

Department of Mathematics & Computer Science
University of Southern Denmark | IMADA
Saturday 29th April, 2023

Compiler for Panda

BADM500: Bachelor Project

Author

KIAN BANKE LARSEN
kilar20@student.sdu.dk

Supervisor

KIM SKAK LARSEN
Professor



Abstract

English This is my very good abstract

Danish Et fantastisk abstract

Contents

1	Introduction	1
2	Project Basics	2
2.1	Project Structure	2
2.2	Python Lex-Yacc	3
2.3	Design Principles & Patterns	3
2.4	Compiler Usage	4
3	Phases	5
3.1	Scanner	7
3.2	Parsing	8
3.3	Symbol Collection	10
3.4	Desugaring	11
3.5	Code Generation	12
3.5.1	ILOC code	12
3.5.2	Stack Machine	12
3.5.3	Register Allocation	12
4	Testing	13
4.1	System Testing	13
4.2	Coverage	15
5	Performance Comparison	17
6	Evaluation	18
6.1	Language Considerations	18
6.2	Further Development	19
7	Conclusion	20
	References	21

Introduction

This report examines how to make a simple compiler in Python, named Panda for no particular reason. The compiler is simple in the sense that some decisions have been made to ease the process, although the decisions are not necessarily optimal. The aim is to learn different compiler techniques and get a hands-on feel for the different compiler phases by actually implementing a working compiler, targeting X86 assembler, from scratch, using a Flex/Bison equivalent package such as PLY for scanning and parsing.

The language to be compiled is a subset of the imperative language C. This has been chosen because of its simpler syntax and easy to read curly bracket enclosed static scopes. In this project, we are interested in making a language having integers, Booleans and preferably some kind of floats. The language must have control flow constructs in form of `if-else` statements and functions, and iterative constructs such as `for-` and `while`-loops.

A modern compiler is, as is well known, divided into phases. These phases relate to lexical and syntactic analysis, resulting in an abstract syntax tree. Subsequent phases analyze and adorn the abstract syntax tree, building a symbol table and finally generating assembler code.

The main focus in regard to advanced techniques will be local register allocation, using techniques described in Copper and Torczon [2022](#). Handling this efficiently requires data flow analysis via control flow-graph, construction of interference graph, graph coloring and translation back to instructions using a combination of the registers and the stack, when the available registers do not suffice.

Initially, a stack machine will be prepared, which will form the basis for developing a compiler that uses CPU registers. We will take advantage of the split phases property when replacing the stack code generation phase in benefit for one that uses register allocation. This allows us to only worry about ensuring that subsequent phases cope with the changes made in the former phases.

When adding extra complexity such as register allocation, it is important to document the benefit of this choice. Performance of the stack machine and the register machine will therefore be constructively compared.

Project Basics

This section is reserved for articulating some choices made at the very beginning of the project that defined the framework for how to develop and use the compiler.

2.1 Project Structure

The Compiler module uses different python packages – things that belong together must be together. It has been desired to make a clear division between the different phases, and this has been achieved by creating the Python package **phase**. Likewise, **dataclass** is a package that contains internal data structures, in other words, classes used to hold data. **Printer** is a package that is used primarily for debugging, but sometimes it is just nice to consider data structures graphically. **Testing** is placed on the same level as **src** because it has nothing to do with the compiler’s implementation, it is just a QA tool that makes it easier to verify correctness. **compiler.py** takes care of summarizing all functionality, but the module requires arguments from the command line, and those arguments (as well as testing) are handled in **main.py** (the project’s main file). Moreover, a **README.md** has been written as a quick-start guide.

