# Assignment 2

Harley McPhee 27003226

# Grammars Changes

## Change Remove *

Reason: * is not supported, so we need to remove all * and make the non-terminals preceding the * recursive and give them an empty string to stop, or if this non-terminal is used elsewhere without a * we will need to introduce a new non-terminal to allow multiple elements preceding the *, and keep the original non-terminal preceding * as it is

OLD:
1. \<prog>           → \<classDecl>*\<progBody>
2. \<classDecl>   → class id {\<varDecl>*\<funcDef>*};
3. \<varDecl>       → \<type>id\<arraySize>*;
4. \<arraySize>   → [ integer ]
5. \<funcDef>       → \<funcHead>\<funcBody>;
6. \<fParams>       → \<type>id\<arraySize>*\<fParamsTail>* | ε
7. \<fParamsTail> → ,\<type>id\<arraySize>*
8. \<funcBody>     → {\<varDecl>*\<statement>*}
9. \<varDecl>         → \<type>id\<arraySize>*;
10. Add statementLst because statement is used without \<statement>*
11. \<variable>         → \<idnest>*id\<indice>*
12. \<idnest>           → id\<indice>*.
13. \<indice>           → [\<arithExpr>]
14. \<factor>           → \<variable>
                          | \<idnest>*id(\<aParams>)
                          | num
                          | (\<arithExpr>)
                          | not\<factor>
                          | \<sign>\<factor>
15. \<aParams>       → \<expr>\<aParamsTail>* | ε
16. \<aParamsTail>→ ,\<expr>
17. \<statBlock>       → {\<statement>*} | \<statement> | ε

CHG:
1. \<prog>           → \<classDecl> \<progBody>
2. \<classDecl>   → class id { \<varDecl> \<funcDef> } ; \<classDecl> | ε
3. \<varDecl>       → \<type> id \<arraySize> ;
4. \<arraySize>   → [ integer ] \<arraySize> | ε
5. \<funcDef>       → \<funcHead> \<funcBody> ; \<funcDef> | ε
6. \<fParams>       → \<type> id \<arraySize> \<fParamsTail> | ε

7. &lt;fParamsTail&gt; → , &lt;type&gt; id &lt;arraySize&gt; &lt;fParamsTail&gt; | ε
8. &lt;funcBody&gt;   → { &lt;varDecl&gt; &lt;statementLst&gt; }
9. &lt;varDecl&gt;      → &lt;type&gt; id &lt;arraySize&gt; ; &lt;varDecl&gt; | ε
10. &lt;statementLst → &lt;statement&gt; &lt;statementLst&gt; | ε
11. &lt;variable&gt;      → &lt;idnest&gt; id &lt;indiceLst&gt;
12. &lt;idnest&gt;         → id &lt;indiceLst&gt; . &lt;idnest&gt; | ε
13. &lt;indice&gt;         → [ &lt;arithExpr&gt; ]
14. &lt;factor&gt;         → &lt;variable&gt;
                      | &lt;idnest&gt; id ( &lt;aParams&gt; )
                      | num
                      | ( &lt;arithExpr&gt; )
                      | not &lt;factor&gt;
                      | &lt;sign&gt; &lt;factor&gt;
15. &lt;aParams&gt;     → &lt;expr&gt; &lt;aParamsTail&gt; | ε
16. &lt;aParamsTail&gt;→ ,&lt;expr&gt; &lt;aParamsTail&gt; | ε
17. &lt;statBlock&gt;     → { &lt;statementLst&gt; } | &lt;statement&gt; | ε
18. &lt;indiceLst&gt;     → &lt;indice&gt; &lt;indiceLst&gt; | ε

## Change Remove left-recursion

Reason: Left recursion causes non-determinism in a left to right parser. To remove the left recursion we add a new non-terminal for each RHS that has left recursion.

