

Programming Assignment 2

LALR Parsing and Building ASTs

The goal of this assignment is to understand how to use LALR parsing for more complex grammars with expressions, and how to construct *Abstract Syntax Trees* (ASTs) during parsing. You will:

- Write a more complex context-free grammar and implement it using an LALR parser.
- Resolve ambiguities and experiment with LALR conflicts.
- Build an AST in the semantic actions of the parser.

In addition to JFlex, Beaver, and other tools you used in the previous assignment, you will use the *JastAdd* tool to define AST classes. For documentation of JastAdd, see <https://jastadd.cs.lth.se>.

You will define a language called *SimpliC* (Simplistic C), and generate a compiler that parses SimpliC into ASTs. To your help, there is a demonstration example *CalcAST* showing how to build ASTs for the Calc language.

Try to solve all parts of this assignment before going to the lab session. If you get stuck, use the forum to ask for help. If this does not help, you will have to ask at the lab session. Make notes about answers to questions in the tasks, and about things you would like to discuss at the lab session.

- Major tasks are marked with a black triangle, like this.

1 The CalcAST Demo

The CalcAST demonstration example implements an extended version of the *Calc* language from Assignment 1, introducing some new features that will illustrate the use of lists and optionals. The language now allows a list of variable bindings in the `let` expression, and a new `ask user` expression with an optional default value has been added.

Here is an example of a Calc program:

```
let
  radius = ask user [1.0]
  pi = 3.14
in
  2.0 * pi * radius
end
```

The idea is that when this program is run, the `ask user` expression will ask the user to input a value, using the default value of 1.0. An `ask user` expression without the optional default value would be written simply as `ask user`, without the square bracket clause.¹

This example program is named `example2.in` and can be found in the `testfiles/ast` directory.

- Download the CalcAST example. Run the tests and see that they pass.

¹Note, however, that we will not actually run Calc or SimpliC programs until Assignment 5 and 6.

1.1 The abstract grammar

The *Context Free Grammar* (CFG) for the new Calc language is

```
program → exp <EOF>
exp     → factor | ( exp '*' factor )
factor  → let | numeral | id | ask
let     → 'let' binding+ 'in' exp 'end'
binding → id '=' exp
numeral → <NUMERAL>
id      → <ID>
ask     → 'ask' 'user' [ '[' exp ']' ]
```

To represent an AST in Java we will use a class hierarchy generated by JastAdd, and defined in an *abstract grammar*, in this case `calc.ast`. The abstract grammar is similar to the CFG:

```
Program ::= Expr;

abstract Expr;
Mul: Expr ::= Left:Expr Right:Expr;
Numeral: Expr ::= <NUMERAL:String>;
IdUse: Expr ::= <ID:String>;

Let: Expr ::= Binding* Expr;
Binding ::= IdDecl Expr;
IdDecl ::= <ID:String>;

Ask: Expr ::= [Default:Expr];
```

Here, an abstract grammar rule

```
A: B ::= children;
```

corresponds to a Java class A which extends the superclass B and has a number of *children* (composition relationships in UML). The children can be of the following kinds:

```
C      A child of type C
C*     A list of children of type C
[C]    An optional child of type C
<T:S>  A token named T of type S
```

Children may be explicitly named, for example:

```
N:C    A child named N of type C
```

and if no explicit name is given, the type name is used for the child.

1.2 Designing the abstract grammar

In comparing the abstract grammar and the CFG, both similarities and differences can be noted.

A nonterminal with several alternatives, like `exp`, can be modelled as an abstract class, like `Expr` in the abstract grammar. The alternatives on the right-hand side in the CFG have been modelled using subclassing in the abstract grammar. This is a common pattern when using object-orientation for representing ASTs: replace alternatives with subclassing.