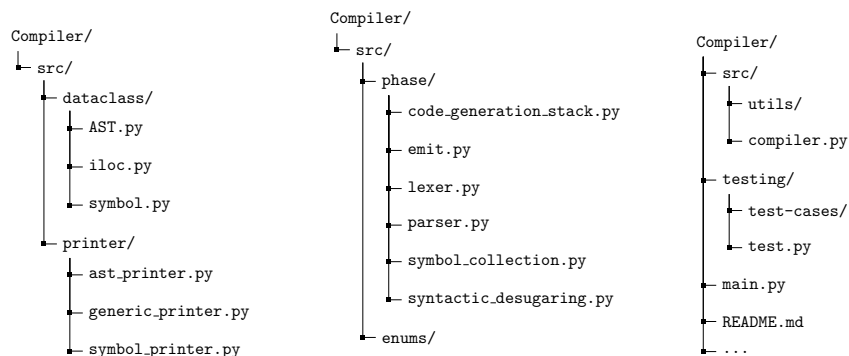


Figure 2.1: Project file tree.

2.2 Python Lex-Yacc

PLY is a native Python tool, relying on reflection, used to automatically generate scanners and LALR(1) parsers. The package is well documented at the following source: Beazley, David M. 2018. Usage of the package will be described when reviewing the compiler phases in isolation. It has been chosen to use PLY in order to reserve more time for the compiler itself, though studies show that most compilers use hand-coded scanners (Copper and Torczon 2022, p. 69). However, tool-generated parsers are more common than hand-coded parsers (*ibid.*, p. 85).

2.3 Design Principles & Patterns

When starting a new project, it is important to make some basic thoughts about the architecture. Sensible choices at the beginning can increase code readability and make maintainability easier. It is particularly important to consider design principles and design patterns, as this will have a big effect on, i.e., how data structures are traversed and code testability. In this context, design principles refer to SOLID, and design patterns refer to Gang of Four's 23 design patterns.

We will start by considering design principles. SOLID is a mnemonic acronym for a set of design principles concerning software development in object-oriented languages: **S**ingle Responsibility, **O**pen Closed, **L**iskov's Substitution, **I**nterface Segregation and **D**ependency Inversion. The principles are in many ways obvious when rehearsed, but not necessarily followed as it requires active consideration. Single responsibility is particularly expressed in the project by the sharp division of phases and their interfaces between them. Open/closed is not particularly used in the project, as inheritance cases sparsely appear, though crucial in the printer package. It will never be necessary to edit the generic printer because functionality to print a specific data structure is first implemented upon extension. Liskov's substitution principle is accommodated in the AST data class, since any subtree is a valid tree and all nodes are AST nodes. Dependency inversion principle simply means that you must program against an interface and not an implementation. There is actually an incident where the project does not live up to this principle, and that is when using the hidden method `._value2member_map_` on `Enum` in the parsing phase. Other examples apply, of course, as every class must accommodate every principle. This was just a quick review.

One of the big decisions regarding behavioral design patterns has been whether to use Visitors, just like in the well known SCIL compiler from the DM565 course. It was decided not to use visitors, because Python 3.10 comes with a new cool feature, namely match statements, which makes it possible to exploit the benefits of structural pattern matching. Although it is nice to let the data structure decide its iteration, I still prefer having everything written explicitly when learning to write a compiler. Using match statements requires one to repeat the iteration logic for every new operation, but that kind of also makes it easier to implement new logic, as one does not have to remember the visitor pattern. Another useful creational design pattern is the Singleton pattern used in `label_generator.py`. This makes it possible to retrieve the label generator in any class, while preserving the state on `count`, without having to pass the object through all the phases manually.

2.4 Compiler Usage

The main file handles the instantiation and therefore also the running of the `PandaCompiler` class. Python `argparse` is used to take command line arguments because it adds a lot of user-friendliness to the compiler. `argparse` provides the opportunity to query the compiler usage in the terminal, thus showing what options are available. Doing so yields the result stated below:

```
1 Compiler$ python3.10 main.py --help
2
3 usage: Compiler for Panda [-h] [-o OUTPUT] [-c] [-d] [-f FILE] [-t] [-r]
4
5 Compiles source code to assembly
6
7 options:
8   -h, --help            show this help message and exit
9   -o OUTPUT, --output OUTPUT
10                        Specify name of assembly output file
11   -c, --compile          Set this flag if the output file should be compiled
12                        ↪ with gcc
13   -d, --debug            Set this flag for debugging information, i.e., ILOC and
14                        ↪ Graphviz
15   -f FILE, --file FILE  Path to input file, otherwise stdin will be used
16   -t, --runTests         Run tests
17   -r, --run             Run compiled program
```

