

Department of Mathematics & Computer Science
University of Southern Denmark | IMADA
Monday 29th May, 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	Parser	8
3.3	Symbol Collection	12
3.4	Desugar	14
3.5	Code Generation	16
3.5.1	Stack Machine	21
3.5.2	Register Machine	22
3.6	Allocator	24
3.7	Emit	29
4	Testing	32
4.1	System Testing	32
4.2	Coverage	37
5	Performance Comparison	39
6	Evaluation	41
6.1	Retrospect	41
6.2	Further Development	42
6.3	Single Point Of Contact	43
7	Conclusion	44

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.

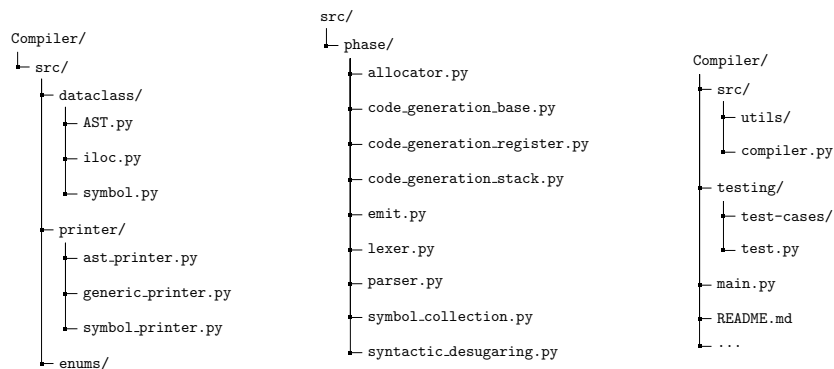


Figure 2.1: Project file tree.

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 somewhat 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] [-s]
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                        Name of assembly output file
11   -c, --compile          Compile output with gcc
12   -d, --debug            Debugging information, i.e., ILOC and Graphviz
13   -f FILE, --file FILE  Path to input file; default is stdin.
14   -t, --runTests        Run tests
15   -r, --run             Run compiled program
16   -s, --stack           Use stack only; default is registers
```

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

Clearly most time is spent on design time, because it is at design time the compiler has been created in the development environment. This includes time spent on designing an interface and writing the code. Compile time is equally important, as it can be coded using more or less efficient algorithms. Perhaps the compiler depends on certain packages being available on the machine, just as various python packages are used in this project. Some code has a very short lifetime in compile time before spending the rest of its life in runtime. It is therefore worth spending substantially more time during compile time in order to perform code analysis and subsequent optimization, such that the code running can be more efficient in some metric. The code could for example be optimized to run faster or save power. Fast programs are wanted when analyzing big data or fast response is needed. Focusing on power saving programs is key when targeting portable devices.

A typical three-phase compiler is designed as shown in Figure 3.1.

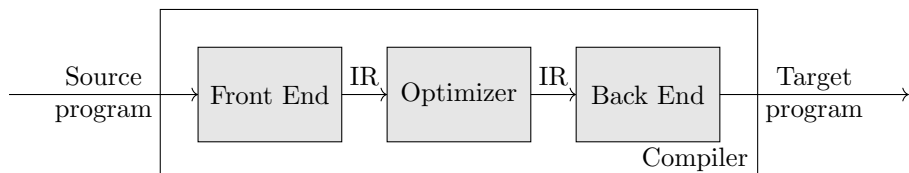


Figure 3.1: Three-phase compiler.

There is no restriction on what the individual stages can contain, as it depends entirely on architecture and program needs. The front end’s responsibility is to understand the input program, such that it is possible to generate correct and meaningful code in the back end. The breakdown of Panda is presented in Figure 3.2.

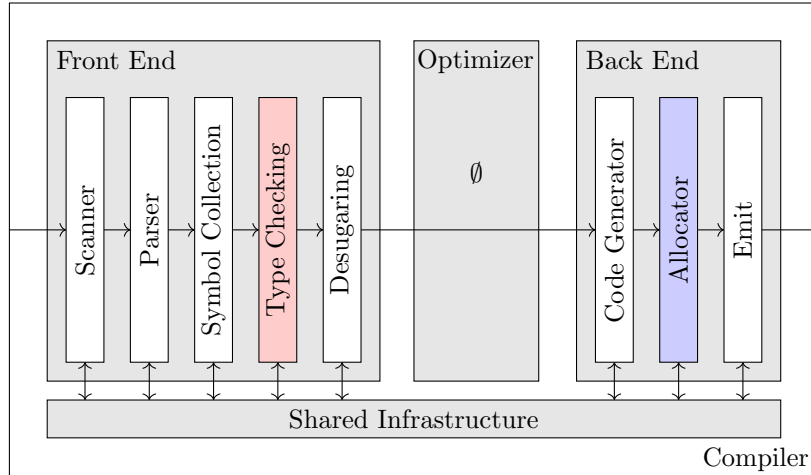


Figure 3.2: Internal structure of Panda.

Type checking has been colored red because it has not been implemented yet, though it would be preferable to have functionality to refrain users from doing something nonsensical, for example, assigning a function to an integer variable.

