

Compiler 2020 Manual

关于课程：

编译器设计是ACM班的传统课程，这门课程旨在锻炼大家的编程能力和工程能力。往年的课程都是进行天梯制度赋分，由于大家都写累了好的数据点就贡献的越来越少了，往后逐渐变成了面向数据编程，从2018级开始回归到编译本身可能更为重要。因此，我们修正了原有语言中描述不清晰的部分，按照Standard C++和Java的语言定义方式给出一个定义，并且按照两个语言的标准编译测试集合制定出一个属于我们的语言的标准集合。欢迎大家提出修改意见和建议。

注：本文参照ISO/IEC 14882:2017 Programming Language C++以及往年的Manual做出修改。

部分术语定义：

1. 未定义行为：指规范并没有定义该情况发生时语言的表现。初衷是为了给同学们提供一些自己发挥的空间，在测试数据里，这些没有定义的情况是不会发生的。可以认为未定义行为是类似于运行时会错误的东西，由于在编译阶段无法确定，因此我们就保证我们的代码不会出现。

例子：对长度超过1M的代码的编译是未定义的。

解释：我们的测试集中没有长度超过1M的代码。

2. 语法错误：指代码违反规范的行为，你的编译器应返回非0返回值作为编译错误指示信息（必须，作为评测之一）以及你的提示信息（可选，给自己看的）。
3. 源代码：你的编译器即将编译的代码。
4. 预留位

语言基本结构：

一个标准的Mx*语言包含有以下部分：大于等于1个函数定义，（可选）类定义，（可选）全局变量声明。

其中：

（1）函数定义中有且仅有1个的名字可以为main，main函数的定义仅可为`int main()`，不符合此定义的main函数或者没有定义main函数均视为语法错误。

（2）在Mx*语言中没有接口的定义，所有的函数必须有对应的函数体，反之视为语法错误。

（3）预留

文法规则：

1 编码与符号：

我们称为这个语言是Mx*，这个语言对大小写敏感，可以使用的符号集合如下：

标识符（包括变量标识符、函数标识符、类对象标识符）：26个小写英语字母，26个大写英语字母，0，1，2，3，4，5，6，7，8，9，下划线（_）；

标准运算符：加号（+），减号（-），乘号（*），除号（/），取模（%）；

关系运算符：大于（>），小于（<），大于等于（>=），小于等于（<=），不等于（!=），等于（==）；

逻辑运算符：逻辑与（&&），逻辑或（||），逻辑取反（!）；

位运算符：算术右移（>>），算术左移（<<），按位与（&），按位或（|），按位异或（^），按位取反（~）；

赋值运算符：赋值（=）；

自增运算符：自增（++），自减（--）；

分量运算符：对象（.）；

下标运算符：取下标对象（[]）；

优先级运算符：括号（()）；

特殊符号：空格（ ），换行符（'\n'），制表符（'\t'），注释标识符（//）。

不包括在以上符号集合内的符号出现在源代码中视为语法错误。

2 关键字：

**int bool string null void true false if for while break continue return
new class this**

3 空白字符：

空白字符、制表符、换行符在Mx*语言中除了区分词素（Token）以外没有作用。

4 注释：

行注释：从“//”开始到这一行末尾的所有内容都会被作为注释，编译的时候应当自动忽略。

该语言中有且仅有一种这样的注释，剩余在C++语言中用于表示注释的方法在我们的Mx*中被认为是未定义的行为（欢迎做一些尝试实现，可以作为Presentation阶段的演讲内容）。

5 标识符：

标识符的第一个字符必须是英语字母（26个大写英语字母和26个小写英语字母，下同）中的一个。第二个字符开始可以是英语字母、数字或者下划线（_）中的。标识符区分大小写并且长度超过64个字符的标识符是未定义的。

6 常量：

注：没有在以下定义的常量都是未定义的。

6.1 逻辑常量

定义**true**为真，**false**为假。

6.2 整数常量

整数常量以十进制表示，整数常量不设负数，负数可以由正数取负号得到。

编译器至少应该能处理大小范围在 $[-2^{31}, 2^{31})$ 内的整数，首位为0的整数常量是未定义的（整数0除外），大小超过上述范围的整数是未定义的。

6.3 字符串常量

字符串常量是由双引号括起来的字符串。字符串长度最小为0，长度超过255的字符串是未定义的。

字符串中的所有字符必须是可示字符（printable character），空格或者转义字符中的一种。

转义字符有三个：\n表示换行符，\\表示反斜杠，\"表示双引号。

其余出现在C++语言里的转义字符是未定义的。

6.4 空值常量

定义**null**为引用类型没有指向任何值。

7 变量：

注：没有在以下定义的变量类型都是未定义的。

7.1 基础类型

1. **bool**类型：**true**为真，**false**为假。
2. **int**类型：大小范围在 $[-2^{31}, 2^{31})$ 内的整数。
3. **void**类型：表明函数没有返回值的特殊类型，仅仅可以用于函数返回值。
4. **string**类型：字符串是引用类型，可以改变它的值但是本身不能被改变（immutable）。