This informs the user that one can specify an input file or provide input directly in the command line. Specifying the name of the output file is optional. Furthermore, one can control whether the file should be automatically compiled with `gcc` and directly run on the system. Many of the arguments act as flags to the compiler, such as `--debug` or `--runTests`, which are flags that specify whether a certain piece of code should be executed. Much of the setup of `argparse` is omitted, in this example, but the essential part of the functionality is listed below. All pip requirements needed for running `main.py` can be installed using `pip install -r requirements.txt` – file located in the root of the project.

```
1 args = argparse.parse_args()
2
3 if args.runTests:
4     runner = unittest.TextTestRunner(verbosity=2)
5     result = runner.run(testing.test.load_tests(args))
6 else:
7     PandaCompiler(args).compile()
```

If the `--runTests` flag is set, then it will take priority over the regular compiler functionality. However, it is still possible to specify `--debug`, as debugging information may be useful in case some tests fail. Debugging information is information such as graphical representation of data structures and pretty printed ILOC code – sequential assembly IR. Testing will be explained in depth later.

One aspect that is important to consider is time. Keeping track of when various things happen is hard.

“Some decisions are made when the compiler is designed, at design time. Some algorithms run when the compiler is built, at build time. Many activities take place when the compiler itself runs, at compile time. Finally, the compiled code can execute multiple times, at runtime.”

— Copper and Torczon [2022](#), p. 8

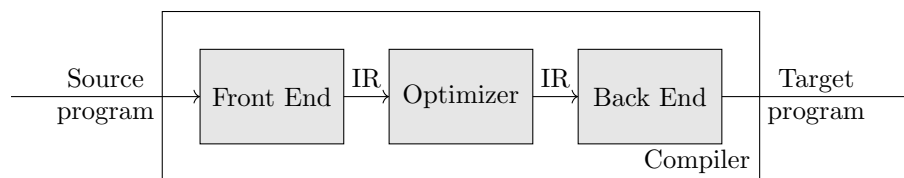


Figure 3.1: Three-phase compiler.

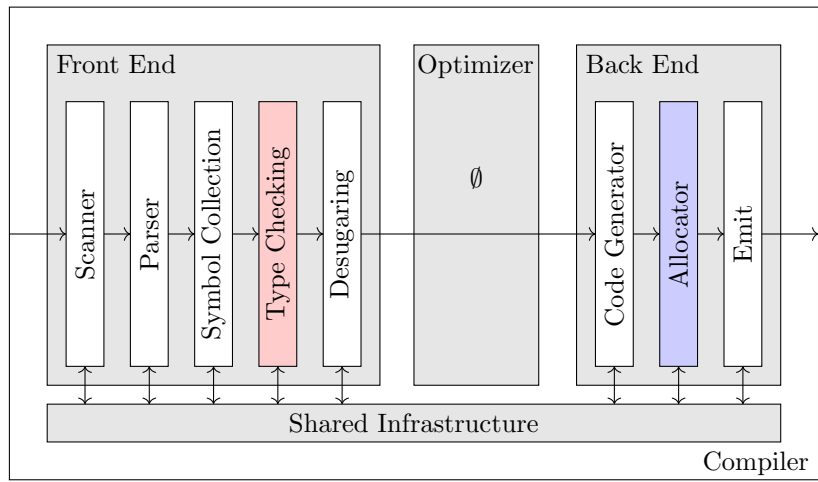


Figure 3.2: Internal structure of Panda.