```
1  int a; int main(){ a = main;
```

This error will cause the compiler to raise a value error during code generation, as the match statement for case `AST.StatementAssignment` is designed to do exactly that for any symbol with `NameCategory` not parameter or variable – function is neither. It is wanted that the compiler stops compiling on such error, but it should never crash because of an uncaught exception. The user must always be informed properly about what went wrong, without having to deal with an indifferent stack trace.

Optimizations like instruction selection, instruction scheduling and peephole etc. could have been interesting to implement, but this is unfortunately beyond the scope of this project, since time did not allow. The stage is therefore skipped.

Panda can either produce stack machine code or register allocation code. The allocator phase is colored blue because it is clearly only used upon compiling source code to assembler code utilizing CPU registers.

The individual phases presented in Figure 3.2 will be reviewed in detail in the subsequent sections of this chapter. As a consequence of having studied the SCIL compiler, it will be clear that several things are done similarly, or at least inspired by SCIL to some extent.

3.1 Scanner

The scanner, or lexical analyzer, is the first phase of the compiler's front end. The scanner reads a stream of characters and produces a stream of words by aggregating the characters. For each word, it determines if the word is valid in the source language, and each valid word is assigned a syntactic category (corresponding to terminal symbols in the language's grammar) used by the parser.

The scanner is the only phase in the compiler that touches every character of the source program. Because grouping characters is a simple task, scanners lend themselves to fast implementations. The Lex part of PLY is an automatic scanner generation tool. Lex requires a token list and specification of token values in order to produce a recognizer. The recognizer is based on a mathematical description of the language's lexical syntax in form of regular expressions and finite automata. What method Lex uses to produce the scanner is unknown, but it can be achieved using Kleene's construction as presented in Copper and Torczon 2022, p. 45.

The documentation for Lex can be found in section 4 Beazley, David M. 2018, and the guide is almost identical to the steps implemented in SCIL. The author of the tool tried to stay faithful to the way in which traditional Lex/Yacc tools work, so there was nothing surprising about how it should be done.

First, it is necessary to make a list of reserved words, as otherwise it will not be possible to separate them from regular identifiers, because reserved words are a subset of legal identifiers. The reserved words are something like `if`, `else`, `return` and so on. These are specified in Python using a dictionary with word as key and token as value.

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

The immutable list of tokens is the union of tokens and reserved words. The list constitutes all valid words of the source program.

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

Rules for individual tokens are specified with a raw regex string as shown below:

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

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

The prefix `t_` is used to indicate that it defines a token. The string must be compatible with Python's `re` module.

A token can also be specified as a function if some kind of action needs to be performed. This is particularly useful in connection with separating identifiers from reserved words. In that case, it will be possible to lookup a given identifier in reserved words. If the relevant key is not found, then it can be concluded that the token in question is just a regular identifier; otherwise, the found reserved token will be returned.

```
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
```

This approach is recommended as it greatly reduces the number of regular expression rules.

3.2 Parser

Parsing is the second phase of the compilers front end. The parser's task is to determine whether a stream of tokenized words produced by the scanner constitutes a valid sentence in the programming language. The parser uses a context-free grammar to derive a syntactic structure for the program, fitting the tokenized words into the grammatical model. If the parser determines that the tokenized program is a valid program, it builds an intermediate representation (shortened IR). The resulting IR of this phase is a graphical IR, acyclic parse tree, because it encodes the program structure well.^[1] Choosing the correct IR has major consequences on what information that can be encoded. It is necessary to have information encoded explicitly, as it is not possible to modify the implicit meaning of some structure.

The Yacc part of PLY uses a LALR(1) parser, which is a bottom up parser, meaning that the parser will attempt to build a derivation bottom up. Though we do not have to deal with the parser implementation, it is worth mentioning that top-down parsers are usually more intuitive, but bottom up parsing is more practical and efficient for complex languages.

The parser is used as a pull parser, meaning that the parser will request a new token whenever it is ready for more input. This is set up by letting the parser control both the input program and the lexer.

```
1 src.phase.parser.parser.parse(
2     user_program,
3     lexer=src.phase.lexer.lexer
4 )
```

^[1] We will consider the properties of a linear IR in section 3.5.

The parsing result is provided by side effect, which is why the graphical IR can be retrieved by reading the variable called `interfacing_program` located in `utils/interfacing_parser.py`.

The parser implementation has been done following the guidelines presented in section 6 Beazley, David M. 2018. Moreover, the grammar implemented is heavily inspired by how it was done in SCIL.

The most prominently used idea is to store the entire program as a function that can be executed by the operating system by calling `main` – the programs single access point. In this way, the idea of stack frames is already established, and the content of the program is thereby just the function body, which in turn can contain other function calls. The starting grammar rule is stated below:

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

Since the source program must be structured as specified by `Body`, it could be interesting to see what fields are available.

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

The `Body` class is annotated with `dataclass`. Python data classes are specially structured class optimized for storage and representation, and it allows for concise class declaration since the constructor is generated automatically.

To make the language semantics easy to understand, it has been decided to constraint that declarations of variables must come before any usage. Declarations and statements are both optional (body can be empty), though useless without statements. Mixing declarations and statements can be confusing to handle, and it will be elaborated in section 3.3 why it encourages trouble.

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

Associativity of tokens can be encoded directly in the grammar, but can also be specified as precedence directives for the parser. It makes the grammar simpler when one can ignore this kind of trouble. The precedence directives stated below specify that the parser must reduce when encountering `TIMES`, `DIVIDE` etc. and shift when encountering `EQ`. The directives go from lowest to highest precedence. The `nonassoc` directive used on comparison based operators are very useful, because it refrains the user from specifying some like `3 < x < 7`. The semantic

of this statement is powerful, although difficult to express directly in code, as it is syntactic sugar for two subsequent statements. The `nonassoc` directive will cause the parser to throw an error whenever it encounters such case.

