

词法分析与语法分析

程序功能(必做+选做1.1)

数组访问错误处理

对于实验说明中,

输入文件的同一行中保证不出现多个错误, 你的程序需要将这些错误全部报告出来, 每一条错误提示信息在输出中单独占一行。

程序对数组访问错误的处理方式和讲义中稍有不同。

```
int main()
{
    float a[10][2];
    int i;
    a[5,3] = 1.5;
    if (a[1][2] == 0) i = 1 else i = 0;
}
```

程序将报告三个错误, 其中第4行有两个错误, 一个是“]”缺失, 一个是语句不合理。

```
Error type B at Line 4: Missing "]"
Error type B at Line 4: Unvalid statement.
Error type B at Line 5: Missing ";"
```

这是为错误恢复提供便利, 如果单纯将错误认定为“]”缺失, 程序将无法处理如下输入文件, 因为不能确保任何“]”缺失都可以有相应的恢复策略。

```
int main()
{
    int a[10];
    int b[10];
    int c;
    c = b[a[1];
}
```

输出:

```
Error type B at Line 6: Missing "]"
```

C--语言和C语言的文法差异处理

对于C--语言中没有的标识符 `for` , `double` 等进行了恢复处理, 对输入文件中类似的标识符进行识别, 如下,

```
int main()
{
    int float = 1;
    float lerr = 1.0;
}
```

输出:

```
Error type B at Line 3: Unvalid name of variable.
Error type B at Line 4: Unvalid identifier.
```

实现思路

文件组织

Code目录下:

```
.
├─ Makefile
├─ lexical.l
├─ main.c
├─ parse.c
├─ parse.h
└─ syntax.y
```

数据结构

定义 `TreeNode` 类型结构体, 建立语法多叉树。

```
typedef struct TreeNode {
    char name[32]; //词法单元
    char text[32]; //词素
    int lineno;    //行号
    int childsum;  //子节点数目
    struct TreeNode* child[MAX_CHILD_NUM];
} Node;
```

主要功能函数如下,

```
Node* createNode(char*, char*);
void insertNode(int, Node*, ...);
void printNode(char*, char*);
void printTree(Node *, int);
```

实现过程

1. **检验词法分析程序** 根据C--文法构造相应的正则表达式，并额外定义了常见错误形式的正则表达式(主要针对C--和C语言的文法差异)，充分利用了正则表达式的匹配规则，权衡了正确性了代码的可读性，如下，
 - 不合法的十进制
 - 不合法的八进制
 - 不合法的十六进制
 - 不合法的单精度浮点数
 - +=, -=, /*, */, //
2. **添加文法的产生式** 根据C--语言文法添加产生式，并通过 `%left`、`%right` 和 `%nonassoc` 以及它们的顺序对终运算符的优先级(precedence)以及结合性(associativity) 进行规定。
3. **恢复解决文法歧义** 悬空else问题的解决方案参照讲义中增加 `LOWER_THAN_ELSE` 算符完成，其他二义文法均使用词法单元error进行错误恢复处理。
4. **构建并打印语法树** 先序遍历多叉树，按照讲义要求递归打印语法树节点。

代码风格

- 把功能的底层实现封装到函数中，增加代码的可读性。
- 开启宏定义 `#define PRINT_DEBUG` 进入调试模式，可以额外打印结点的词法单元类型，并在语义恢复之后输出语法树(即使存在语法错误)。

写在最后

和室友聊天说起编译原理实验的讲义，我们都觉得相对PA和OS，这本指导手册的可读性更高一些。坦白的说，之前的两个实验使我对稍具规模的实验任务产生了自然而然的畏惧感。不过欣喜的是，我感受到了实验任务和理论课程的直接联系非常紧密，这使我在完成实验一的过程中有不小的成就感。

最后非常感谢助教学长和许畅老师对我的问题详尽的解答！

语义分析

程序功能(必做+选做2.3)

