 }
```

输出:

样例输入中数组访问符中出现了非整型常数“1.5”，因此你的程序可以输出如下的错误提示信息:

```
Error type 12 at Line 4: "1.5" is not an integer.
```

样例13:

输入:

```
1 struct Position
2 {
3     float x, y;
4 };
5
6 int main()
7 {
8     int i;
9     i.x;
10 }
```


输出:

样例输入中变量“i”非结构体类型变量,因此你的程序可以输出如下的错误提示信息:

```
Error type 13 at Line 9: Illegal use of ".".
```

样例14:

输入:

```
1 struct Position
2 {
3     float x, y;
4 };
5
6 int main()
7 {
8     struct Position p;
9     if (p.n == 3.7)
10         return 0;
11 }
```

输出:

样例输入中结构体变量“p”访问了未定义的域“n”,因此你的程序可以输出如下的错误提示信息:

```
Error type 14 at Line 9: Non-existent field "n".
```

样例15:

输入:

```
1 struct Position
2 {
3     float x, y;
4     int x;
5 };
6
7 int main()
8 {
9 }
```

输出:

样例输入中结构体的域“x”被重复定义,因此你的程序可以输出如下的错误信息:

```
Error type 15 at Line 4: Redefined field "x".
```

样例16:

输入:

```
1 struct Position
2 {
3     float x;
4 };
5
6 struct Position
7 {
8     int y;
9 };
```

```
10
11 int main()
12 {
13 }
```

输出:

样例输入中两个结构体的名字重复, 因此你的程序可以输出如下的错误信息:

```
Error type 16 at Line 6: Duplicated name "Position".
```

样例17:

输入:

```
1 int main()
2 {
3     struct Position pos;
4 }
```

输出:

样例输入中结构体“Position”未经定义, 因此你的程序可以输出如下的错误信息:

```
Error type 17 at Line 3: Undefined structure "Position".
```

3.1.7 样例 (选做要求)

这节列举选做要求样例。

样例1:

输入:

```
1 int func(int a);
2
3 int func(int a)
4 {
5     return 1;
6 }
7
8 int main()
9 {
10 }
```

输出:

如果你的程序需要完成要求2.1, 这个样例输入不存在任何词法、语法或语义错误, 因此不需要输出。

如果你的程序不需要完成要求2.1, 这个样例输入存在语法错误, 因此你的程序可以输出如下的错误提示信息:

```
Error type B at Line 1: Incomplete definition of function "func".
```

样例2:

输入:

```
1 struct Position
2 {
3     float x,y;
4 };
5
6 int func(int a);
7
8 int func(struct Position p);
9
10 int main()
11 {
12 }
```

输出:

如果你的程序需要完成要求2.1, 这个样例输入存在两处语义错误: 一是函数“func”的两次声明不一致; 二是函数“func”未定义, 因此你的程序可以输出如下的错误提示信息:

```
Error type 19 at Line 8: Inconsistent declaration of function "func".
Error type 18 at Line 6: Undefined function "func".
```

注意, 我们对错误提示信息的顺序不做要求。

如果你的程序不需要完成要求2.1, 这个样例输入存在两处语法错误, 因此你的程序可以输出如下的错误提示信息:

```
Error type B at Line 6: Incomplete definition of function "func".
Error type B at Line 8: Incomplete definition of function "func".
```

样例3:

输入:

```
1 int func()
2 {
3     int i = 10;
4     return i;
5 }
6
7 int main()
8 {
9     int i;
10    i = func();
11 }
```

输出:

如果你的程序需要完成要求2.2, 这个样例输入不存在任何词法、语法或语义错误, 因此不需要输出。

如果你的程序不需要完成要求2.2, 样例输入中的变量“i”被重复定义, 因此你的程序可以输出如下的错误提示信息:

```
Error type 3 at Line 9: Redefined variable "i".
```

样例4:

输入:

```
1 int func()
2 {
3     int i = 10;
4     return i;
5 }
6
7 int main()
8 {
9     int i;
10    int i, j;
11    i = func();
12 }
```

输出:

如果你的程序需要完成要求2.2, 样例输入中的变量 “i” 被重复定义, 因此你的程序可以输出如下的错误提示信息:

```
Error type 3 at Line 10: Redefined variable "i".
```

如果你的程序不需要完成要求2.2, 样例输入中的变量 “i” 被重复定义了两次, 因此你的程序可以输出如下的错误提示信息:

```
Error type 3 at Line 9: Redefined variable "i".
Error type 3 at Line 10: Redefined variable "i".
```

样例5:

输入:

```
1 struct Temp1
2 {
3     int i;
4     float j;
5 };
6
7 struct Temp2
8 {
9     int x;
10    float y;
11 };
12
13 int main()
14 {
15     struct Temp1 t1;
16     struct Temp2 t2;
17     t1 = t2;
18 }
```

输出:

如果你的程序需要完成要求2.3, 这个样例输入不存在任何词法、语法或语义错误, 因此不需要输出。

如果你的程序不需要完成要求2.3, 样例输入中的语句 “t1 = t2;” 其赋值号两边变量的类型不匹配, 因此你的程序可以输出如下的错误提示信息:

```
Error type 5 at Line 17: Type mismatched for assignment.
```

样例6:

输入:

```
1 struct Temp1
2 {
3     int i;
4     float j;
5 };
6
7 struct Temp2
8 {
9     int x;
10 };
11
12 int main()
13 {
14     struct Temp1 t1;
15     struct Temp2 t2;
16     t1 = t2;
17 }
```

输出:

如果你的程序需要完成要求2.3, 样例输入中的语句 “t1 = t2;” 其赋值号两边变量的类型不匹配, 因此你的程序可以输出如下的错误提示信息:

```
Error type 5 at Line 16: Type mismatched for assignment.
```

如果你的程序不需要完成要求2.3, 应该输出与上述一样的错误提示信息:

```
Error type 5 at Line 16: Type mismatched for assignment.
```

3.2 实验指导

除了词法和语法分析之外，编译器前端所要进行的另一项工作就是对输入程序进行语义分析。进行语义分析的原因很简单：一段语法上正确的源代码仍可能包含严重的逻辑错误，这些逻辑错误可能会对编译器后面阶段的工作产生影响。首先，我们在语法分析阶段所借助的理论工具是上下文无关文法，从名字上就可以看出上下文无关文法没有办法处理一些与输入程序上下文相关的内容（例如变量在使用之前是否已经被定义过，一个函数内部定义的变量在另一个函数中是否允许使用等）。这些与上下文相关的内容都会在语义分析阶段得到处理，因此也有人将这一阶段叫做**上下文相关分析（Context-sensitive Analysis）**。其次，现代程序设计语言一般都会引入类型系统，很多语言甚至是强类型的。引入类型系统可以为程序设计语言带来很多好处，例如它可以提高代码在运行时刻的安全性，增强语言的表达力，还可以使编译器为其生成更高效的目标代码。对于一个具有类型系统的语言来说，编译器必须要有能力检查输入程序中的各种行为是否都是类型安全的，因为类型不安全的代码出现逻辑错误的可能性很高。最后，为了使之后的阶段能够顺利进行，编译器在面对一段输入程序时不得不从语法之外的角度进行理解。比如，假设输入程序中有一个变量或函数 x ，那么编译器必须要提前确定：

- 1) 如果 x 是一个变量，那么变量 x 中存储的是什么内容？是一个整数值、浮点数值，还是一组整数值或其它自定义结构的值？
- 2) 如果 x 是一个变量，那么变量 x 在内存中需要占用多少字节的空间？
- 3) 如果 x 是一个变量，那么变量 x 的值在程序的运行过程中会保留多长时间？什么时候应当创建 x ，而什么时候它又应该消亡？
- 4) 如果 x 是一个变量，那么谁该负责为 x 分配存储空间？是用户显式地进行空间分配，还是由编译器生成专门的代码来隐式地完成这件事？
- 5) 如果 x 是一个函数，那么这个函数要返回什么类型的值？它需要接受多少个参数，这些参数又都是什么类型？

以上这些与变量或函数 x 有关的信息中，几乎所有都无法在词法或语法分析过程中获得，即输入程序能为编译器提供的信息要远超过词法和语法分析能从中挖掘出的信息。

从编程实现的角度看，语义分析可以作为编译器里单独的一个模块，也可以并入前面的语法分析模块或者并入后面的中间代码生成模块。不过，由于其牵扯到的内容较多而且较为繁杂，我们还是将语义分析单独作为一块内容。我们下面先对语义分析所要用到的属性文法做简要介绍，然后对C语言编译中的符号表和类型表示这两大重点内容进行讨论，最后提出帮助

顺利完成实验二的一些建议。

3.2.1 属性文法

在词法分析过程中，我们借助了正则文法；在语法分析过程中，我们借助了上下文无关文法；现在到了语义分析部分，为什么我们不能在文法体系中更上一层楼，采用比上下文无关文法表达力更强的上下文相关文法呢？

之所以不继续采用更强的文法，原因有两个：其一，识别一个输入是否符合某一上下文相关文法，这个问题本身是P-Space Complete¹的，也就是说，如果使用上下文相关文法那么编译器的复杂度会很高；其二，编译器需要获取的很多信息很难使用上下文相关文法进行编码，这就迫使我们为语义分析寻找其它更实用的理论工具。

目前被广泛使用的用于语义分析的理论工具叫做**属性文法 (Attribute Grammar)**，它是由Knuth在50年代所提出。属性文法的核心思想是，为上下文无关文法中的每一个终结符或非终结符赋予一个或多个属性值。对于产生式 $A \rightarrow X_1 \dots X_n$ 来说，在自底向上分析中 $X_1 \dots X_n$ 的属性值是已知的，这样语义动作只会为 A 计算属性值；而在自顶向下分析中， A 的属性值是已知的，在该产生式被应用之后才能知道 $X_1 \dots X_n$ 的属性值。终结符号的属性值通过词法分析可以得到，非终结符号的属性值通过产生式对应的语义动作来计算。

属性值可以分成不相交的两类：**综合属性 (Synthesized Attribute)**和**继承属性 (Inherited Attribute)**。在语法树中，一个结点的综合属性值是从其子结点的属性值计算而来的，而一个结点的继承属性值则是由该结点的父结点和兄弟结点的属性值计算而来的。如果对一个文法 P ， $\forall A \rightarrow X_1 \dots X_n \in P$ 都有与之相关联的若干个属性定义规则，则称 P 为**属性文法**。如果属性文法 P 只包含综合属性而没有继承属性，则称 P 为**S属性文法**。如果每个属性定义规则中的每个属性要么是一个综合属性，要么是 X_j 的一个继承属性，并且该继承属性只依赖于 $X_1 \dots X_{j-1}$ 的属性和 A 的继承属性，则称 P 为**L属性文法**。