```

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

```

Parsing is a relatively large operation, but it was a representative sample.

To gain better insight into the resulting IR of the parsing phase, one can advantageously use the printer class to create a graphical representation of the DAG data structure. We will consider the input program stated in Figure 3.3.

```

1 int j = 5;
2 for(int i = 1; i < 5; i = i + 1){}

```

Figure 3.3: Example program.

The resulting abstract syntax tree is given in Figure 3.4.

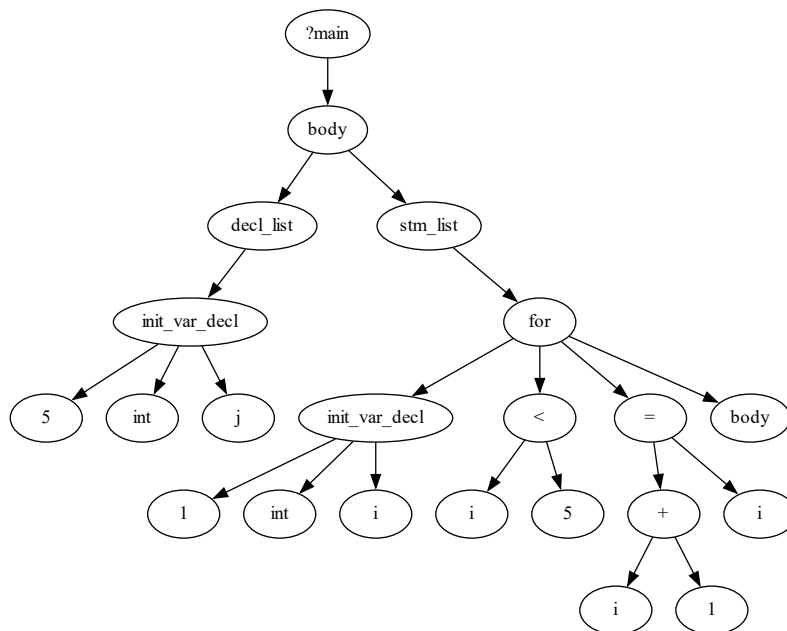


Figure 3.4: Abstract Syntax Tree.

The input program stated in Figure 3.3 is a good example because it is a small program with an associated small parse tree, although the parse tree is more verbose than it necessarily should be. Furthermore, there are two different uses of the `init_var_decl` vertex – this vertex contains syntax such as `int i = 1;`. All usages of that vertex is syntactic sugar for a declaration and simultaneous assignment. The iterative `for` construction uses it as a parameter, so that can just be transformed directly when visited. The story is a little more complicated for other usages of this vertex, since code for declarations is first generated after all statements belonging to a `body` vertex have been visited. It has therefore been necessary to add a desugaring phase to transform the AST, as we will see later in section 3.4.

Panda is designed in such a way that every curly bracket enclosed environment creates a new scope. The default scope is global scope. Every scope consists of a single `body` vertex – the grammar is ‘`new_scope : LCURL body RCURL`’. Recall that a `body` vertex can contain optional declarations followed by optional statements. Just to mention boolean values, the convention used is that the integer value 0 is `false`; otherwise `true`.

The syntax is very C-like, but below are some examples of the different language constructions. Starting with a `for` loop:

```
1  for(int i = 1; i < 5; i = i + 1){...}
```

Then a `while` loops follows naturally:

```
1  while(...){...}
```

The `if-else` flow control statement:

```
1  if(...){...} else{...}
```

The classical `print` statement:

```
1  print(3 + 5);
```

And finally function calls:^[2]