3.1 Scanner

```
1 reserved = {  
2     'print': 'PRINT',  
3     'return': 'RETURN',  
4     ...  
5 }
```

```
1 tokens = (  
2     'IDENT', 'INT', 'FLOAT',  
3     'PLUS', 'MINUS', 'TIMES', 'DIVIDE',  
4     ...  
5 ) + tuple(reserved.values())
```

```
1 t_PLUS = r'\+'
```

```
2 t_MINUS = r'\-'  
3 ...
```

```
1 def t_IDENT(t):  
2     r'[a-zA-Z_][a-zA-Z_0-9]*'  
3     t.type = reserved.get(t.value, 'IDENT')  
4     return t
```

3.2 Parsing

```
1 def p_program(t):
2     'program : body'
3     interfacing_parser.the_program = AST.Function(
4         "?main", None, t[1], t.lexer.lineno)
```

```
1 @dataclass
2 class Body(AstNode):
3     decls: DeclarationList
4     stm_list: StatementList
5     lineno: int
```

```
1 precedence = (
2     ('nonassoc', 'NEQ', 'LT', 'GT', 'LTE', 'GTE'),
3     ('right', 'EQ'),
4     ('left', 'PLUS', 'MINUS'),
5     ('left', 'TIMES', 'DIVIDE')
6 )
```

```
1 def p_body(t):
2     'body : optional_declarations optional_statement_list'
3     t[0] = AST.Body(t[1], t[2], t.lexer.lineno)
```

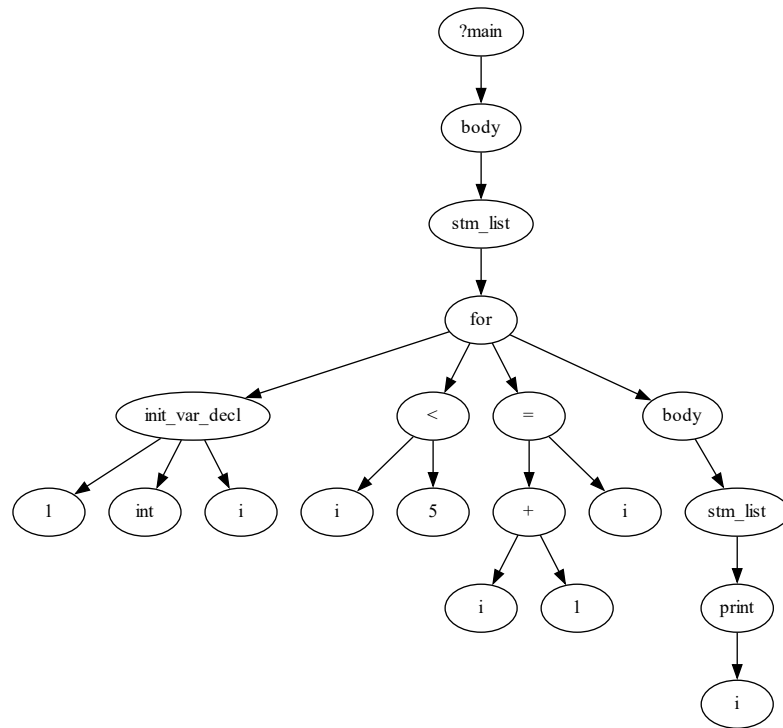


Figure 3.3: Abstract Syntax Tree.

3.3 Symbol Collection

```
1 @dataclass
2 class Symbol:
3     type: str
4     kind: NameCategory
5     info: int
6     SR: int = None
7     escaping: bool = False
```

```
1 @dataclass(init=False)
2 class SymbolTable:
3     level: int
4     parent: SymbolTable
5     _tab: dict
6
7     def __init__(self, parent: SymbolTable) -> SymbolTable:
8         self._tab = {}
9         self.level = parent.level + 1 if parent else 0
10        self.parent = parent
```

```
1 case AST.DeclarationFunction(type, func, lineno):
2     symval = Symbol(type, NameCategory.FUNCTION, func)
3     self._current_scope.insert(func.name, symval, lineno)
4     self._current_scope = SymbolTable(self._current_scope)
5     self._build_symbol_table(func)
6 case ...
```

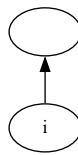


Figure 3.4: Symbol collection.