实现思路

数据结构

符号表的设计

符号表使用Hash表，将基本类型、数组类型、结构类型、函数类型变量存储在一个统一的符号表中，提高插入和查找的效率。

符号表：

```
struct SymbolTable {
    char name[32];
    Type type;
    int lineno;
    bool occupied;
} symboltable[0x4000];
```

变量类型：

```
typedef struct Type_
{
    enum { BASIC, ARRAY, STRUCTURE, FUNCTION } kind;
    union
    {
        int basic;
        struct { Type elem; int size; } array;
        FieldList structure;
        Function function;
    } u;
} Type_;
```

其中，枚举类型变量 BASIC，ARRAY，STRUCTURE，FUNCTION 变量类型，在 getNodeType 和 handleExp 中作为返回类型的标识。实现假设每种类型的变量个数不会超过10个，因此使用了 类型 * 10 + 序号 作为标识。

函数类型：

```
typedef struct Function
{
    enum { INT, FLOAT } returnvalue;
    int argsum;
    int argbasic[8];
} Function;
```

结构体类型：

```
typedef struct FieldList_
{
    char name[32];
    Type type;
    FieldList tail;
} FieldList_;
```

结构体的管理

使用链表结构存储结构体，包括结构体名称和其中的成员变量。

```
typedef struct StructNode_ {
    int no;
    char name[32];
    FieldList children;
    StructNode next;
} StructNode_;
```

文法属性的处理

重新遍历实验一中建立的树，对文件进行语义分析。使用地址传递函数参数和返回值相结合的方式对文法的继承属性和综合属性进行处理。

```
// ExtDefList -> ExtDef ExtDefList | ε
void handleExtDefList(Node* root) {
    if (Childsum == 0) return;
    // Iteration is better than recursion.
    while (root != NULL) {
        handleExtDef(Child(0));
        root = Child(1);
    }

#ifdef HANDLE_DEBUG
    printf("handleExtDefList bingo\n");
#endif
}
```

全局变量的使用

由于函数传递参数的缺陷，使用全局变量进行相关辅助。

- `PARA` 用于标识所在位置是否属于函数的参数定义。
- `IN_STRUCT` 用于标识结构体嵌套的深度。
- `FLAG` 用于标识 `ExtDef` 的产生式的类型。

实现过程

1. 设计数据结构，规划整体实现流程
2. 将变量加入符号表，处理相关语义错误
3. 进行变量声明和类型检查
4. 构造测试用例进行调试，尽量补充遗漏的特例

代码风格

优点：

- 把功能的底层实现封装到函数中，增加代码的可读性。
- 在实验一的基础上增加宏定义 `#define PHASE_SEM` 进入调试模式，可以打印插入符号表的相关信息。（将实验一中的宏 `PRINT_DEBUG` 更名为 `PRINT_TRACE` 。）
- 使用多个参数的宏定义，子孙节点的访问和控制，最多支持6个参数。

有待提高：

- 对指针的操作不是很熟练，相关代码实现结构性相对较差，给实现和阅读带来不便。
- 对于结构体嵌套和多维数组以及一些局部变量的处理存在硬编码。

写在最后

报告至此，我的心情十分复杂。实现编译器的语义分析的复杂度远超出我之前的想象。体会最深的一点是有必要在动手之前设计出合理的数据结构，不然将会带来极大的困难。设计模式是否友好奠定了整个工程的基础。大学三年，第一次在写实验的过程中感受如此深刻。多次的代码重构使我有点烦躁，但不得不承认，加深了我对理论知识的理解，提高了自己的动手能力，感觉收获很大。

中间代码生成

程序功能(必做+部分选做)

程序接口说明

```
./parser in out
```

其中 `in` 为输入文件，转换生成的中间代码将存储在 `out` 文件中。输出目标文件缺省值为 `out.ir`，另外还会生成 `original.ir` 文件，写入内容为未经过优化的代码。