以属性文法为基础可衍生出一种非常强大的翻译模式，我们称之为**语法制导翻译 (Syntax-Directed Translation或SDT)**。在SDT中，人们把属性文法中的属性定义规则用计算属性值的语义动作来表示，并用花括号“{”和“}”括起来，它们可被插入到产生式右部的任何合适的位置上，这是一种语法分析和语义动作交错的表示法。事实上，我们在之前使用Bison时已经用到了属性文法和SDT。

¹ <https://www.princeton.edu/~achaney/tmve/wiki100k/docs/PSPACE-complete.html>。

3.2.2 符号表

符号表对于编译器至关重要。在编译过程中，编译器使用符号表来记录源程序中各种名字的特性信息。所谓“名字”包括：程序名、过程名、函数名、用户定义类型名、变量名、常量名、枚举值名、标号名等，所谓“特性信息”包括：上述名字的种类、具体类型、维数、参数个数、数值及目标地址（存储单元地址）等。

符号表上的操作包括**填表**和**查表**两种。当分析到程序中的说明或定义语句时，应将说明或定义的名字，以及与之有关的特性信息填入符号表中，这便是填表操作。查表操作则使用得更广泛，需要使用查表操作的情况有：填表前查表，包括检查在输入程序的同一作用域内名字是否被重复定义，检查名字的种类是否与说明一致，对于那些类型要求更强的语言，则要检查表达式中各变量的类型是否一致等；此外生成目标指令时，也需要查表以取得所需要的地址或者寄存器编号等。符号表的组织方式也有多种，可以将程序中出现的所有符号组织成一张表，也可以将不同种类的符号组织成不同的表（例如，所有变量名组织成一张表，所有函数名组织成一张表，所有临时变量组织成一张表，所有结构体定义组织成一张表，等等）。你可以针对每个语句块、每个结构体都新建一张表，也可以将所有语句块中出现的符号全部插入到同一张表中。符号表可以仅支持插入操作而不支持删除操作（此时如果要实现作用域则需要将符号表组织成层次结构），也可以组织一张既可以插入又可以删除的、支持动态更新的表。不同的组织方式各有利弊，你可仔细思考并为实验二做出决定。

至于在符号表里应该填些什么，这与不同程序设计语言的特性相关，更取决于编译器的设计者本身。只要觉得方便，可以向符号表里填任何内容！毕竟符号表就是为了支持编写编译器而设置的。就实验二而言，对于变量至少要记录变量名及其类型，对于函数至少要记录其返回类型、参数个数以及参数类型。

至于符号表应该采用何种数据结构实现，这个问题同样没有统一的答案。不同的数据结构有不同的时间复杂度、空间复杂度以及编程难度，我们下面讨论几种最常见的选择。

线性链表：

符号表里所有的符号（假设有 n 个，下同）都用一条链表串起来，插入一个新的符号只需将该符号放在链表的表头，其时间复杂度是 $O(1)$ 。在链表中查找一个符号需要对其进行遍历，时间复杂度是 $O(n)$ 。删除一个符号只需要将该符号从链表里摘下来，不过在摘之前由于我们必须执行一次查找操作以找到待删除的结点，因此时间复杂度也是 $O(n)$ 。

链表的最大问题是它的查找和删除效率太低，一旦符号表中的符号数量较大，查表操作将

变得十分耗时。不过，使用链表的好处也是显而易见：它的结构简单，编程容易，可以被快速实现。如果你事先能够确定表中的符号数目较少（例如，在结构体定义中或在面向对象语言的一些短方法中），链表是一个非常不错的选择。

平衡二叉树：

相对于只能执行线性查找的链表而言，在平衡二叉树上进行查找天生就是二分查找。在一个典型的平衡二叉树实现（例如AVL树、红黑树或伸展树¹等）上查找一个符号的时间复杂度是 $O(\log n)$ 。插入一个符号相当于进行一次失败的查找而找到待插入的位置，时间复杂度也是 $O(\log n)$ 。删除一个符号可能需要做更多的维护操作，但其时间复杂度仍然维持在 $O(\log n)$ 的级别。

平衡二叉树相对于其它数据结构而言具有很多优势，例如较高的搜索效率（在绝大多数应用中 $O(\log n)$ 的搜索效率已经完全可以接受）以及较好的空间效率（它所占用的空间随树中结点的增多而增长，不像散列表那样每张表都需要大量的空间）。平衡二叉树的缺点是编程难度高，成功写完并调试出一个能用的红黑树所需要的时间不亚于你完成实验二所需的时间。不过如果你真的想要使用类似于红黑树的数据结构，也可以从其它地方（例如Linux内核代码中）寻找别人写好的红黑树源代码。

散列表：

散列表是一种可以达到搜索效率极致的数据结构。一个好的散列表实现可以让插入、查找和删除的平均时间复杂度都达到 $O(1)$ 。同时，与红黑树等操作复杂的数据结构不同，散列表在代码实现上也很简单：申请一个大数组，计算一个散列函数的值，然后根据该值将对应的符号放到数组相应下标的位置即可。对于符号表来说，一个最简单的散列函数（即hash函数）可以把符号名中的所有字符相加，然后对符号表的大小取模。你可以寻找更好的hash函数，这里我们提供一个不错的选择，由P.J. Weinberger²所提出：