OLD:
1. &lt;arithExpr&gt;     →  &lt;arithExpr&gt; &lt;addOp&gt; &lt;term&gt; | &lt;term&gt;
2. Add &lt;arithExprD&gt; to remove left recursion from &lt;arithExpr&gt;
3. &lt;term&gt;          → &lt;term&gt;&lt;multOp&gt;&lt;factor&gt; | &lt;factor&gt;
4. Add &lt;termD&gt; to remove left recursion from &lt;term&gt;

CHG:
1. &lt;arithExpr&gt;     →  &lt;term&gt;&lt;arithExprD&gt;
2. &lt;arithExprD&gt; → &lt;addOp&gt;&lt;term&gt;&lt;arithExprD&gt; | ε
3. &lt;term&gt;          → &lt;factor&gt;&lt;termD&gt;
4. &lt;termD&gt;        → &lt;multOp&gt;&lt;factor&gt;&lt;termD&gt; | ε

## Change Allow &lt;varDecl&gt; or &lt;funcDef&gt; before or after each other in &lt;classDecl&gt;

Allow <classDecl> to have either <varDecl> and <funcDef> before or after each other and remove ambiguity that'll be created by it because either <varDecl> and <funcDef> both have id in their first sets

Reason: A variable can technically be declared after a function inside a class, the grammar as it is does not allow that, it requires that variables be declared before, also <classDecl> will then contain an ambiguity, the first set of <varDecl> and <funcDef> contains (int, float, id), which leads to non determinism in <classBody>

OLD:
1. <classDecl>   → class id { <varDecl><funcDef>};<classDecl> | ε

CHG:
1. <classDecl>   → class id {<classBody>};<classDecl> | ε
2. <classBody>   →  <classInDecl><classBody> | ε
3. <classInDecl> → <type>id<postTypeId>
4. <postTypeId> → <arraySize>;  | (<fParams>)<funcBody>


## Change Remove ambiguity caused by <idnest> and <variable> in <factor>

Reason:
These production cause non-determinism because <idnest> and <variable> first set will contain id and empty string, and both are the first non-terminals in one of factors RHS.

OLD:
1. <variable>      →  <idnest> id <indiceLst>
2. <idnest>        → id <indiceLst> . <idnest> | ε
3. <factor>        → <variable>
                     | <idnest> id (<aParams>)
                     | num
                     | ( <arithExpr> )
                     | not <factor>
                     | <sign> <factor>


CHG:
1. <variable>       → id <variableIndice>
2. <idnest>         → . id <indice> <idnest> | ε
3. <factor>         → id <factorVarOrFunc> // Replace the two first rules with this
                      |  num

       | ( &lt;arithExpr&gt; )
       | not &lt;factor&gt;
       | &lt;sign&gt; &lt;factor&gt;
4. &lt;factorVarOrFunc&gt; → &lt;factorVarArray&gt; | &lt;varOrFuncIdNest&gt; | ε
5. &lt;factorVarArray&gt; →  &lt;indice&gt; &lt;indiceLst&gt; &lt;factorVarArrayNestId&gt;
6. &lt;factorVarArrayNestId&gt; → . id &lt;factorVarOrFunc&gt; | ε
7. &lt;varOrFuncIdNest&gt; →  ( &lt;aParams&gt; ) | . id &lt;factorVarOrFunc&gt;

More explanation:
&lt;factorVarOrFunc&gt; is used to determine if the statement is going to be a function call, or just a variable access. The reason we need to separate them is because a function call can not be accessed via an array access, ex: id.id[integer]() is not allowed.


## Change terminal num in &lt;factor&gt; to non-terminal and add a new &lt;num&gt; non-terminal

Reason: My scanner does not have a num token, it has only integer or float so &lt;num&gt; will contain &lt;float&gt; or &lt;integer&gt;

OLD:
  1. &lt;factor&gt;   → id &lt;factorVarOrFunc&gt; // Replace the two first rules with this
          |  num
          | ( &lt;arithExpr&gt; )
          | not &lt;factor&gt;
          | &lt;sign&gt; &lt;factor&gt;

CHG:
  1. &lt;factor&gt;   → id &lt;factorVarOrFunc&gt; // Replace the two first rules with this
          |  num
          | ( &lt;arithExpr&gt; )
          | not &lt;factor&gt;
          | &lt;sign&gt; &lt;factor&gt;
  2. &lt;num&gt;   → &lt;integer&gt; | &lt;float&gt;