Note, however, that the nonterminal `factor` in the CFG does not appear at all in the abstract grammar: this nonterminal is only used for disambiguating the CFG so that we can parse text into a tree in

an unambiguous way. Once we have the tree, the `factor` node is uninteresting. For this reason, all expressions (`Mul`, `Numeral`, `IdUse`) are direct subclasses of `Expr`.

Another difference is that the nonterminal `id` in the CFG has been replaced with two classes, `IdUse` and `IdDecl`, in the abstract grammar. This separation is useful when implementing name analysis, where name uses are bound to name declarations. We will implement name analysis later on, in another assignment.

A third difference is that the "one or more" repetition in the CFG (+) is replaced in the abstract grammar by a "zero or more" (*). It is rarely useful to make any distinction between these cases in the semantic analysis, and JastAdd does in fact only support * and not +. To check that there is at least one element can be done easily by the parsing rules. It can also be done by adding an attribute in the semantic analysis (see Assignment 4). In particular, make sure to not replace `X+` in the CFG by `X X*` in the abstract grammar, as this will only lead to an unnecessarily complex API for the semantic analysis.

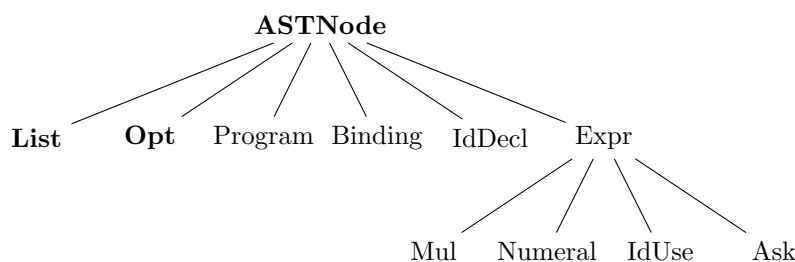
- To sum up, when designing the abstract grammar, it is important to think about what data structure will be useful for doing semantic analysis. Favor a simple structure. Don't include parsing-specific details in the abstract grammar. Compare the CFG and the abstract grammar for Calc. What similarities and differences do you find?

1.3 Generated AST classes

The generated classes will contain accessor methods for the children and tokens, as the following simplified API illustrates for the classes `Mul` and `IdUse`.

```
public class Mul extends Expr {
    Expr getLeft();
    Expr getRight();
    void setLeft(Expr e);
    void setRight(Expr e);
}
public class IdUse extends Expr {
    String getID();
    void setID(String id);
}
```

The Java classes generated by JastAdd will have the following class hierarchy:



`ASTNode`, `Opt` and `List` are predefined implicit AST classes. They are generated by JastAdd without any mention in the `.ast` file. `Opt` and `List` are used to represent optional and list components in the abstract syntax tree. `ASTNode` serves as the common superclass of all AST nodes.²

- Look at the abstract grammar specification (`src/jastadd/calc.ast`) and make sure you understand it. Look at the corresponding generated classes in `src/gen/lang/ast`. The classes contain lots of generated methods, but you should be able to find, for example, `Mul`'s `getLeft()` and `getRight()` methods mentioned above. Familiarize yourself briefly with the JastAdd reference manual³ so that you know where to look for details later.

²`ASTNode`, `Opt` and `List` are generated into the package `lang.ast`. Note that this generated `List` class is different from the standard class `java.util.List`.

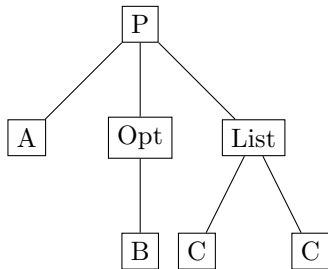
³<https://jastadd.cs.lth.se/web/documentation/reference-manual.php>

1.4 Well-formed ASTs

When building a JastAdd AST, care must be taken to build it so that it is *well-formed*. Suppose we have an abstract grammar class $P ::= A [B] C^*$. A P node is well-formed if

- an ordinary child, like A , is represented by a node of type `A`.
- an optional child, like $[B]$, is represented by a node of type `Opt` which has zero or one children of type `B`.
- a list child, like C^* , is represented by a node of type `List` which has zero or more children of type `C`.
- Each node must be a real object, i.e., not `null`.
- The node objects must form a tree, i.e., it is not allowed to reuse a node in several places in the same AST, or in different ASTs.