```
1 unsigned int hash_pjw(char* name)
2 {
3     unsigned int val = 0, i;
4     for (; *name; ++name)
5     {
6         val = (val << 2) + *name;
7         if (i = val & ~0x3fff) val = (val ^ (i >> 12)) & 0x3fff;
8     }
9     return val;
10 }
```

¹ 《数据结构与算法分析——C语言描述》，Mark Allen Weiss著，冯舜玺译，机械工业出版社，第80、351和89页，2004年。

² http://en.wikipedia.org/wiki/Peter_J._Weinberger。

需要注意的是，代码第7行的常数（0x3fff）确定了符号表的大小（即16384），用户可根据实际需要调整此常数以获得大小合适的符号表。如果散列表出现冲突，则可以通过在相应数组元素下面挂一个链表的方式（称为open hashing或close addressing¹方法，推荐使用），或再次计算散列函数的值而为当前符号寻找另一个槽的方式（称为open addressing或者rehashing²方法）来解决。如果你还知道一些更酷的技术，如multiplicative hash function以及universal hash function³，那将会使你的散列表的元素分布更加平均一些。由于散列表无论在搜索效率和编程难度上的优异表现，它已经成为符号表的实现中最常被采用的数据结构。

Multiset Discrimination:

虽然散列表的平均搜索效率很高，但在最坏情况下它会退化为 $O(n)$ 的线性查找，而且几乎任何确定的散列函数都存在某种最坏的输入。另外，散列表所要申请的内存空间往往比输入程序中出现的符号的数量还要多，较为浪费。如果我们能只为输入程序中出现的每个符号单独分配一个编号和空间，那岂不是既省空间又不会有冲突吗？Multiset discrimination⁴就是基于这种想法。在词法分析部分，我们先统计输入程序中出现的符号（包括变量名、函数名等），然后把符号按照名字进行排序，最后申请一张与符号总数量一样大的符号表，查表功能可通过基于符号名的二分查找实现。

3.2.3 支持多层作用域的符号表

如果你的编译器不需要支持变量的作用域（即不需要实现要求2.2），那可以跳过本节内容，不会对实验二的完成产生负面的影响。否则，请考虑下面这段代码：