3.4 Desugaring

```
1 def _desugar_AST(self, ast_node: AST.AstNode) -> None:
2     match ast_node:
3         case AST.Body(decls):
4             self._desugar_AST(decls)
5             decl_var_init_list = self._collect_decl_var_init(decls)
6             assignments = self._trans_var_init_list_to_stm(
7                 decl_var_init_list)
8             self._insert_stm(assignments, ast_node)
9             self._desugar_AST(ast_node.stm_list)
10        case ...
```

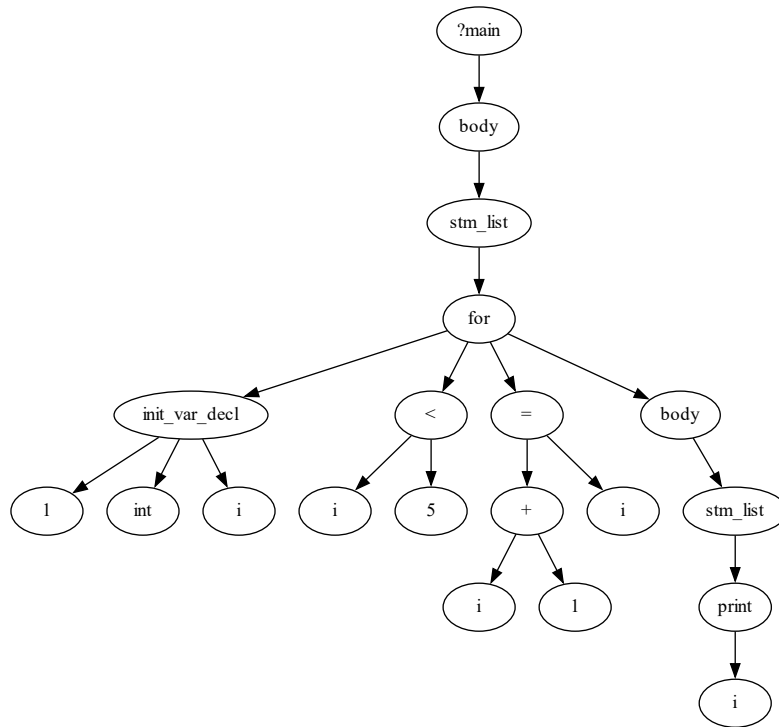


Figure 3.5: Desugaring tree.

3.5 Code Generation

3.5.1 ILOC code

3.5.2 Stack Machine

3.5.3 Register Allocation

The main benefit of testing is the identification and subsequent mitigation of errors. It is to be expected to find errors in larger projects, which is why there must be a way to quickly and dynamically run tests during development. Testing helps software developers compare expected and actual output in order to improve quality. In this case, only system testing is carried out, although a higher granularity of testing would only enrich the project. There exist many tools to solve this task, which is why the big challenge lies in choosing the one that best suits the problem. The most advanced tool is not always the best solution. The chosen tools in this project are the testing framework `unittest` and the coverage tool `coverage.py`. In addition, SonarLint is used to perform static code analysis within the IDE to provide best-practice hints, cognitive complexity metrics and more – not important but worth mentioning.

4.1 System Testing

The most common choice would have been to use `pytest` because it is significantly more popular than `unittest`, but for the task `unittest` seemed like the right choice. I am not able to judge whether the same could have been achieved using `pytest` with parametrized tests, but setting up `unittest` was straight forward.