The example below shows a well-formed AST for $P ::= A [B] C^*$



These conditions are not checked when generating the parser, and several of them are not even checked when the parser code is compiled. So take care in building the AST correctly, as you might otherwise run into runtime errors and problems that are difficult to debug.

- What would be examples of other well-formed ASTs for the $P ::= A [B] C^*$ class? What would be examples of ill-formed ASTs?

1.5 Building the AST

The abstract syntax tree is built by the parser, as specified in the `.beaver` file, using so called *semantic actions* of the production rules. We have to instruct Beaver what kind of JastAdd AST node is built by the productions for each Beaver nonterminal. This will become the return type of the semantic action in the generated parser. We specify this using Beaver's `%typeof` directive, as the following code fragment illustrates.

```
%typeof program = "Program";
%typeof exp = "Expr";
%typeof factor = "Expr";
%typeof let = "Let";
...
```

Note that to the left in each clause, we have a nonterminal in the Beaver CFG, like `program`. To the right, we have a class from the JastAdd abstract grammar, like `Program`.

1.5.1 Semantic Actions

To build AST nodes, the productions in the parser specification need to be extended with semantic actions. Semantic actions are code snippets that are run during parsing. Consider the production for the `program` nonterminal:

```
program = exp.a {: return new Program(a); :} ;
```

Note the code inside the `{:` and `}`. This is the semantic action code – ordinary Java code for constructing the AST node representing the parsed production. In this case, the semantic action simply creates a new node of type `Program` with an expression as its child. Note that `program` is a nonterminal in the parser specification, whereas `Program` is an AST node type, specified in the abstract grammar. The expression inside the semantic action must have the same type as specified in the `%typeof` directive.

In Beaver, a right-hand side component can be *named* using dot notation: in the example above, the nonterminal `exp` has been named `a`, by writing `exp.a`. The name `a` refers to the AST subtree constructed by the `exp` nonterminal, and can be used inside the semantic action to refer to that subtree. In this case, `a` is used in the semantic action as an argument to the `Program` constructor.

1.5.2 Optional Components

For most productions, building the AST is straightforward as explained in the previous section. But if there are optional or list components in a production, an extra `Opt` or `List` node must be built.

Although Beaver does have some support for optionals and lists, it is easier to build the AST if we first rewrite the grammar to eliminate optionals and lists, using explicit alternatives. Recall the new `Ask` construct which has an optional default expression. The context-free rule is

$$\text{ask} \rightarrow \text{'ask' 'user' ['[' exp '']]}$$

and the abstract grammar rule is

```
Ask: Expr ::= [Default:Expr];
```

To implement this in Beaver, we first rewrite the context-free rule to eliminate the optional notation. We can do this by introducing two alternatives for `ask`: one without and one with the default expression:

$$\text{ask} \rightarrow \text{'ask' 'user' | 'ask' 'user' exp}$$

The Beaver specification can now be written as follows:

```
ask =
  ASK USER                               {: return new Ask(new Opt()); :}
  | ASK USER LBRACKET exp.a RBRACKET     {: return new Ask(new Opt(a)); :}
  ;
```

Note that the semantic action always builds an `Opt` node to match the optional expression, and that different `Opt` constructors are used depending on if the optional expression is present or missing.

1.5.3 List Components

In order to handle list components (repetition), we use a similar technique as for optionals, rewriting the parsing grammar to use explicit alternatives. Here, we need also to introduce a new nonterminal to replace repetition with recursion.