## Change Allow &lt;varDecl&gt; or &lt;statementLst&gt; before or after each other in &lt;funcBody&gt;

Allow &lt;funcBody&gt; to have either &lt;varDecl&gt; and &lt;statementLst&gt; before or after each other and remove ambiguity that'll be created by it because both &lt;varDecl&gt; and &lt;statementLst&gt; have id in their first sets

Reason: A variable assignment, control flow statements, loops, and other types of statements can be used before a variable declaration, as long as the variable isn't used before its declaration but that's the purpose of the semantic analyzer to verify so the syntax analyzer should allow either to be before or after one another.

OLD:
   **1.** &lt;funcBody&gt;  → { &lt;varDecl&gt; &lt;statementLst&gt; }
   2. &lt;statement&gt;  →  &lt;assignStat&gt; ;
                     | if ( &lt;expr&gt; ) then &lt;statBlock&gt; else &lt;statBlock&gt; ;
                     | for ( &lt;type&gt; id &lt;assignOp&gt; &lt;expr&gt; ; &lt;relExpr&gt; ; &lt;assignStat&gt; )
                     &lt;statBlock&gt; ;
                     | get ( &lt;variable&gt; ) ;
                     | put ( &lt;expr&gt; ) ;
                     | return ( &lt;expr&gt; ) ;

New:
   1. &lt;funcBody&gt;    → { &lt;funcInBody&gt; }
   2. &lt;funcInBodyLst&gt; →  &lt;funcInBody&gt; &lt;funcInBodyLst&gt; | ε
   3. &lt;funcInBody&gt; → id &lt;varOrStat&gt;
                     | &lt;statementRes&gt;
                     | &lt;numType&gt; id &lt;arraySize&gt; ; // Add for int or float variable declarations

   4. &lt;varOrStat&gt;   → id &lt;arraySize&gt; ; | &lt;indiceLst&gt; &lt;idnest&gt; &lt;assignOp&gt;
   5. &lt;statement&gt;   → &lt;assignStat&gt; ; | &lt;statementRes&gt;
   6. &lt;statementRes&gt; → if ( &lt;expr&gt; ) then &lt;statBlock&gt; else &lt;statBlock&gt;;
                     | for ( &lt;type&gt; id &lt;assignOp&gt; &lt;expr&gt; ; &lt;relExpr&gt; ; &lt;assignStat&gt; )
                     &lt;statBlock&gt; ;
                     | get ( &lt;variable&gt; ) ;
                     | put ( &lt;expr&gt; ) ;
                     | return ( &lt;expr&gt; ) ;

More Explanation:

&lt;varOrStat&gt; is added because a variable declaration or assignment operation both can begin with an id, so we will need to look further ahead to see what it is.
&lt;statement&gt; was causing ambiguity, because it contains &lt;assignStat&gt; so it was removed from &lt;funcInBody&gt; and &lt;funcInBody&gt; used the new &lt;statementRes&gt; to use the reserved words statements, &lt;statementRes&gt; was also added to &lt;statement&gt; so whatever rules that were using &lt;statement&gt; would still have access to the reserved word statements.

&lt;varDecl&gt; is no longer used anywhere so it is removed

## Change <arithExpr> and <relExpr> in <expr> both derive <term>

Reason**:** These production cause non-determinism because <arithExpr> and <relExpr> both derive <term> so they will have the same first set

OLD:
1. <expr>          →  <arithExpr> | <relExpr>


CHG
1. <expr>         →  <term> <relOrAri>
2. <relOrAri>     → <relOp> <arithExpr> | <arithExprD>


## Change <statementRes> RHS if production to not use <statBlock>

Reason: An if statement doesn't always need an else and when it doesn't have an else the <statBlock> must have a semicolon after it.