7.2 数组类型

注：该部分的<typename>指的是类型，可以是基础类型（除外**void**）也可以是类。

<identifier>指的是变量标识符。

数组是一种可以动态创建的引用类型，长度不需要在声明的时候确定。声明语句的语法要求为

<typename>[] <identifier>(<initial sentence>);

例如：`bool[] flag;` 是一句合法声明语句，不加创建的数组在创建后对应变量值为`null`，此时访问数组下标是未定义的。

创建数组可以用`new`关键字创建，创建数组的语句语法要求是

```
(<typename>[]) <identifier> = new <typename>[arraySize];
```

例如：`flag = new bool[10];` 是一种合法的创建方式。创建数组必须制定数组的长度，方括号中仅可以传入一个整型的数。数组长度一定小于 $2^{31} - 1$ 。

在我们的`Mx*`中，所有的数组都是通过动态创建的，我们不支持静态确定数组长度的数组，因此形如：`<typename>[] <identifier>[arraySize]`都是未定义的。

数组内建方法：`<identifier>.size()`返回数组的长度，函数返回值为`int`。对`null`的数组对象执行`.size()`返回数组的长度是未定义的。

多维数组：我们采用交错数组来达到多维数组的效果，交错数组就是数组的数组。声明方法和C#语言保持一致，可以理解为C++语言中`vector`套`vector`的效果。声明交错数组的语句语法要求为

```
<typename>[]... <identifier>(<=<initial sentence>);
```

例如声明一个2维数组的语句可以是：

```
int[][] graph;
```

声明创建交错数组的语法为：

```
(<typename>[]...) <identifier> = new <typename>[outerSize][]..;
```

创建交叉数组需要先创建最外层数组的空间，然后再创建内层数组空间。类似于C++的`std::vector`。

例如声明创建一个2维数组的语句可以是：

```
int[][] graph = new int[3][];  
graph[0] = null; // Valid  
graph[1] = new int[10];  
graph[2] = new int[30];
```

交叉数组的声明创建还可以有一种简单的方法：

```
<typename>[]... <identifier> = new <typename>[size1][size2]..;
```

例如声明创建一个大小为`3*4`的2维数组的语句可以是：

```
int[][] graph = new int[3][4];
```

这个创建方法在Java中支持，并且看上去也比较简洁。

常量数组：这个在我们的语言中是未定义的，但是可以做一下尝试。

8 类：

我们的语言需要面向对象，类的定义的方式如下：

```
int[][] graph = new int[3][4];  
class <ClassIdentifier> {  
    <Type 1> <MemberIdentifier 1>;  
    <Type 2> <MemberIdentifier 2>, <MemberIdentifier 3>..;  
    <Type 3> <FunctionIdentifier>(<FunctionParameterList>){  
        <Expressions and Statements>  
    }  
    <ClassIdentifier>(){ // Can be ignored  
        <Expressions and Statements>  
    }  
}
```

8.1 类成员变量

对于类成员变量，要求必须在构造函数中赋初值。对于没有赋值的成员变量是未定义行为，访问没有赋值的类成员变量也是未定义的。所有的类成员变量都是`public`的，我们对于`private`的对象是未定义的行为。

8.2 类方法

对于类方法，要求和第九部分函数的要求相同（除了构造函数），语法如下：

```
<Type> <FunctionIdentifier>(<FunctionParameterList>){  
    <Expressions and Statements>  
}
```

8.3 类成员访问

对于类成员不论是方法还是变量，都可以用对象标识符.取对象，对于除了字符串`string`的基本类型`int`，`bool`返回一个实值，剩下返回的应当是一个引用。语法如下：

```
<ClassObjectIdentifier>.<ClassMember>;  
OR  
<ClassObjectIdentifier>.<ClassMethod>(<FunctionParameterList>);
```

8.4 类构造函数

构造函数的定义和C++相同，无返回值无参数（有参数的未定义），可以没有构造函数，语法如下：

```
<ClassIdentifier>(){ // Can be ignored  
    <Expressions and Statements>  
}
```

8.5 this指针

this指针返回某个类的引用对象，关键字仅在类作用域内可以使用。不在类作用域内的this应当视为语法错误，this指针作为左值视为语法错误。

| | |
|--|--|
| <pre>class foo { int a; int b; int c; foo test(){ return this; } }</pre> | <pre>class foo { int a; int b; int c; } foo test(){ return this; }</pre> |
| 语法正确 | 语法错误 |

8.6 类之中的未定义行为

析构函数、虚函数、类的继承、接口、权限标示、抽象类、成员的默认初始化表达式、函数重载。

9 函数：

9.1 函数定义

标准的函数定义应该满足如下语法：

```
<ReturnType> <FunctionIdentifier>(<FunctionParameterList>){  
    <Expressions and Statements>  
}
```