Recall the abstract grammar rules for bindings in let-expressions:

```
Let: Expr ::= Binding* Expr;
Binding ::= IdDecl Expr;
```

To implement this in Beaver, we first rewrite the context-free grammar rule for `Let`:

$$\text{let} \rightarrow \text{binding}^+$$

replacing the repetition with a new nonterminal `binding_list`:

$$\text{binding_list} \rightarrow \text{binding} \mid \text{binding_list binding}$$

The Beaver specification for `binding_list` can be written as follows:

```
%typeof binding_list = "List";
...
binding_list =
    binding.a          {: return new List().add(a); :}
  | binding_list.a binding.b {: return a.add(b); :}
;
```

Note bene! Don't forget to declare the type of the new nonterminal, using a `%typeof` directive, as above.

- Look at the parser specification (`src/parser/parser.beaver`) and make sure you understand how the AST is built. Where are the semantic actions? How are optionals and lists handled?

1.6 The DumpTree Aspect

The CalcAST example contains code for printing an AST, using indentation to show the tree structure. This is useful when writing test cases and for debugging, to check that the AST is well-formed and correct.

The code is located in an “aspect file” `DumpTree.jrag` in the `src/jastadd` directory in the CalcAST project. A JastAdd *aspect* may contain methods of AST classes, but specified in a separate file. The JastAdd system *weaves* those methods into the generated AST classes. The file `DumpTree.jrag` contains a single aspect declaration:

```
aspect DumpTree {
    public String ASTNode.dumpTree() {
        ByteArrayOutputStream bytes = new ByteArrayOutputStream();
        dumpTree(new PrintStream(bytes));
        return bytes.toString();
    }

    public void ASTNode.dumpTree(PrintStream out) {
        [...]
    }

    [...]
}
```

The methods in the aspect have an extra qualifier before the method name, e.g., `ASTNode.dumpTree`. The JastAdd code generator will weave the method into the generated `ASTNode` class. This way of writing a declaration in an aspect instead of in the class in which it belongs is called an *inter-type declaration*. Later in the course, we will look at other kinds of inter-type declarations you can write in an aspect.

The `compiler.jar` main class is `Compiler.java` and it will parse in a calc program and then print it using the `dumpTree` method.

- ▶ Try out `compiler.jar` on one of the test programs. For example:

```
./gradlew jar
java -jar compiler.jar testfiles/ast/example2.in
```

1.7 Testing framework

The CalcAST example uses a testing framework that is more advanced than the one used in Assignment 1. The new test framework uses a JUnit feature called *parameterized tests*. With parameterized tests you don't need to write a test method for each test case: it is sufficient to add an input and expected output file to the `testfiles/ast` directory.

In CalcAST, this testing technique is used for the AST tests to check that correct ASTs are built when parsing different test programs. The framework starts one test per `.in`-file in the `testfiles/ast` directory. Each test file is parsed and the tree is printed to a corresponding `.out` file. If there exists a matching `.expected` file then the contents of the output is checked against the expected output.

When adding new tests it is usually too tedious to hand-write each `.expected` file. Instead, you can add just the input file, run the tests to get the `.out` file, and then rename the `.out` file to use the `.expected` extension. Of course, it is important to inspect the `.out` file before renaming it, to make sure it really is correct!

- ▶ Look at the testcases in the `testfiles/ast` directory. The `.expected` files contain the expected output after parsing the corresponding `.in`-file and printing the AST using the `dumpTree` method. Make sure you understand what AST an `.expected` file corresponds to.
- ▶ Add a new AST test case in the CalcAST project.