```
1  int a; int b = 1;
2  int fee(int arg){... return ...;}
3  void foo(){...}
4
5  a = fee(b);
6  foo();
```

More examples can be found in the `/testing/test-cases/` folder.

^[2] Remark: it is not valid to assign `void` function calls to any variable.

3.3 Symbol Collection

Symbol collection is the third phase of the compilers front end. The compiler discovers the names and properties of many entities during parsing. For each name used in the program, the compiler needs a variety of information before it can generate code to manipulate that entity. The Symbol Collection phase implements a recursive traversal of the AST in order to collect symbols into symbol tables.

Each scope has its own symbol table and each symbol within that scope is stored in that symbol table. The sole purpose of scopes is to encapsulate symbols so that logically separate parts of the program do not interfere. The introduction of scopes makes it possible to reuse symbol names, which is convenient, and a necessity when using recursion. The scoping mechanism implemented in Panda is static nested scoping. “Static” because it can be determined where a symbol is accessible and visible by reasoning about the lexical structure of the source program – this is often a cognitive comprehensive task when reading programs based on dynamic scoping. “Nested” refer to the hierarchical separation of scopes. The primary purpose of a symbol table is to resolve names. When the compiler finds a symbol, it needs a mechanism that maps that symbol back to its declaration; otherwise we would not be able to locate the entity in memory.

The compiler for Panda only supports simple 64 bit scalar variables and functions, and some additional information such as source register (SR) and `escaping` which will be elaborated in the code generation phase (section 3.5). The `Symbol` class is defined as stated in Figure 3.5.

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

Figure 3.5: Symbol to be stored in symbol table.

Each symbol encountered is stored in a symbol table associated with information about its `type`, `kind` (variable, parameter or function) and `info` (offset such that the variable can be located on the stack). The following has been inserted, to provide insight into what happens when encountering a function declaration:

```
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 ...
```

In this case, `_current_scope` is the active symbol table. It is thereby possible to instantiate a `Symbol` with information about the function and insert it into the table. Before recursively collecting symbols within that function, a new symbol table is instantiated for that scope.

The way the symbol table is declared is shown below:

```

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

```

The symbol table itself takes care of keeping track of its lexical `level`, its `parent`, and the actual symbol table named `_tab`. The `parent` field makes it possible to search upwards in the structure towards the global scope, if the symbol searched for is not at the current level.

Scopes create a tree structure, but where the leafs point to the root and not the other way around. Considering the input program stated in Figure 3.3 again, we see, in Figure 3.6, that `j` is declared in the global scope (root, level 0) and `i` is declared at level 1 immediately incident to global scope.

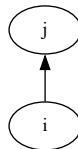


Figure 3.6: Symbol collection.

Because the example just presented is relatively minimal, a more exciting input program is considered below:

```

1  int a;
2  if(3<4){int b;} else{int c; if(3){int c;} else{int d;}}
3  while(4){int d;}

```

It is here clearly shown how to quickly introduce some scopes by using `if-else` statements, since all curly bracket enclosed environments constitute a new scope. Furthermore, it is also noted that two versions of `d` are accessible in two scopes respectively.

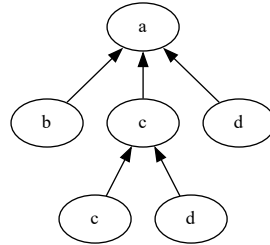


Figure 3.7: Another symbol collection example.

The symbol collection figures are automatically generated using `printer/-symbol_printer.py`. The feature is only provided when executing the compiler with the debug option.

To return to the issue of mixing declarations and statements, it is quickly realized that this places greater demands on parser flexibility, since it will no longer be possible to make a discrete segregation of these groups. As such, it makes no difference to the symbol collection phase, as we already have to recurse through both declarations and statements. I will though make code generation more complicated, because code for functions declared within the scope will be generated simultaneously with statements, if this is not handled with some special case logic. Besides, the mixing also makes it confusing what variable that accessed whenever, since some declaration can appear near the end of the scope – the declaration is valid for the entire scope and not just after the line it appears – this is a hypothetical implementation, open for discussion, of course.

To give an example of what problems may arise with this implementation, please consider the code fragment below. At first sight, one will assume that the input program is valid, but as it is considered, one will realize that `a` is not defined before `print(a)` is called. This would have been harder to spot of if the distance between print and declaration of `a` spanned multiple lines.

```

1 void main(){print(a); int a = 8;}
2 a = 5;
3 main();
4 int a;

```

3.4 Desugar

Desugar is the fourth, and last, phase of the compiles front end. Desugar is a phase that has been included because it was realized that the way nodes were visited during code generation was inexpedient, although necessary. The problem is only associated with the grammar:

```
variable.init.declaration : type IDENT ASSIGN expression SEMICOL
```

Code for statements is produced before producing code for declarations. This is because we want to strangle declared functions such that they do not appear within the function executed by the operating system. The functions would otherwise execute without being called by the user.

The way the problem has been solved is by traversing declarations and collecting all nodes of this type. The nodes are then transformed to **StatementAssignment** nodes and inserted, in the beginning and in the same order, as assignment statements too. The transformed nodes could have been advantageously removed from the tree, to minimize the IR, but since they do nothing, it has been decided to leave them as they are. The algorithmic approach is stated below:

```

1 case AST.Body(decls):
2     self._desugar_AST(decls)
3     decl_var_init_list = self._collect_decl_var_init(decls)
4     assignments = self._trans_var_init_list_to_stm(decl_var_init_list)
5     self._insert_stm(assignments, ast_node)
6     self._desugar_AST(ast_node.stm_list)

```

Reusing the input program from Figure 3.3, we can observe that the yellow nodes have been inserted, in Figure 3.8, as a result of the desugaring phase.

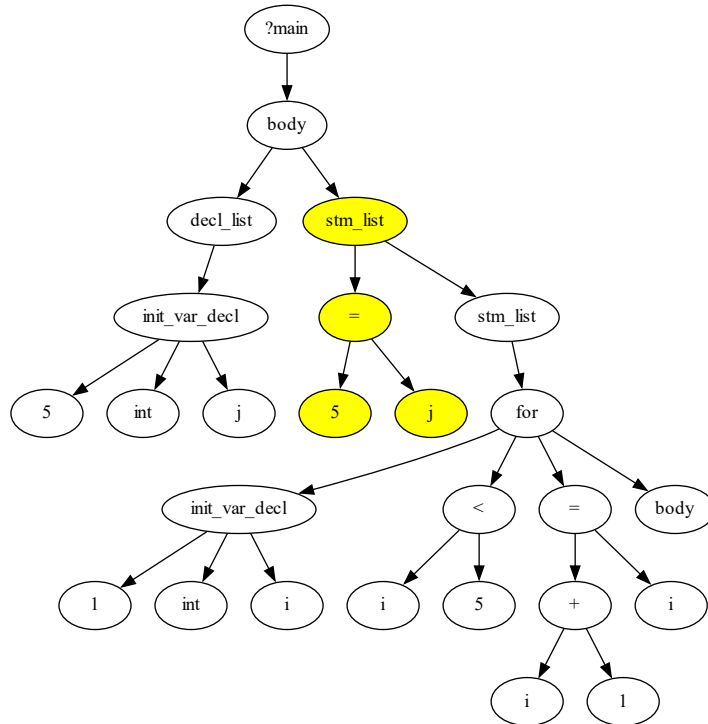


Figure 3.8: Desugaring tree.

3.5 Code Generation

Every process running on the computer is given a virtual address space. Figure 3.9 shows this. Here we note that there are four sections whose significance we should consider. Executable code lives in the *code* section. Static contains statically defined entities – not used for anything other than the **form** string specified in **emit.py**. The Panda programming language does not allow the programmer to define something statically. The *heap* is for dynamically allocating memory, which unfortunately is also not possible in the Panda language. Our primary concern is therefore the *stack*, since the stack contains all variables and information associated with procedure linkage.

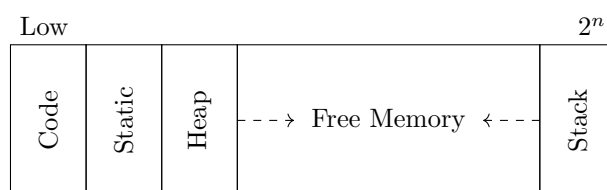


Figure 3.9: Virtual address-space layout.

The stack is a memory area where it is possible to **push** and **pop** elements as needed. One can also choose to statically reserve an amount of memory for data storage. The stack only brings value when used properly and judiciously. It is therefore up to the compiler-writer how the stack is to be used. The actual stack architecture does not matter as long as one decide on a protocol and act accordingly and consistently. However, the final stack architecture is very much controlled by procedure linkage, which limits the number of free choices, since certain things must be located on the stack, although the scheme is a free choice. The chosen stack architecture used in Panda is represented in Figure 3.10. The stack has been strategically divided into sections belonging to the different phases of the procedure linkage. In addition, addressing of significant pointers are shown. RBP (base pointer) and RSP (stack pointer) are used for immediate addressing of entities.

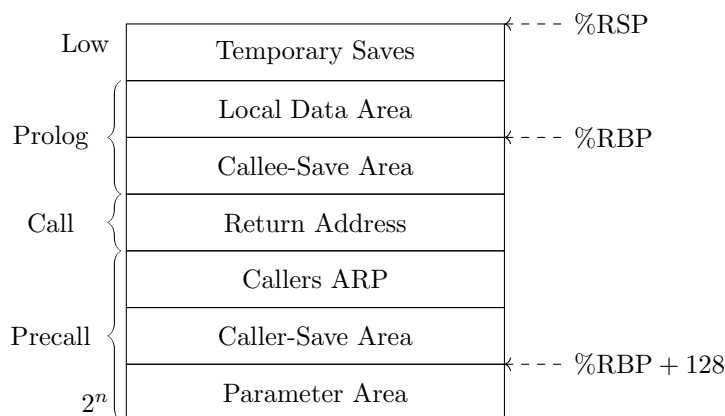


Figure 3.10: Stack Frame.

The first thing the operating system does when executing an assembly program is to call its main function. As mentioned earlier, the entire program's entry point is wrapped as a function, which is why it will be necessary to focus on procedure linkage already now. There are clearly defined rules for which registers caller and callee must preserve. This agreement makes it possible to utilize all registers, since we are guaranteed that callee registers are restored to the same state as when they were borrowed when the caller needs them. Figure 3.11 shows how procedure linkage is implemented in Panda. **Prolog** is responsible for saving callee-save registers and allocating space for local variables on the stack. **Precall** is responsible for putting arguments on the stack, if there are any, saving caller-save registers and making sure to set up static link pointer such that variables declared in other scopes can be accessed. **Postreturn** and **Epilog** are really just inverse processes of **Prolog** and **Precall** to restore registers and free allocated stack space.

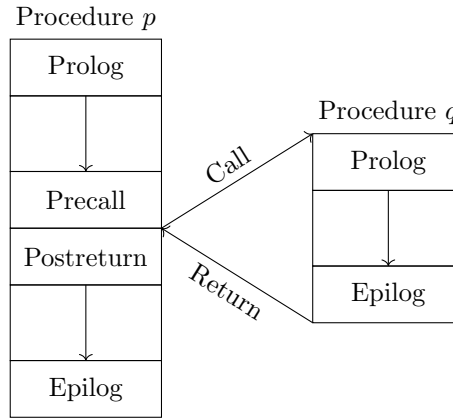


Figure 3.11: Procedure linkage.

The entire procedure linkage process is quite essential for the register machine, whereas the stack machine only needs a subset of the functionality. Because the stack machine uses the stack for everything, it is not necessary to save registers other than RBP, however all register are saved anyway according to convention. The reason is that the final emit phase is reused for both the stack machine and the register machine.

We schematically know how things should be done, we just have to express this in code. The graphical IR is no longer sufficient, which is why we need a linear IR. The linear IR, called ILOC in Copper and Torczon 2022, p. 19, is very similar to the assembly code we want to output. The literature describes ILOC as an instruction set designed for the RISC architecture, since it uses the three-operand format,

add $r_1, r_2 \rightarrow r_3$ // example instruction

where r_1 and r_2 are used for calculation and the result is stored in r_3 . However, CISC is based primarily on two-operand format (second operand is destructive

using AT&T). Therefore, a closely related IR is used instead. The actual implementation is given below. It should be no secret that this is exactly the one used in SCIL, however it has been rewritten to use Python `dataclasses`. An instruction thus consists of an opcode and some operands:

```

1  @dataclass(init=False)
2  class Instruction:
3      opcode: Op
4      args: list(Operand)
5
6      def __init__(self, *args) -> Instruction:
7          self.opcode = args[0]
8          self.args = args[1:]

```

An operand consists of a `target` and an addressing `mode`. Valid targets and modes are those specified in Enum T and Enum M, respectively. These classes are found in the file: `/src/enums/code_generation_enum.py`. This makes it possible to address, for example, registers directly or immediately.

```

1  @dataclass
2  class Operand:
3      target: Target
4      addressing: Mode

```

Translating the AST to the linear ILOC IR thus yields a list of instructions as shown below. The smart thing about this IR is that needed information is explicitly represented. It can be information such as which temporaries are used and their addressing. One can therefore easily perform analyzes on the code by iterating through the list of instructions and using structural pattern matching to extract information. The register machine will pass this IR through the allocator phase in order to assign temporaries to physical registers before emitting code. The stack machine does not require further analysis, instead the list of instructions will simply be passed on to the emit phase which will emit assembly code.

```

1  [
2      Instruction(opcode=<Op.LABEL: '7'>,
3          args=(Operand(target=Target(spec=<T.MEM: 2>, val='main'),
4              addressing=Mode(mode=<M.DIR: 1>, offset=None))),),
5      Instruction(opcode=<Op.META: '8'>, args=(<Meta.PROLOG: 2>,)),
6      ...
7  ]

```

It's time to provide some examples of code generation. Therefore, we consider how to access entities in external scopes. `_follow_static_link`, shown below, solves this problem by repeatedly dereferencing the parent's base pointer (called ARP in Figure 3.10, an acronym for Activation Record Pointer). After the function has returned, RSL can be offset to access the desired entity. The

only parameter the function requires is `symbol_level`, and then it makes sure to follow the static link pointer the required number of times (based on the `_current_scope` field).

```

1 def _follow_static_link(self, symbol_level: int) -> int:
2     level_difference = self._current_scope.level - symbol_level
3     self._append_instruction(
4         Instruction(Op.MOVE,
5                     Operand(Target(T.RBP), Mode(M.DIR)),
6                     Operand(Target(T.RSL), Mode(M.DIR)))
7     )
8     self._get_code_block_to_extend().extend(
9         [Instruction(Op.MOVE,
10                    Operand(Target(T.RSL), Mode(M.IRL, -7)),
11                    Operand(Target(T.RSL), Mode(M.DIR)))
12          for _ in range(level_difference)]
13     )
14
15     return level_difference

```

`_get_code_block_to_extend` is an auxiliary function used to solve a problem that arose in the code generation for the register machine, because it is necessary to preserve blocks of code such that liveness can be performed on one block of instructions at a time. Both `_append_instruction` and `_get_code_block_to_extend` are annotated as abstract in `code_generation_base.py` and therefore defined individually in the classes that uses them. This is a consequence of the register machine and the stack machine not sharing the same structure of the generated code, precisely because code generated by the register machine is not completely linear, but instead a mixture of instructions and lists – elaborated later.

A similar problem is removing stack frames. In the case of returns from deeply nested scopes, it is necessary to be able to remove several stack frames simultaneously, otherwise saved register values and the return address will not be in the expected place, which will cause a segmentation fault. This is solved by iteratively generating code to tear down the individual stack frames by updating RBP. Now we just need to tear down the stack by moving the RBP into the RSP, and that's done as part of the epilogue. The code snippet below is an excerpt of `StatementReturn`.