I wanted a solution where it is possible to load tests dynamically, based on test files in a given folder. Hence, it would not be necessary to write more code just because tests are added to the test suite – writing test cases must not be a burden; otherwise it will not be done. The idea is as follows (in pseudo Python):

```
1 class TestCase(unittest.TestCase):
2     def __init__(self, <args>) -> TestCase:
3         super().__init__()
4
5     def runTest(self):
6         pass
```

First, A class inheriting from `unittest.TestCase` is declared, which makes it possible to define a test case with the required interface. A single test case can then be created by creating an instance of `TestCase`. A test case is created for every pair of `file.panda` and `file.eop` by walking the directory `test-cases`.

```
1 def load_tests(args: argparse.Namespace) -> unittest.TestSuite:
2     test_cases = unittest.TestSuite()
3
4     for src, res in test-cases:
5         test_cases.addTest(TestCase(src, res))
```

The method named `runTest` will be executed when running the test, unless otherwise specified. `runTest` has the responsibility to report whether an individual test succeeds or fails. My `runTest` does the following:

1. Compile source code;
2. Execute compiled code and pipe `std:out` to file;
3. Assert whether output and expected output is identical.

Of course, some exception handling needs to be done, but that is the basics. An interesting exception handling is when running a test developed to fail, as the exception handling in the compiler will execute `os.exit(1)` (it does not make sense to continue), so that exception must be caught in the running test in order to retrieve whatever is printed to `std:err`.

The test suite can be run as follows when the desired tests have been added:

```
1 runner = unittest.TextTestRunner(verbosity=2)
2 result = runner.run(testing.test.load_tests(args))
```

The above logic is contained within `main.py` and wrapped in an `if-else` statement. Running the tests therefore requires calling `main.py` with option `--runTests`, or simply `-t`, as shown below:

```
1 Compiler$ python3.10 main.py --runTests
2
3 runTest (testing.test.TestCase)
4 Testing testing/test-cases/fibonacci_classic.panda ... ok
5 runTest (testing.test.TestCase)
6 Testing testing/test-cases/static_nested_scope.panda ... ok
7 ...
8 -----
9 Ran 21 tests in 2.261s
10
11 OK
```

The test suite is automatically run on every pull request or push to GitHub as part of the CI workflow using GitHub Actions. Configuration can be found in

`/.github/workflows/unittest.yml`. A `ValueError` is raised when a test fails, otherwise the GitHub Actions workflow will not detect this.

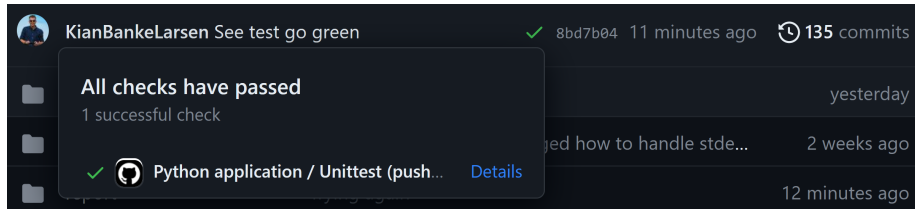


Figure 4.1: GitHub Actions for unit testing.

That way, one know exactly which push or pull request caused the tests to fail, and associated code is directly available from the workflow via commit ID. History for workflow is available too.

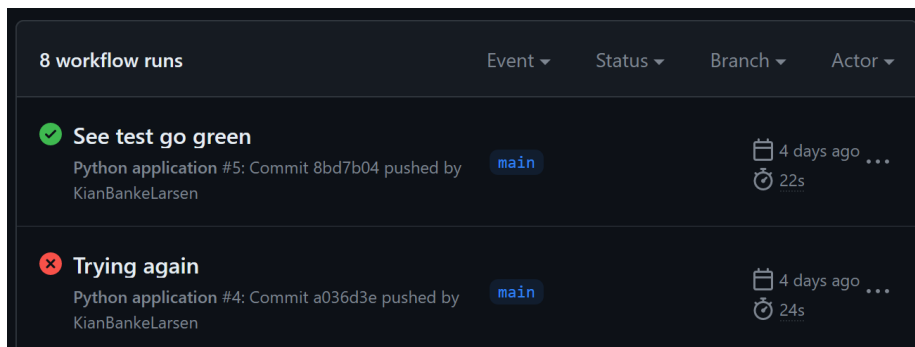


Figure 4.2: Workflow history.