1.8 Build Script

The CalcAST Gradle script (`build.gradle`) contains some more elements and tasks compared to the script we used in Assignment 1. To be able to run JastAdd, it uses the `org.jastadd` plugin and a new `generateJava` task:

```
compileJava.dependsOn 'generateJava', 'generateScanner', 'generateParser'
task generateJava(type: org.jastadd.JastAddTask) {
    description 'Generates the compiler AST classes from JastAdd specifications.'
    ...
}
```

2 Troubleshooting

Many things can go wrong when you generate a parser. Your grammar might be non-LALR in which case the parser generator will report conflicts. Some errors in the `.beaver` specification are not discovered until you compile the generated Java code, or until you run the generated parser. It is very likely that you will run into problems like these, and it will be much easier to deal with them if you have encountered them before in a controlled setting. So in this section, we will experiment with the CalcAST example to provoke some common errors so that you understand them better.⁴

2.1 Provoke an LALR conflict in the parsing specification

- ▶ Add a conflict to the CalcAST parser specification (`parser.beaver`) by changing `factor.b` to `exp.b` in the rule

```
exp =
    factor
    | exp.a MUL factor.b {: return new Mul(a, b); :}
    ;
```

After this change the grammar is no longer LALR(1), and there will be a conflict in the parsing table. In this case, it is a shift-reduce conflict.

- ▶ Rerun the build (`./gradlew test`) to see how Beaver handles this problem. Remember to “refresh” your Eclipse project (if using Eclipse).⁵

Note that the build goes through and the tests are run, even if there is a conflict. Scroll back to look at the Beaver output. Note that Beaver gives a *warning*, rather than an error. This is because it by default resolves shift-reduce conflicts by selecting shift. It goes on and generates a parser that will shift in the conflicting situation. In this case, shift is not the right way to resolve the conflict, and because we happen to have a good test suite, several tests fail.

Important advice. When you get a conflict, try to figure out what the problem is, and change the grammar to resolve it. Even if Beaver happens to resolve the conflict in the right way, make sure to change the grammar so that you don’t get any warnings and default resolutions by Beaver. Otherwise, it will be very difficult for you to sort out new conflicts later on.

To find out what the problem is, study the error message given by Beaver, study any test cases that fail, and think about your last change to the grammar. If you can construct an example program that illustrates the conflict, this will help you understand the problem.

In this particular case, the error message is:

```
WARNING: resolved SHIFT/REDUCE conflict on [MUL] by selecting SHIFT:
  REDUCE exp = exp MUL exp
  SHIFT MUL
Context:
  exp = exp . MUL exp [END]
  exp = exp MUL exp . [MUL]
  exp = exp . MUL exp [EOF]
  exp = exp . MUL exp [MUL]
  exp = exp . MUL exp [RBRACKET]
  exp = exp . MUL exp [IN]
  exp = exp . MUL exp [ID]
```

This means that if a `MUL` token appears after parsing the right hand side of

⁴The examples in this section are tested on NeoBeaver 0.1.5. If you are using another version, error handling might differ.

⁵When you run unix commands, Eclipse does not automatically know that files have been added, changed, or removed. When you do a “refresh” in Eclipse, it will scan the files in the project directory to look for changes.


```
exp = exp.a MUL exp.b
```

both shift and reduce would fit the grammar, and Beaver will choose to shift.

We can try to understand what happens by constructing an example corresponding to that situation:

```
x * y * ...
```

The conflict occurs after `y`, when reading the second `MUL`. We can either reduce `x * y`, or shift the `MUL`. If we reduce `x * y`, a left-recursive tree will be built. If we shift the `MUL`, we will later have to reduce `y * ...`, and a right-recursive tree will be built.

- ▶ Undo the change and make sure all tests pass again.

2.2 Unintended empty productions

A common mistake is to use an extra `|` in a rule, which results in an extra empty production. This can easily happen when copying and pasting lines with production rules. As an example, consider the following rule:

```
binding_list =
  binding.a      {: return new List().add(a); :}
  | binding_list.a binding.b {: return a.add(b); :}
  ;
```

When writing this rule, we might by mistake let the first production start with `|` as follows:

```
binding_list =
  | binding.a      {: return new List().add(a); :}
  | binding_list.a binding.b {: return a.add(b); :}
  ;
```