OLD:
1. <statementRes>        →  if(<expr>)then<statBlock>else<statBlock>;


CHG
1. <statementRes>        →  if(<expr>)then<statBlock>else<statBlock>;
2. <statBlock>           → { <statementLst> } ; | <statement> | EPSILON
3. <statThenBlock>       → { <statementLst> } <statIfElseBlock> | <statement>
   <statElseBlock>
4. <statIfElseBlock>     → else <statBlock>  | ;
5. <statElseBlock>       → else <statBlock> | EPSILON


# Final Grammar

Please note that my tokens are a bit different from the terminals displayed in the grammar, I find that the terminals given are easier to process so I will leave them as is but a mapping of the tokens to terminals will be found at the end of this report. The syntax parser will use the tokens, but output the derivations in the terminals below.

| | |
|---|---|
| <prog> | → <classDeclLst> <progBody> |
| <classDeclLst> | → <classDecl> <classDeclLst> \| ε |
| <classDecl> | → class id { <classBody> } ; |
| <classBody> | → <classInDecl> <classBody> \| ε |
| <classInDecl> | → <type> id <postTypeId> ; |
| <postTypeId> | → <arraySize> \| ( <fParams> ) <funcBody> |
| <progBody> | → program <funcBody> ; <funcDefLst> |
| <funcHead> | → <type> id ( <fParams> ) |
| <funcDefLst> | → <funcDef> <funcDefLst> \| ε |
| <funcDef> | → <funcHead> <funcBody> ; |
| <funcBody> | → { <funcInBodyLst> } |
| <funcInBodyLst> | → <funcInBody> <funcInBodyLst> \| ε |
| <funcInBody> | → id <varOrStat> \| <statementRes> \| <numType> id <arraySize> ; |
| <varOrStat> | → id <arraySize> ; \| <indiceLst> <idnest> <assignOp> <expr> ; |
| <statementLst> | → <statement> <statementLst> \| ε |
| <statement> | → <assignStat> ; \| <statementRes> |
| <statementRes> | → if ( <expr> ) then <statThenBlock> |
| | \| for ( <type> id <assignOp> <expr> ; <relExpr> ; <assignStat> ) <statBlock> |
| | \| get ( <variable> ) ; |
| | \| put ( <expr> ) ; |
| | \| return ( <expr> ) ; |
| <statThenBlock> | → { <statementLst> } <statIfElseBlock> \| <statement> <statElseBlock> |
| <statIfElseBlock> | → else <statBlock> \| ; |
| <statElseBlock> | → else <statBlock> \| ε |
| <assignStat> | → <variable> <assignOp> <expr> |
| <statBlock> | → { <statementLst> } \| <statement> \| ε |
| <expr> | → <term> <relOrAri> |
| <relOrAri> | → <relOp> <arithExpr> \| <arithExpr'> |
| <relExpr> | → <arithExpr> <relOp> <arithExpr> |
| <arithExpr> | → <term> <arithExpr'> |
| <arithExpr'> | → <addOp> <term> <arithExpr'> \| ε |
| <sign> | → + \| - |
| <term> | → <factor> <term'> |
| <term'> | → <multOp> <factor> <term'> \| ε |
| <factor> | → id <factorVarOrFunc> \| <num> \| ( <arithExpr> ) \| not <factor> \| <sign> <factor> |
| <factorVarOrFunc> | → <factorVarArray> \| <varOrFuncIdNest> \| ε |
| <factorVarArray> | → <indice> <indiceLst> <factorVarArrayNestId> |
| <factorVarArrayNestId> | → . id <factorVarOrFunc> \| ε |
| <varOrFuncIdNest> | → ( <aParams> ) \| . id <factorVarOrFunc> |
| <variable> | → id <variableIndice> |