```

1  ...
2  int f()
3  {
4      int a, b, c;
5      ...
6      a = a + b;
7      if (b > 0)
8      {
9          int a = c * 2;
10         b = b - a;
11     }
12     ...
13 }
```

¹ 《计算机程序设计艺术 第3卷 排序与查找》，Donald E. Knuth著，苏运霖译，国防工业出版社，第496页，2002年。

² 《计算机程序设计艺术 第3卷 排序与查找》，Donald E. Knuth著，苏运霖译，国防工业出版社，第501页，2002年。

³ 《算法导论》，Thomas H. Corman等著，潘金贵、顾铁成、李成法和叶懋译，机械工业出版社，第138和139页，2007年。

⁴ 《Engineering a Compiler》，第2版，Keith D. Cooper和Linda Torczon著，Morgan Kaufmann出版社，第256和751页，2011年。

函数 f 中定义了变量 a ，在 if 语句中也定义了一个变量 a 。如果要支持作用域，那么：第一，编译器不能在“ $int\ a = c * 2;$ ”这个地方报错；第二，语句“ $a = a + b;$ ”中的 a 的值应该取外层定义中 a 的值，语句“ $b = b - a;$ ”中的 a 的值应该是 if 语句内部定义的 a 的值，而这两个语句中 b 的值都应该取外层定义中 b 的值¹。那么如何使得我们的符号表支持这样的行为呢？

第一种方法是维护一个符号表栈。假设当前函数 f 有一个符号表，表里有 a 、 b 、 c 这三个变量的定义。当编译器发现函数中出现了一个被“{”和“}”包含的语句块（在C—中就相当于发现了CompSt语法单元）时，它会将 f 的符号表压栈，然后新建一个符号表，这个符号表里只有变量 a 的定义。当语句块中出现任何表达式使用到某个变量时，编译器先查找当前的符号表，如果找到就使用这个符号表里的该变量，如果找不到则顺着符号表栈向下逐个符号表进行查找，使用第一个查找成功的符号表里的相应变量。如果查遍所有的符号表都找不到这个变量，则报告当前语句出现了变量未定义的错误。每当编译器离开某个语句块时，会先销毁当前的符号表，然后从栈中弹一个符号表出来作为当前的符号表。这种符号表的维护风格被称为**Functional Style**。该维护风格最多会申请 d 个符号表，其中 d 为语句块的最大嵌套层数。这种风格比较适合于采用链表或红黑树数据结构的符号表实现。假如你的符号表采用的是散列表数据结构，申请多个符号表无疑会占用大量的空间。

另一种维护风格称作**Imperative Style**，它不会申请多个符号表，而是自始至终在单个符号表上进行动态维护。假设编译器在处理到当前函数 f 时符号表里有 a 、 b 、 c 这三个变量的定义。当编译器发现函数中出现了一个被“{”和“}”包含的语句块，而在这个语句块中又有新的变量定义时，它会将该变量插入 f 的符号表里。当语句块中出现任何表达式使用某个变量时，编译器就查找 f 的符号表。如果查找失败，则报告一个变量未定义的错误；如果查表成功，则返回查到的变量定义；如果出现了变量既在外层又在内层被定义的情况，则要求符号表返回最近的那个定义。每当编译器离开某个语句块时，会将这个语句块中定义的变量全部从表中删除。

Imperative Style对符号表的数据结构有一定的要求。图1是一个满足要求的基于十字链表和open hashing散列表的**Imperative Style**的符号表设计。这种设计的初衷很简单：除了散列表本身为了解决冲突问题所引入的链表之外，它从另一维度也引入链表将符号表中属于同一层

¹ 我们通常使用的程序设计语言（包括C、C++以及Java）其作用域规则都来源于Algol，即内层的变量定义总会覆盖外层的变量定义。

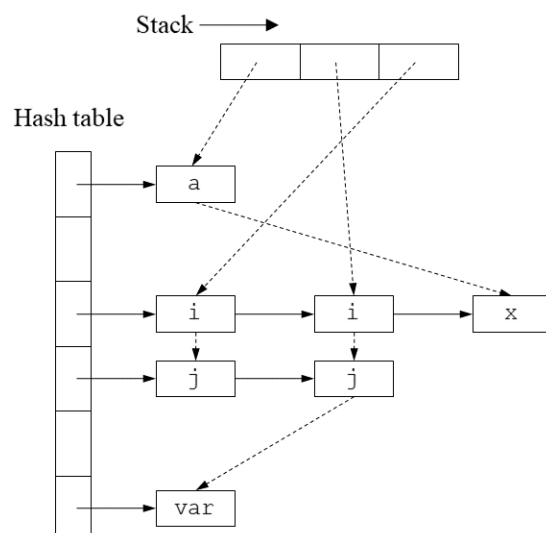


图1. 基于十字链表和open hashing散列表的符号表。

作用域的所有变量都串起来。在图中，a、x同属最外层定义的变量，i、j、var同属中间一层定义的变量，i、j同属最内层定义的变量。其中i、j这两个变量有同名情况，被分配到散列表的同一个槽内。每次向散列表中插入元素时，总是将新插入的元素放到该槽下挂的链表以及该层所对应的链表的表头。每次查表时如果定位到某个槽，则按顺序遍历这个槽下挂的链表并返回这个槽中符合条件的第一个变量，如此一来便可以保证：如果出现了变量既在外层又在内层被定义的情况，符号表能够返回最内层的那个定义（当然最内层的定义不一定在当前这一层，因此我们还需要符号表能够为每个变量记录一个深度信息）。每次进入一个语句块，需要为这一层语句块新建一个链表用来串联该层中新定义的变量；每次离开一个语句块，则需要顺着代表该层语句块的链表将所有本层定义变量全部删除。