But this means that before the first `|`, there is actually an empty production. If we do want an empty production, it is good practice to add a comment to make it more prominent, like this:

```
binding_list =
  /* epsilon */
  | binding.a      {: return new List().add(a); :}
  | binding_list.a binding.b {: return a.add(b); :}
  ;
```

But in this case, we don't want the empty production, and running the grammar above through Beaver gives a shift-reduce conflict with the following message:

```
WARNING: resolved SHIFT/REDUCE conflict on [ID] by selecting SHIFT:
  REDUCE binding_list =
  SHIFT ID
```

We see that it has found a shift-reduce conflict that would occur after parsing the right-hand side of an empty production for `binding_list`.

- ▶ Experiment with this problem by adding an extra `|` to the `binding_list` production in the `.beaver` file to get an empty production. Rebuild, and confirm that there is now a shift-reduce conflict. Do all the tests still pass? If so, why?
- ▶ Undo the change and make sure all tests pass again.

2.3 Errors in the generated Java code

The code you write in Beaver’s semantic actions is Java code. But Beaver does not analyze it, so you will not see compile-time errors in it until you compile the generated code. It can be a bit tricky to understand how compile-time or runtime errors correspond to the code you originally wrote in the semantic actions.

Try a couple of common mistakes:

- ▶ Remove one of the `%typeof` directives in the `.beaver` file. Note that the parser is generated without any errors, but when you compile the generated code, the Java compiler will complain about both incompatible types and missing symbols (declarations).
- ▶ In one of the semantic actions, remove the `return` keyword. Note that the parser is generated without any errors. Depending on what version of Beaver/NeoBeaver you are running, you might get a compile-time error when you run the Java compiler, or a run-time error when you run the parser.

2.4 Cast exception when running the generated parser

If you leave out a semantic action in the parser specification, a default semantic action will be generated. The default action will return the result of parsing the last symbol in the production (or the last variable, in case the rule contains variables). This is sometimes what you want, but not always. Consider the following example:

```
factor =
    let
    | ask
    | numeral
    | id_use
    ;

let = LET binding_list IN exp END;
```

For `factor`, the default behavior is what we want. However, for the `let` rule, the last symbol is the `END` token. This is, of course, completely wrong—what we would like to return is a new `Let` subtree. Because of the way Beaver is implemented, the error will not show at compile-time, but instead at runtime, when running the parser. The generated parser will at some point cast an object to what it thinks is an expression, but since it is actually an `END` token, a class cast exception is thrown.

- ▶ Replace the `let` rule in the Beaver specification by the one shown above (without semantic action or variables). Confirm that you get a class cast exception when you run the compiler on `testfiles/ast/let.in`.

2.5 Add print statements

As always, inserting print statements is a primitive but useful way to help with debugging. The semantic actions are just Java code inserted into the generated parser at different places. If you add more code than building ASTs, it will of course be run too.

- ▶ As an experiment, add some print statement, for example `System.out.println("Hello")`, in one of the semantic actions. Generate the parser and run the tests.

3 The SimpliC Language

Your main task in this assignment is to implement a parser for a small C-like language named SimpliC. You should do this by specifying a JastAdd abstract grammar, a Beaver parser, and a JFlex scanner. The parser should produce a well-formed abstract syntax tree which follows the abstract grammar. Furthermore, there should be automated tests showing that the parser produces the correct AST for a number of examples. There should also be some negative tests, showing that the parser detects parsing errors in erroneous input files.

3.1 Informal specification

A SimpliC program consists of a list of function declarations which all return an integer result (although the return statement can be omitted). There are also two built-in functions: `read` which reads an integer from standard in, and `print` which prints its argument to standard out. These are not declared in the program, but can be called from the program.