| | |
|---|---|
| &lt;variableIndice&gt; | → &lt;indiceLst&gt; &lt;idnest&gt; |
| &lt;idnest&gt; | → . id &lt;indiceLst&gt; &lt;idnest&gt; \| ε |
| &lt;indice&gt; | → [ &lt;arithExpr&gt; ] |
| &lt;indiceLst&gt; | → &lt;indice&gt; &lt;indiceLst&gt; \| ε |
| &lt;arraySize&gt; | → [ integer ] &lt;arraySize&gt; \| ε |
| &lt;type&gt; | → &lt;numType&gt; \| id |
| &lt;numType&gt; | → int \| float |
| &lt;fParams&gt; | → &lt;type&gt; id &lt;arraySize&gt; &lt;fParamsTail&gt; \| ε |
| &lt;aParams&gt; | → &lt;expr&gt; &lt;aParamsTail&gt; \| ε |
| &lt;fParamsTail&gt; | → , &lt;type&gt; id &lt;arraySize&gt; &lt;fParamsTail&gt; \| ε |
| &lt;aParamsTail&gt; | → , &lt;expr&gt; &lt;aParamsTail&gt; \| ε |
| &lt;assignOp&gt; | → = |
| &lt;relOp&gt; | → == \| <> \| < \| > \| <= \| >= |
| &lt;addOp&gt; | → + \| - \| or |
| &lt;multOp&gt; | → * \| / \| and |
| &lt;num&gt; | → integer \| float |

**Operators and additional lexical conventions**

| | |
|---|---|
| &lt;assignOp&gt; | → = |
| &lt;relOp&gt; | → == \| <> \| < \| > \| <= \| >= |
| &lt;addOp&gt; | → + \| - \| or |
| &lt;multOp&gt; | → * \| / \| and |

# First and follow sets

| Non-Terminals | First Set | Follow Set |
|---|---|---|
| &lt;prog&gt; | {class, program} | { $ } |
| &lt;classDeclLst&gt; | {class, ε} | { program } |
| &lt;classDecl&gt; | {class} | {class, program} |
| &lt;classBody&gt; | {int, id, float, ε} | { } } |
| &lt;classInDecl&gt; | {int, id, float} | { } , int, id, float} |
| &lt;postTypeId&gt; | { [ , (, ε} | { ; } |

| | | |
|---|---|---|
| <progBody> | {program} | { $ } |
| <funcHead> | {int, id, float} | { { } |
| <funcDefLst> | {int, id, float, ε} | { $ } |
| <funcDef> | {int, id, float} | { id, int, float, $ } |
| <funcBody> | { { } | { ; } |
| <funcInBodyLst> | {id, int, float, if, for, get, put, return, ε} | { } } |
| <funcInBody> | {id, int, float, if, for, get, put, return} | { id, int, float, if, for, get, put, return, } } |
| <varOrStat> | { id, [, =, . } | { id, if, for, get, put, return, int, float, } } |
| <statementLst> | {id, if, for, get, put, return, ε} | { } } |
| <statement> | {id, if, for, get, put, return} | { id, if, for, get, put, return, else, ;, } } |
| <statementRes> | {if, for, get, put, return} | { id, if, for, get, put, return, else, ;, }, int, float } |
| <assignStat> | {id} | { ), ; ) |
| <statBlock> | { { , id, if, for, get, put, return, ε } | { ;, id, if, for, get, put, return, else, }} |
| <statThenBlock> | { { , id, if, for, get, put, return, ε } | { else } |
| <statThenBlock> | { else, ; } | |
| <expr> | {id, integer, float, (, not, +, - } | { , , ), ; } |
| <relOrAri> | { ==, <>, <, >, <=, >=, +, -, or, ε} | { , , ), ; } |
| <relExpr> | {(, id, not, float, integer, +, -} | { ; } |
| <arithExpr> | {(, id, not, float, integer, +, -} | { ], ), ==, <>, <, >, <=, >=, ;, , } |
| <arithExpr'> | {+, -, or, ε} | { ], ), ==, <>, <, >, <=, |

| | | >=, ;, , } |
|---|---|---|
| <sign> | {+ , -} | { id, (, not, integer, float, +, - } |
| <term> | {id, integer, float, (, not, +, - } | { +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <termD> | {*, /, and, ε} | { +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <factor> | {id, integer, float, (, not, +, - } | { *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <factorVarOrFunc> | { [, ., ( ε } | { *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <factorVarArray> | { [ } | { *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <factorVarArrayNestId> | { . , ε } | { *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <varOrFuncIdNest> | { (, . , ε } | { *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <variable> | {id} | { =, ) } |
| <variableIndice> | { [, ., ε } | { =, ) } |
| <idnest> | { ., ε} | { =, ) } |
| <indice> | { [} | { [, ., =, *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |
| <indiceLst> | { [, ε } | {., =, ), *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ;, ,} |
| <arraySize> | { [, ε } | { ,, ;, ) } |
| <type> | {int, float, id} | { id } |