如何处理作用域是语义分析的一大重点也是难点。考虑到实现难度，实验二并没有对作用域作过多要求，但现实世界中的动态作用域将更难实现，某些与作用域相关的问题甚至涉及代码生成与运行时刻环境！

3.2.4 类型表示

“类型”包含两个要素：一组值，以及在这组值上的一系列操作。当我们在某组值上尝试去执行其不支持的操作时，类型错误就产生了。一个典型程序设计语言的类型系统应该包含如下四个部分：

- 1) 一组基本类型。在C—语言中，基本类型包括int和float两种。
- 2) 从一组类型构造新类型的规则。在C—语言中，可以通过定义数组和结构体来构造新

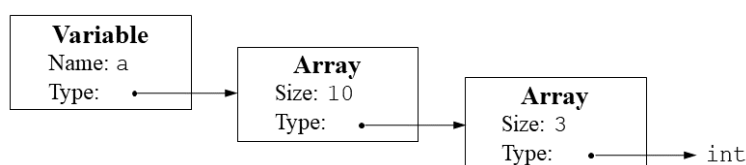


图2. 多维数组的链表表示示例。

的类型。

3) 判断两个类型是否等价的机制。在C语言中，默认要求实现名等价，如果你的程序需要完成要求2.3，则需实现结构等价。

4) 从变量的类型推断表达式类型的规则。

目前程序设计语言的类型系统分为两种：**强类型系统（Strongly Typed System）**和**弱类型系统（Weakly Typed System）**。前者在任何时候都不允许出现任何类型错误，而后者可以允许某些类型错误出现在运行时刻。强类型系统的语言包括Java、Python、LISP、Haskell等，而弱类型系统的语言最典型的代表就是C和C++语言¹。

编译器尝试去发现输入程序中的类型错误的过程被称为是**类型检查**。根据进行检查的时刻的不同，类型检查可被划分为两类，即**静态类型检查（Static Type Checking）**和**动态类型检查（Dynamic Type Checking）**。前者仅在编译时刻进行类型检查，不会生成与类型检查有关的任何目标代码，而后者则需要生成额外的代码在运行时刻检查每次操作的合法性。静态类型检查的好处是生成的目标代码效率高，缺点是粒度比较粗，某些运行时刻的类型错误可能检查不出来。动态类型检查的好处是更加精确与全面，但由于在运行时执行了过多的检查和维护工作，故目标代码的运行效率往往比不上静态类型检查。

关于什么样的类型系统更好，人们进行了长期、激烈而又没有结果的争论。动态类型检查语言更适合快速开发和构建程序原型（因为这类语言往往不需要指定变量的类型²），而使用静态类型检查语言写出来的程序通常拥有更少的错误（因为这类语言往往不允许多态）。强类型系统语言更加健壮，而弱类型系统语言更加高效。总之，不同的类型系统特点不一，目前还没有哪种选择在所有情况下都比其它选择来得更好。

¹ 有关类型系统强弱的定义在不同的文献中不尽相同，例如另一种说法是，强类型系统要求每个变量在定义时都必须赋予一个类型，并且语言本身很少做隐式类型转换。按照这种标准，C和C++语言就应该算是强类型语言，而那些类型系统比C还弱的像Basic、JavaScript才算是弱类型语言。

² 对于那些对变量没有类型限制的语言，有一种生动形象的说法是，这类语言采用了“鸭子类型系统”（duck typing）：如果一个东西看起来像一只鸭子、叫起来也像一只鸭子，那么它就是一只鸭子（if it walks like a duck and quacks like a duck, it's a duck）。

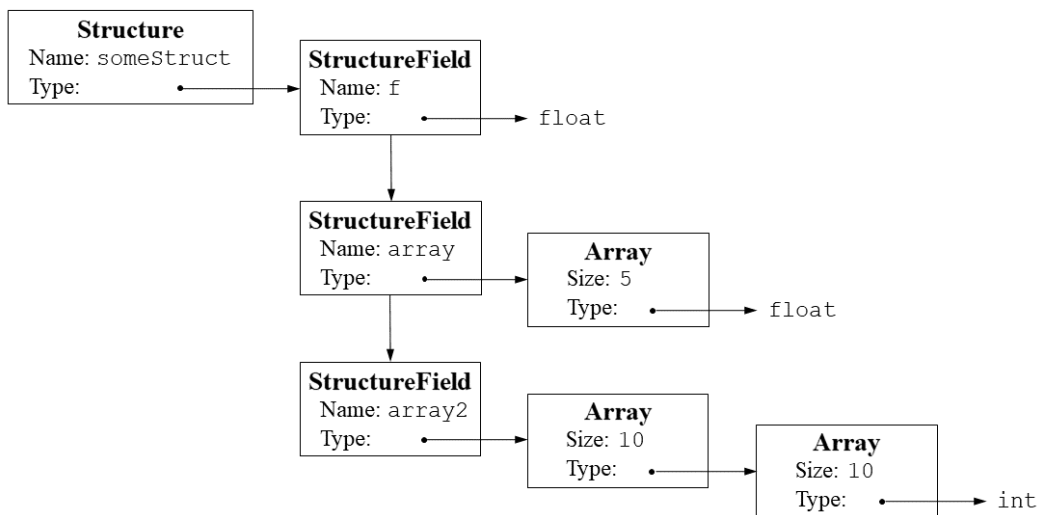


图3. 结构体的链表表示示例。

介绍完基本概念后，我们来考察实现上的问题。如果整个语言中只有基本类型，那么类型的表示将会极其简单：我们只需用不同的常数代表不同的类型即可。但是，在引入了数组（尤其是多维数组）以及结构体之后，类型的表示就不那么简单了。想像一下如果某个数组的每一个元素都是结构体类型，而这个结构体中又有某个域是多维数组，那么该如何去表示呢？

最简单的表示方法还是链表。多维数组的每一维都可以作为一个链表结点，每个链表结点存两个内容：数组元素的类型，以及数组的大小。例如，`int a[10][3]`可以表示为图2所示的形式。

结构体同样也可以使用链表保存。例如，结构体`struct SomeStruct { float f; float array[5]; int array2[10][10]; }`可以表示为图3所示的形式。

在代码实现上，你可以使用如下定义的Type结构来表示C语言中的类型：