There are no global variables, only local variables and function parameters. All variables and parameters are integer typed. Local variables may optionally be given an initial value. Statements supported include assignments, `return` statements, function call statements, `if`, and `while` statements. Expressions supported include identifier uses, function calls, and common binary arithmetic expressions like `+`, `-`, `*`, `/`, `%`, the comparison expressions `==`, `!=`, `<=`, `<`, `>=`, and `>`, and unary minus. Binary expressions should be represented as trees, where each node corresponds to an operator having its operands as left and right subtrees. The usual priority and associativity rules should apply. For example, “`<`” should have lower priority than `+` and be non-associative. The remainder operator, `%`, has the same priority and associativity as multiplication and division. Parentheses can be used in expressions. SimpliC also supports end-of-line comments.

The SimpliC language is simple enough that it would be possible to prevent some type errors already in the parser. However, you should **not** try to prevent type errors in the parser. They should instead be caught in the semantic analysis phase, which we will do in later assignments. This makes the compiler design simpler and more extensible. For example, programs like the following should be parsable, even if it is not type correct.

```
...
a = (b > c) - 10;
a = (b > c) * x;
a = x / -(b == c);
if (b) {
    a = 3;
}
```

The semantic analyzer will check the types of all subexpressions anyway, so nothing is gained by trying to prevent some of these cases in the parser. It only makes the design unnecessarily complex.

Some additional points to keep in mind:

- Function declarations may not be nested.
- It should be possible to make function calls both as expressions and as statements. *Hint!* Introduce a `FunctionCall` which is an expression, and a `FunctionCallStmt` which is a statement that has a `FunctionCall` as a child.
- Don't introduce special syntax for calls to `read` and `print`. Use the normal function call syntax. This will simplify semantic analysis later on.
- Functions and blocks are allowed to be empty.
- Make sure lists work correctly. For example, it should **not** be allowed to write an extra comma, like
`p(x, y,)`

- Variable declarations like `int x;` and `int x=3;` should be represented by the same node type, e.g., `VarDecl` with an optional expression component. (Don't use two different node types for these constructs, since that would complicate semantic analysis later on.)

3.2 SimpliC Code Samples

The example program below shows large parts of the SimpliC syntax. The `gcd1` and `gcd2` functions implement two different algorithms for finding the greatest common divisor (GCD) of two positive integers. The `main` function tests both algorithms with inputs provided by the user:

```
int gcd1(int a, int b) {
    while (a != b) {
        if (a > b) {
            a = a - b;
        } else {
            b = b - a;
        }
    }
    return a;
}

int gcd2(int a, int b) {
    if (b == 0) {
        return a;
    }
    return gcd2(b, a % b);
}

// Test two different GCD algorithms and print the results.
// The results should be equal when both inputs are positive.
int main() {
    int a;
    int b;
    a = read();
    b = read();
    print(gcd1(a, b));
    print(gcd2(a, b));
    int diff = gcd1(a, b) - gcd2(a, b);
    print(diff);
    return 0;
}
```

3.3 Implementing the parser

In implementing the parser, you should work in small steps: first implement a parser for a small subset of the language, then add a few constructs, and so on. Many things can go wrong, so this way you can easier relate any new problems to your latest changes.

- For each of the implementation steps below, this is how you should work: Write down a test program that you want to be able to parse. Then design the abstract grammar (using the JastAdd `.ast` format), and a corresponding concrete grammar in EBNF (the latter can be on paper). Try to make the EBNF as similar as possible to the `.ast` grammar. However, it must be unambiguous, for example dealing with priorities and associativity, which might call for adding more rules. Then rewrite the EBNF grammar to eliminate optionals and lists (as discussed in sections 1.5.2 and 1.5.3), and code it up in a `.beaver` specification. Get the parser to work, and automate the test. Run all your tests to make sure you have not broken anything that worked previously. Review your implementation to make sure you use good

names for classes and components in the abstract grammar. If you get stuck, take a look at the CalcAST example for inspiration.

To your help there is an example called MinimalAST that you can download and use as a starting point. MinimalAST is similar to CalcAST, but supports a very small language with just one token and one production.