中间代码生成

将中间代码输出成线性结构，存储在双向链表的结构中。

中间代码优化

实现了常量折叠和代数运算简化的优化功能。具体说来，先定位将常量赋值给临时变量的赋值语句，并在其后使用这个临时变量时将其替换为对应的常量；再对可能会由此产生一些可以直接进行计算的四则运算表达式进行求值。之后迭代进行常量赋值表达式消除和代数运算简化，直到产生的中间代码条数不发生变化，具体实现在 `optimize.c` 文件中的 `interCodeOptimize()` 函数。

实现思路

相关源码文件

```
.
├─ ir.c
├─ translate.c
└─ optimize.c
```

数据结构设计

对于每一条中间代码，使用 `InterCode` 结构体类型，设计 `Operand_` 结构体类型表示每一个对应操作的相关信息。定义在 `ir.h` 文件中。

```

typedef struct Operand_ {
    enum { OP_VARIABLE, OP_TEMPORARY, OP_CONSTANT, OP_ADDRESS, OP_LABEL, OP_FUNCTION, OP_VALUE, OP_SIZE } kind;
    union {
        int no;
        int value;
    } u;
} Operand_;

typedef struct InterCode {
    enum {
        IR_ADD, IR_SUB, IR_MUL, IR_DIV, IR_IFGOTO,
        IR_LABEL, IR_FUNC, IR_GOTO, IR_RET, IR_DEC, IR_ARG, IR_PARAM, IR_READ, IR_WRITE,
        IR_ASSIGN, IR_GETADD, IR_GETVAL, IR_SETVAL, IR_CALL
    } kind;

    union {
        struct { Operand x; } op;
        struct { Operand x, y; } biop;
        struct { Operand x, y, z; } triop;
        struct { Operand x, y, z; Relop relop; } ifop;
    };
} InterCode;

```

使用双向链表的结构存储中间代码，具体实现参考了Linux内核的 `list_head` 结构。

```

typedef struct InterCodes
{
    InterCode intercode;
    struct list_head list;
    bool isLeader;
} InterCodes;

```

翻译流程的架构设计

类似实验二中的自顶向下分析方法，从 `Program` 开始进行逐层函数调用，在底层进行中间代码生成，必要的时候返回调用者 `place` 参数，具体功能实现参照翻译模式。临时变量和标号的计数使用了全局变量。对于变量和函数的符号表，原则上可以沿用语义分析中的设计，我对之前的版本进行了完善。

数组的翻译模式设计

为提高编译器运行效率，使用循环而非递归进行数组的翻译模式设计，具体实现在 `translate.c` 中的 `translate_Array`。

以四维数组 `int array[7][8][9][10]` 为例说明设计的翻译模式，设 `&v1` 表示 `array` 的地

址, `v2`, `v3`, `v4`, `v5` 分别对应 `i`, `j`, `k`, `s` 的值。计算数组元素 `array[i][j][k][s]` 的地址, 有

$$\begin{aligned} & \&v1 + v2 * 8 * 9 * 10 * 4 + v3 * 9 * 10 * 4 + v4 * 10 * 4 + v5 * 4 \\ & = \&v1 + v2 * 8 * 9 * 10 * 4 + v3 * 9 * 10 * 4 + v4 * 10 * 4 + v5 * 4 \\ & = \&v1 + ((v2 * 8 + v3) * 9 + v4) * 10 + v5) * 4 \end{aligned}$$

对应的中间代码翻译为

```
t2 := #8
t3 := #9
t4 := #10

t5 := v2 * t2
t6 := t5 * t3
t7 := t6 * t4
t8 := t7 * #4
t15 := &v1 + t8

t9 := v3 * t3
t10 := t9 * t4
t11 := t10 * #4
t16 := t15 + t11

t12 := v4 * t4
t13 := t12 * #4
t17 := t16 + t13

t14 := v5 * #4
t18 := t17 + t14
```