4.2 Coverage

Coverage cannot and should not be a measure of test quality, because the order in which functions and flow are executed has a great influence on the result. However, code coverage gives a direct warning if subsets of the code have not been tested at all. Thus, it is clear that more tests need to be written.

The code coverage measuring tool `coverage.py` is used to perform code coverage statistics on this project. Code coverage is performed in the following way:

```
1 Compiler$ python3.10 -m coverage run main.py --runTests -d
2
3 Name                               Stmts  Miss  Cover
4 -----
5 main.py                             21      1   95%
6 src/compiler.py                     59      3   95%
7 src/dataclass/AST.py               125      0  100%
8 src/dataclass/iloc.py               22      0  100%
9 src/dataclass/symbol.py             35      2   94%
```

10	src/enums/code_generation_enum.py	39	0	100%
11	src/enums/symbols_enum.py	5	0	100%
12	src/phase/code_generation_base.py	57	3	95%
13	src/phase/code_generation_register.py	232	8	97%
14	src/phase/code_generation_stack.py	196	3	98%
15	src/phase/emit.py	128	6	95%
16	src/phase/lexer.py	44	8	82%
17	src/phase/parser.py	101	3	97%
18	src/phase/parsetab.py	18	0	100%
19	src/phase/symbol_collection.py	86	0	100%
20	src/phase/syntactic_desugaring.py	65	0	100%
21	src/printer/ast_printer.py	141	3	98%
22	src/printer/generic_printer.py	17	0	100%
23	src/printer/symbol_printer.py	40	0	100%
24	src/utils/error.py	5	0	100%
25	src/utils/interfacing_parser.py	1	0	100%
26	src/utils/label_generator.py	9	0	100%
27	src/utils/x86_instruction_enum_dict.py	2	0	100%
28	testing/test.py	73	0	100%
29	...			
30	...			
31	...			
32	...			
33	-----			
34	TOTAL	1521	40	97%

Based on this output, it can be assessed that the prepared tests in the `test-cases` folder cover the code well. Note that code coverage is run with the debug flag, `-d`. This is because code in the printer files are only executed when debug is desired.

It is possible to convert the coverage data to an HTML report with the command `python3.10 -m coverage html`. The advantage of having the data in report form is that it is possible to clearly see which lines have been executed and which are missing, as shown in the figure below:

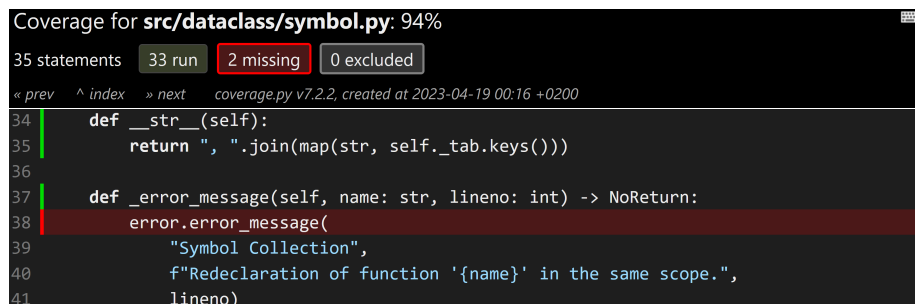


Figure 4.3: Coverage HTML report.

The HTML report also makes it possible to carry out filtering etc. in the file overview table, and thus get a nicer and more user-friendly interface.

5

Performance Comparison

6

Evaluation

6.1 Language Considerations

6.2 Further Development

7

Conclusion

Bibliography

- Beazley, David M. (2018). *PLY (Python Lex-Yacc)*.
<https://www.dabeaz.com/ply/ply.html>. Online; accessed 16. April 2023.
- Copper, Keith D. and Linda Torczon (Oct. 2022).
Engineering a Compiler. English. 3. Morgan Kaufmann Publishers.
ISBN: 978-0-12-815412-0.