```

1 typedef struct Type_* Type;
2 typedef struct FieldList_* FieldList;
3
4 struct Type_
5 {
6     enum { BASIC, ARRAY, STRUCTURE } kind;
7     union
8     {
9         // 基本类型
10        int basic;
11        // 数组类型信息包括元素类型与数组大小构成
12        struct { Type elem; int size; } array;
13        // 结构体类型信息是一个链表
14        FieldList structure;
15    } u;
16 };
17
18 struct FieldList_
19 {
20     char* name; // 域的名字
21     Type type; // 域的类型

```

```
22   FieldList tail; // 下一个域
23 };
```

同作用域一样，类型系统也是语义分析的一个重要的组成部分。**C**语言属于强类型系统，并且进行静态类型检查。当我们尝试着向**C**语言中添加更多的性质，例如引入指针、面向对象机制、显式/隐式类型转换、类型推断等时，你会发现实现编译器的复杂程度会陡然上升。一个严谨的类型检查机制需要通过将类型规则转化为形式系统，并在这个形式系统上进行逻辑推理。为了控制实验的难度我们可以无需这样费事，但应该清楚实用的编译器内部类型检查要复杂的多。

3.2.5 语义分析提示

实验二需要在实验一的基础上完成，特别是需要在实验一所构建的语法树上完成。实验二仍然需要对语法树进行遍历以进行符号表的相关操作以及类型的构造与检查。你可以模仿SDT在Bison代码中插入语义分析的代码，但我们更推荐的做法是，Bison代码只用于构造语法树，而把和语义分析相关的代码都放到一个单独的文件中去。如果采用前一种做法，所有语法结点的属性值请尽量使用综合属性；如果采用后一种做法，就没有这些限制。

每当遇到语法单元ExtDef或者Def，就说明该结点的子结点们包含了变量或者函数的定义信息，这时候应当将这些信息通过对子结点们的遍历提炼出来并插入到符号表里。每当遇到语法单元Exp，说明该结点及其子结点们会对变量或者函数进行使用，这个时候应当查符号表以确认这些变量或者函数是否存在以及它们的类型是什么。具体如何进行插入与查表，取决于你的符号表和类型系统的实现。实验二要求检查的错误类型较多，因此你的代码需要处理的内容也较复杂，请仔细完成。还有一点值得注意，在发现一个语义错误之后不要立即退出程序，因为实验要求中有说明需要你的程序有能力查出输入程序中的多个错误。

实验要求的必做内容共有17种语义错误需要检查，大部分只涉及到查表与类型操作，不过有一个错误例外，那就是有关左值的错误。简单地说，左值代表地址，它可以出现在赋值号的左边或者右边；右值代表数值，它只能出现在赋值号的右边。变量、数组访问以及结构体访问一般既有左值又有右值，但常数、表达式和函数调用一般只有右值而没有左值。例如，赋值表达式 $x = 3$ 是合法的，但 $3 = x$ 是不合法的； $y = x + 3$ 是合法的，但 $x + 3 = y$ 是不合法的。简单起见，你可以只从语法层面来检查左值错误：赋值号左边能出现的只有ID、Exp LB Exp RB以及Exp DOT ID，而不能是其它形式的语法单元组合。最后5种语义错误都与结构体有关，结构体我们前面提到过，可以使用链表进行表示。

要求2.1与函数声明有关，函数声明需要你在语法中添加产生式，并在符号表中记录每个函数当前的状态：是被实现了，还是只被声明未被实现。要求2.2涉及作用域，作用域的实现方法前文已经讨论过。要求2.3为实现结构等价，对于结构等价来说，你只需要在判断两个类型是否相等时不是直接去比较类型名，而是针对结构体中的每个域逐个进行类型比较即可。