写在最后

个人觉得中间代码生成部分的必做实验难度要略低于语义分析部分。按照讲义上的指导一步步完成相对比较顺利。其中, 耗时最长的数组的翻译模式设计, 我使用了从特殊到一般的归纳法。先写出低维数组的翻译函数, 然后类比推理出任意维的数组的翻译函数。另外, 对于中间代码的优化, 虽然我基本理解了课上老师介绍的使用DAG进行基本块内部优化的理论, 但是在实践过程中遇到了重重的困难。经过两天的尝试, 为了尽量保证中间生成代码的正确性, 我最后选择了相对保守的、实现复杂度较低的优化策略, 在实现复杂度、优化效果和正确性之前进行了权衡。

目标代码生成

程序功能(必做)

程序接口说明

```
./parser in out
```

其中 `in` 为输入文件，转换生成的目标代码文件将存储在 `out` 中。输出目标文件缺省值为 `out.s`，另外还会生成 `out.ir` 文件，写入内容为实验三的中间代码。

为了便于测试，同时利用目标代码的语法特点，在输出文件中将所有中间代码保留并在写入文件之前加 `#` 作为注释。

实现思路

相关文件

主要的功能函数在 `generate.c` 中完成，其对应的头文件为 `generate.h`，`setupTable` 函数对应第一次扫描，`codeGenerate` 函数对应第二次扫描，其中 `setTempVar`、`setArrayVar` 和 `getRegNo` 函数分别实现函数符号表中变量的相关信息存储和读取。

数据结构

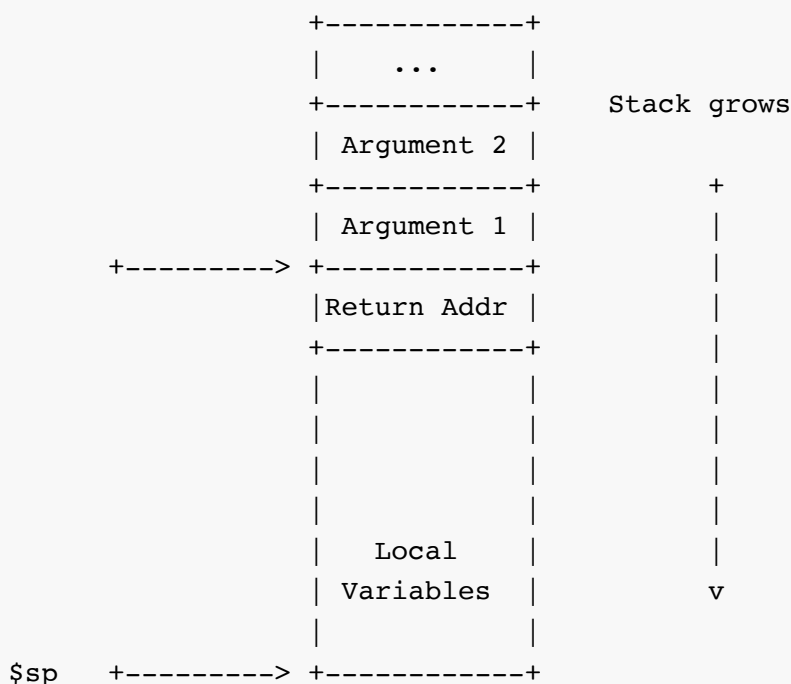
对之前设计的函数符号表略加修改和扩充，记录目标代码中每个函数相关的变量对应的伪寄存器(假设寄存器数目无限多)。

```
typedef struct Symbol {
    int no[512];
    int type[512];
    int scale[512];
    int size;
} Symbol;

typedef struct FunctionTable {
    SymbolTable symboltable;
    Symbol para;
    Symbol temp;
    int argno;
} FunctionTable;
```

函数栈帧结构

为了简化实现，我不仅使用了最朴素的寄存器分配算法，还没有使用 `$fp` 指针，也没有使用 `$s0-$s7` 系列寄存器，同时把调用函数使用的所有参数都存储在函数的栈帧中。函数的活动记录结构简化如下。



实现过程

1. 确定指令选择机制

首先根据讲义中的表11等进行中间代码与伪指令之间的对应，假设寄存器有无限多个，对双向链表中存储的中间代码进行逐条翻译。由于C++语言和实验的限制，涉及的指令较少，因此按照给定模式进行翻译的难度并不大。