| | | |
|---|---|---|
| <numType> | {int, float} | { id } |
| <fParams> | { int, float, id, ε} | { ) } |
| <aParams> | {id, integer, float, (, not, +, -, ε } | { ) } |
| <fParamsTail> | { , , ε} | { ) } |
| <aParamsTail> | { , , ε} | { ) } |
| <assignOp> | { = } | { id, (, not, integer, float, +, -, if, for, get, put, return, int, } } |
| <relOp> | { ==, <>, <, >, <=, >= } | { id, (, not, integer, float, +, - } |
| <addOp> | { +, -, or } | { id, (, not, integer, float, +, - } |
| <multOp> | { *, /, and} | { id, (, not, integer, float, +, - } |
| <num> | {integer, float} | { *, /, and, +, -, or, ==, <>, <, >, <=, >=, ], ), ;, , } |

## Token Mapping with terminals

| Tokens | Lexemes |
|---|---|
| ID | letter alphanum* |
| INUM | nonzero digit* \| 0 |
| FNUM | integer fraction |
| SCMT | // alphanum* \n |
| CMT | /* (alphanum \| \n)* */ |
| EQUAL | = |

| | |
|---:|---|
| EQUIV | == |
| NOTEQ | <> |
| GTEQ | >= |
| GT | > |
| LTEQ | <= |
| LT | < |
| ADD | + |
| SUB | - |
| MULTI | * |
| DASH | / |
| OPENPARA | ( |
| CLOSEPARA | ) |
| OPENBRA | [ |
| CLOSEBRA | ] |
| OPENCURL | { |
| CLOSECURL | } |
| INT | int |
| FLOAT | float |
| AND | and |
| NOT | not |
| OR | or |
| IF | if |
| THEN | then |
| ELSE | else |

| | |
|---:|:---|
| FOR | for |
| CLASS | class |
| GET | get |
| PUT | put |
| RETURN | return |

# Error Recovery

The syntax parser error recovery/reporter will indicate that an error has occurred when an unexpected token is received. At the beginning of a production procedure, skip_errors is called to check if the lookahead value matches a value in the first and follow set, but the follow set is used only if the first set contains and epsilon. The skip_errors will continue to skip over tokens until it finds a token that's valid. Besides skip_errors, match will also report errors when an unexpected token is received. Match does not always recover from the error though and will stop the parsing from continuing, I added error handling for match in a few places though to check if there is a missing semicolon or a missing token, but this wasn't done everywhere as there's simply too many cases to cover.

# Program Overview

The syntax parser was written in C++11, the main component is the SyntaxParser class, this class can be instantiated and is a recursive descent top down parser, this class contains all the procedures for the non terminals in the grammar.

# Tools, Libraries and Techniques used

## Tools

The following tools were used to help develop the lexical analyzer.

CLion is a cross platform IDE for c++, I used this instead of Visual Studio or Xcode as I work both on a Mac and Windows so I needed a cross platform solution.
https://www.jetbrains.com/clion/

Mingw allowed me to use gcc, a C++ compiler on Windows, this helped with my cross platform solution as gcc is a cross platform compiler for C++.
http://mingw.org/

CMake let me specify how to compile my program and the multiple files required to compile it.
https://cmake.org/

Atocc kFG program was used to help removed ambiguities in the language
http://www.atocc.de/cgi-bin/atocc/site.cgi?lang=de&site=main

Libraries

Google test was used to help test the lexical analyzer, I chose it because it's one of the most popular c++ unit testing frameworks and was easy to setup and use.
https://github.com/google/googletest

Various C++ standard libraries were used to help develop the program, the following includes were used in my program
- #include <iostream>
- #include <vector>
- #include <map>
- #include <fstream>
- #include <algorithm>
- #include <string>


# Techniques

The technique used was the Recursive descent predictive parser.