Remember to commit often to your Git repository to show your progression with the lab. By committing as soon as you have a new working version of your project, you can easily back up to an earlier stable state, in case you get into a too big mess.

3.3.1 Step 1: Simplified function declarations

- ▶ As the first step, implement a parser for programs consisting of only a list of empty function declarations without parameters. For example:

```
int a() { }
int b() { }
```

3.3.2 Step 2: Variable declarations and uses

- ▶ As a next step, extend your parser to support variable declarations and uses, and also assignments and integer literals. **Nota bene!** Make sure to distinguish between declarations and uses of identifiers in the abstract grammar, for example using the classes `IdDec1` and `IdUse`. Otherwise, you will not be able to do name analysis in later assignments. Here is an example test program:

```
int a() {
    int x;
    int y;
    x = 5;
    y = x;
}
```

3.3.3 Step 3: Expressions

- ▶ Expressions is usually the trickiest thing to get right in a parser. The abstract grammar should be simple, but the parser needs to deal with priorities and associativity. Extend your parser to support all the different kinds of expressions. First write down some simple test programs that will make it possible for you to check that you get everything right. Then try to make them work. It is a good idea to break this problem down into several subproblems. For example, start by making `+` and `*` work. Then add the other operators. See section 3.1 for what operators should be supported.

Nota bene! You should write an *unambiguous* parsing grammar (e.g., with `term` and `factor`), and *not* use parsing directives like `%left` and `%right`. This is because it is important to understand how to disambiguate a grammar – directives are not always available, and are not always applicable to ambiguities.

3.3.4 Step 4: Complete the language

- ▶ Complete the parser by adding new test cases and extending the specification. Remember to work in small steps: add one or a few language constructs at a time and make them work, before going on with the next step. Make sure you cover all the constructs discussed in sections 3.1 and 3.2. Make sure the parser generation does not give any warnings about shift-reduce conflicts that are automatically resolved. For handling the dangling-else problem, you may restrict the *then* part of an *if* statement to be a block, as discussed in lecture 4.

3.3.5 Step 5: Look over your test cases and specifications

- ▶ Look over your test cases and specifications. There should be test cases covering all parts of the syntax, and test cases covering tricky interactions between constructs (like priorities, etc.). There should also be negative test cases that check that the parser gives errors for syntactically incorrect input programs. The `gcd` program in section 3.2 should be one of the test cases. Make sure also that your abstract grammar is simple and well designed with good names for the different constructs. It is important to get everything correct and simple, as the following assignments will build on this one.

3.3.6 Step 6: Test your compiler in the "self-grader"

- ▶ To help you debug your compiler, there is a web site where you can upload your compiler and run it on a large number of test cases. See the course page for details.

3.4 Optional additions

If you would like to implement a slightly larger language, here are some optional things you could add:

- Allow multiple variable declarations like `(int a, b;)`.
- Support for `break` and `continue` statements.
- Multi-line comments.

It is recommended to not add these optional parts until you have all the mandatory parts working.

4 What to show and discuss with assignment supervisor

Before calling on your supervisor to show your assignment, check for the following things:

- Your abstract grammar should be simple, and without unnecessary parsing details. It should follow the advice in Sections 1.2 and 3.1. If your abstract grammar is unnecessarily complex or unreadable, the supervisor may ask you to simplify and improve it. (A simple clear abstract grammar will make both your work and the supervisors' work much easier in upcoming labs.)
- As explained in Section 3.3.3, your parsing grammar must not include disambiguation directives like `%left` and `%right`.

When you are ready with the assignment, these are typical things your supervisor may ask you to do:

- Show your test cases.
- Show that you can run your generated parser (`compiler.jar`) on an example
- Show that you can run all your test cases automatically on the generated parser (`./gradlew test`).
- Show that your parser generation does not result in any warnings.
- What does your `.ast`, parser, and scanner specifications look like? Is there anything you would like to ask or comment about them?