2. 确定寄存器分配算法

为简化实现过程，使用朴素寄存器分配算法，将所有变量或临时变量都放在内存里。

对中间代码进行两次扫描。第一次扫描的任务是，完成局部变量在栈中的空间分配，并将记录相关信息。第二次扫描负责根据中间代码和已知的栈帧信息，完成中间代码的逐条翻译。

3. 活动记录的布局方式

完成对 `ARG`、`PARAM`、`RETURN` 和 `CALL` 语句的翻译。在这一部分，为使测试用例尽可能全面，我尝试回忆上个星期许老师在课堂介绍的“编译器测试”专题讲座，设计了包含多个参数的函数、使用递归函数等程序进行调试分析，尽可能提高testcase的覆盖率。

这一部分的要点是关于栈帧指针 `$sp` 位置的计算，一旦发生了函数调用使得指针位置变化，要注意改变局部变量访问时的指针位置。

4. 数组类型变量的处理

与数组类型变量相关的中间代码类型较为单一，仅涉及 `IR_ASSIGN` 和 `IR_ADD` 两种类型。对操作数 (`OP_VALUE`，`OP_TEMPORARY`，`OP_VARIABLE`，`OP_ADDRESS`，`OP_CONSTANT`) 的不同类型

进行讨论即可。至于栈帧空间的分配，和其他简单类型的变量类似，区别在于分配栈帧空间大小。

写在最后

之前写过编译实验的一些同学说目标代码生成是最简单的一个实验，如果使用最朴素的寄存器分配算法。但是即使选用了最简单的实现方式，我依然挣扎了很长时间。写完之后回看，倒是确实是这样hhh，毕竟在不考虑编译器效率的情况下，无需设计繁杂的数据结构和高级的算法，按照给定的模式，一步一步走即可完成。

遇到最大的坑是函数的活动记录，在之前的《计算机系统基础》课程中，对这部分的知识一直是完全不懂的状态，或者说不知道自己懂不懂，直到编译课，感觉自己好像有那么一点点懂了。虽然想来想去其实也没什么，大概的思想是需要用什么就去拿什么，不用的东西先放在一边，确保想用的时候能找到就可以了。到这个时候才稍微理解了当时袁春风老师“做菜取食材”的比喻，感觉自己慢的不只是一怕orz。

整个实验过程中，收获真的是很大，有过怀疑自己，怀疑人生的时候，连续写实验的几天晚上梦里都会惊醒(一点都不夸张)，每天和爸爸妈妈同学聊天开口三句话之内必含“编译”，那时手机输入法打出“编译”自动匹配的下一个词就是“💊”。不过写完之后确实对理论的知识有了更加深刻的理解。这一点我在每一次报告都会说，真的是最深刻的体会，没有之一。之前的PA和OS实验都是连原理没有搞清楚，那又怎么可能写好呢？课程中理论和实践关系，想必该是相辅相成的。实验需要理论知识作为基础，书本上的知识又可以在做完实验之后有更加深刻的理解。可笑的是，之前的自己似乎全然没有意识到。

emmm扯远了，有同学希望编译实验给出框架代码，我个人觉得现在的模式很好，不需要修改。虽然讲义需要读很多遍才能发现一些细微的奥秘，但个中滋味唯努力着清啊。

遗憾总还是有的，比如前两次实验的完成很不好，不够细心，写之前没有很好的整体意识，数据结构不友好，代码风格也不很优雅。另外，对关于编程语言的限制也使我更加深入地了解了一些关于C语言的特性，同时也对C++的特性有了那么一点小小的喜欢😋(觉得实验四用C++可能会好写一点？.....)。Anyway，实验接近尾声啦！还有五门考试等着复习(yuxi)，这个学期的最后一个小假期一如既往地奉献给了编译原理实验。看着朋友圈周游世界的同学们，再看看一次次的commit，假装脸上露出了欣慰的微笑！