```

1     _, symbol_level = self._current_scope.lookup(func.name)
2     level_difference = self._current_scope.level - symbol_level
3
4     self._get_code_block_to_extend().extend(
5         [Instruction(Op.MOVE,
6                     Operand(Target(T.RBP), Mode(M.IRL, -7)),
7                     Operand(Target(T.RBP), Mode(M.DIR)))
8          for _ in range(level_difference-1)]
9     )

```

Figure 3.12: Tear down multiple stack frames if needed.

The example below shows how to generate code for a function. The use of match statements contributes to a good overview because it is possible to observe the entire function at once. The steps come in linear order, exactly as we want the internal representation. A good approach when writing code for a construction is therefore to think about the flow for the construction in question, and then simply apply the individual steps in that particular order.

```

1  case AST.Function(body=body):
2      ...
3      self._ensure_labels(ast_node)
4      self._append_instruction(
5          Instruction(Op.LABEL,
6                      Operand(Target(T.MEM, ast_node.start_label), Mode(M.DIR)))
7      )
8      self._prolog(body)
9      self._generate_code(body.stm_list)
10     self._append_instruction(
11         Instruction(Op.LABEL,
12                     Operand(Target(T.MEM, ast_node.end_label), Mode(M.DIR)))
13     )
14     self._epilog()
15     self._append_instruction(Instruction(Op.META, Meta.RET))
16     self._generate_code(body.decls)
17     ...

```

We do not care about the actual content of the function body because there are prepared other match cases that handle it recursively. We just know that whatever instructions the body contains, they must be specified there. This approach is guaranteed to work by the invariant-based code generation scheme.

Invariants are conditions that we can rely on throughout the compiler. The benefit is that it is easier to implement new constructions when we know that we can rely on certain invariants. Invariants are established and maintained by following certain protocols such as the stack protocol. Below are stated some concrete Panda invariants:

- The stack organization in the form of stack frames.
- For constructions with multiple operands, the operands are processed left-to-right. Actual arguments to functions are processed right-to-left.
- Function calls use the stack for all parameters.
- A function's return value is stored in `RAX`.
- 64 bits are used for all types of values.
- Constructions consuming a number of operands pop them from the stack; the stack top is the last argument. This invariant is mostly related to the stack machine, but a similar approach has been taken in the register machine – we just maintain a stack of temporary register.
- Constructions producing a result: The stack machine pushes the result onto the stack. The register machine moves the result into a new temporary.

3.5.1 Stack Machine

It is not possible to review code generation for all constructions for the stack machine, and therefore the code snippet for `ExpressionBinop` below is representatively selected. We start by noting that the piece of code perfectly complies with the presented invariants, so the content should not be surprising.

A binary operation consists of a left and right child. The invariant promises to process operands left-to-right, and therefore the left child is called recursively before the right child is called recursively. The recursive calls return two values on the stack that the binary operation annotated by the vertex applies. Depending on the AST, it may be that the resulting value must be used to evaluate another vertex higher up the tree, and in that way the result propagates upwards.

Structural pattern matching is used to extract `binop`, `lhs` and `rhs` as their own variables. To limit the matching to arithmetic binary operations (not comparisons), an `if-guard` is used to check that binary operation is one of the enumerations in the list. As said, the two recursive calls return two values on the stack. These two values are popped into register 1 and register 2, respectively, after which the binary operation (opcode) is performed with register 1 and register 2 as operands. Because the operation is destructive, the result is saved in register 2, and this value is finally pushed onto the stack. The stack machine is quite worry-free in the sense that we do not have to keep track of both registers and the stack, as everything we need is just consistently on the stack. It is only required that we push and pop in the correct order.