注意在Mx*中不支持lambda函数表达式，不支持匿名函数，没有方法声明函数的签名，也不支持在一个函数内嵌套申明另一个子函数或类。

9.2 内建函数

以下函数是系统包括的函数，不需要申明就可以使用。

函数：`void print(string str);`

作用：向标准输出流中输出字符串`str`。

函数：`void println(string str);`

作用：向标准输出流中输出字符串`str`，并且在行尾处输出一个换行符。

函数: `void printInt(int n);`

作用: 向标准输出流中输出数字`n`。

函数: `void printlnInt(int n);`

作用: 向标准输出流中输出数字`n`, 并且在行尾处输出一个换行符。

函数: `string getString();`

作用: 从标准输入流中读取一行并且返回。

函数: `int getInt();`

作用: 从标准输入流中读取一个整数, 遇到空格、回车符、制表符作为分隔, 返回这个整数。

函数: `string toString(int i);`

作用: 把整数`i`转换为字符串。

9.3 函数返回值

如果函数声明的返回值不是`void`, 就必须有`return`语句返回函数返回值, 反之语法错误。`main`函数例外, 可以没有返回值, 此时返回值为0。

10 表达式:

10.1 单目表达式

单目表达式有常量, 标识符变量名。等等

10.2 双目表达式

双目表达式的定义和C++类似, 在类型`int, bool`中, 要求表达式两边的对象类型必须一致而表达式两边的对象的常量/变量属性没有特别要求(除了赋值, 参阅左值部分定义)。数组对象可以和`null`比较但是不能运算。类对象的运算符重载是未定义的。字符串部分的参阅字符串部分定义。

特殊的是: 自增自减运算符在前缀加和后缀加意义下表达式返回本身的值+1或-1的值。

例如: `a = 1; b = ++a; // After Execution: a = 2, b = 2`

11 语句:

11.1 变量声明语句

此处的变量不是类成员变量, 类成员变量的定义参阅类的定义。变量声明语句语法如下:

`<Type 1> <MemberIdentifier 1>, <MemberIdentifier 2>;`

变量在使用之前应当被赋值了, 没有赋值的对象直接使用是未定义行为。

11.2 条件语句

条件语句语法要求如下:

```
if (condition) {  
    <Expressions and Statements if true>  
} else {  
    <Expressions and Statements if false>  
}
```

其中`condition`字段必须返回`bool`值, 并且不能为空。如果`condition`返回了非`bool`值或者空应当视为语法错误。一个`if`语句可以没有`else`部分。

11.3 循环语句

while循环语句语法要求如下:

```
while (condition) {  
    <Expressions and Statements if true>  
}
```

如果`condition`返回了非`bool`值或者空应当视为语法错误。

for循环语句语法要求如下：

```
for (init; condition; incr) {  
    <Expressions and Statements if true>  
}
```

如果`condition`返回了非`bool`值视为语法错误，但是可以空。

11.4 跳转语句

`return/break/continue`语法要求如下：

```
return <Expression>;  
break;  
continue;
```

`return`只在函数中有效，不在函数中的`return`视为语法错误。

12 字符串：

12.1 字符串对象

字符串对象赋值为`null`是语法错误。

12.2 字符串的双目运算

`+`表示字符串拼接

`==, !=`比较两个字符串是否完全一致（不是内存地址）

`<, >, <=, >=`用于比较字典序大小

剩余双目运算符都是语法错误，字符串双目运算符要求两边类型相同，不满足则语法错误。

12.3 字符串的内建方法

函数：`int length();`

使用：`<StringIdentifier>.length();`

作用：返回字符串的长度。

函数：`string substring(int left, int right);`

使用：`<StringIdentifier>.substring(left, right);`

作用：返回下标为`[left, right)`的子串。

函数：`int parseInt();`

使用：`<StringIdentifier>.parseInt();`

作用：返回一个整数，这个整数应该是该字符串的最长前缀。如果该字符串没有一个前缀是整数，结果未定义。如果该整数超界，结果也未定义。

函数：`int ord(int pos);`

使用：`<StringIdentifier>.ord(pos);`

作用：返回字符串中的第`pos`位上的字符的ASCII码。下标从0开始编号。
常量字符串不具有内建方法，使用内建方法的常量字符串未定义。

13 作用域：

13.1 作用域规则

1. 在一段语句中，由`{`和`}`组成的块会引进一个新的作用域。
2. 用户定义函数入口会引入一个新的作用域。
3. 用户定义类的入口会引入一个新的作用域，该作用域里声明的所有成员，作用域为整个类。

4. 全局变量和局部变量不支持前向引用，作用域为声明开始的位置直到最近的一个块的结束位置。
5. 函数和类的声明都应该在顶层，作用域为全局，支持前向引用（forward reference）。
6. 不同作用域的时候，内层作用域可以遮蔽外层作用域的名字。

注意：诸如**for**等表达式没有大括号也会引入一个新的作用域，如下：

```
int a = 0;  
for(;;) int a = 0;
```

可以通过clang编译。

14 命名空间：

所有符号共享一个命名空间，所以在同一个作用域里，变量，函数，和class，都不能同名如果重名视为语法错误，注意作用域规则。

15 左值：

由以下方法给出的对象为左值，可以被赋值。

1. 函数的形参。
2. 全局变量和局部变量。
3. 类的一个成员。
4. 数组对象的一个元素。

我们的Mx*要求至少支持上述四种类型的左值。更多的左值是未定义的（注意不是语法错误）。