```
1  case AST.ExpressionBinop(binop, lhs, rhs) if binop in [Op.ADD, Op.SUB, Op.DIV,
   ↪ Op.MUL]:
2      self._generate_code(lhs)
3      self._generate_code(rhs)
4
5      self._append_instruction(
6          Instruction(Op.POP,
7                      Operand(Target(T.REG, 1), Mode(M.DIR)))
8      )
9      self._append_instruction(
10         Instruction(Op.POP,
11                    Operand(Target(T.REG, 2), Mode(M.DIR)))
12     )
13     self._append_instruction(
14         Instruction(binop,
15                     Operand(Target(T.REG, 1), Mode(M.DIR)),
16                     Operand(Target(T.REG, 2), Mode(M.DIR)))
17     )
18     self._append_instruction(
19         Instruction(Op.PUSH,
20                    Operand(Target(T.REG, 2), Mode(M.DIR)))
21     )
```

The idea behind comparison binary operations is very similar, but there is a little extra logic to handle pushing 1 when true and 0 otherwise.

3.5.2 Register Machine

Many of the techniques used in the stack machine can be transferred to the register machine. Instead of having a stack of values in the virtual address space used at runtime, we just maintain a stack of temporary registers at compile time. This makes it possible to keep track of which registers are included in various operations. As a starting point, we pretend that we have an infinite number of registers available, and it is therefore up to the allocation phase to map these to physical registers. This functionality is implemented by enriching the class with the two fields below. `_reg_count` is used to keep track of how many temporaries have been used and which unique temporary to use next. `_reg_stack` is just a list that we can pop and push values to. By only using push and pop we achieve a LIFO priority queue AKA. stack.

```
1 _reg_count: int = 0
2 _reg_stack: list[int] = field(default_factory=list)
```

To compare the stack machine with the register machine, we again give the example of `ExpressionBinop`. Two recursive calls are made to produce two temporaries on the stack. Next, we move the second temporary, `reg2`, into a new temporary because the binary operation is destructive. `push_new_reg_count` takes care of incrementing `_reg_count` and pushing the temporary onto the temporary stack. When the binary operation has been performed, we know that the result is located in the new temporary and that it is at the top of the temporary stack, such that it can be used for subsequent operations.

```
1 case AST.ExpressionBinop(binop, lhs, rhs) if binop in [Op.ADD, Op.SUB, Op.DIV,
↳ Op.MUL]:
2     self._generate_code(lhs)
3     self._generate_code(rhs)
4
5     reg1 = self._reg_stack_pop()
6     reg2 = self._reg_stack_pop()
7
8     self._push_new_reg_count()
9
10    self._append_instruction(
11        Instruction(Op.MOVE,
12                    Operand(Target(T.REG, reg2), Mode(M.DIR)),
13                    Operand(Target(T.REG, self._reg_count), Mode(M.DIR)))
14    )
15    self._append_instruction(
16        Instruction(binop,
17                    Operand(Target(T.REG, reg1), Mode(M.DIR)),
18                    Operand(Target(T.REG, self._reg_count), Mode(M.DIR)))
19    )
```

Figure 3.13: Code generation for binary operation.

Instead of letting values flow through the stack, they just flow through registers. However, these are managed using a stack at compile time, so in that way the methods for the register machine and the stack machine are very similar – at least in a programming aspect.

The goal in this phase is, like in the register machine, to produce a linear IR. However, the allocation phase requires that the bracket structure be preserved, and therefore the result is a hybrid IR. The final IR therefore looks as indicated below. A hybrid IR means to combine both graphic and linear IR's to capture their strengths and avoid their weaknesses. A typical hybrid IR is a control-flow graph that uses a linear to represent blocks of code and a graph to represent the flow of control among those blocks. As a starting point, the resulting hybrid IR from this phase is not a control-flow graph, but it has all the prerequisites to become one. The final modifications will be done in the allocation phase.

```
1  [  
2      instruction_1,  
3      instruction_2,  
4      [  
5          instruction_...,  
6          ...,  
7          instruction_n  
8      ],  
9      instruction_n+1  
10 ]
```

Some things have been changed to achieve the above. `_append_instruction` has been modified to always append the given instruction to the last list in `_code`. Furthermore, symbols and scopes must be handled in order to be able to reuse temporaries for identifiers. To create a new symbol scope use `_create_new_symbol_scope` and `_remove_symbol_scope` to remove it again. For each scope that is opened, in addition to handling symbols, a new empty list is appended to `_code`:

```
1  self._code.append([])
```

As a scope is closed, in addition to processing symbols, instruction lists will be merged. The list we are working with is thus at the end as long as it is active, and when the scope is closed it becomes part of the accumulated IR.

```
1  if len(self._code) >= 2:  
2      list_to_merge = self._code.pop()  
3      self._code[-1].append(list_to_merge)
```

The described IR implementation is quite simple in itself, but it is complicated by the fact that it's administration is included in another method. This is done because both methods are used in exactly the same places, however it would have been more readable to keep them separate.

3.6 Allocator

The liveness implementation is based on the literature Appel 1997, pp. 218–223 and the subsequent register allocation (graph coloring) is based on the literature *ibid.*, pp. 235–238.

There are many steps involved in register allocation. In order to make the process clear and provide an overview, each step has its own method, which also call auxiliary methods. `perform_register_allocation` is the only method exposed by the `Allocator` class. The method is included below:

```
1 def perform_register_allocation(self, code: list[Instruction]) ->
  ↪ list[Instruction]:
2     code = copy.deepcopy(code) # input is left untouched
3
4     self._find_labels(code)
5     self._control_flow(code)
6     self._liveness_analysis(code)
7     graph = self._build_graph(code)
8     colors = self._color_graph(graph)
9     code = self._flatmap(code)
10    self._rename_registers(colors, code)
11
12    return code
```

We are only interested in doing liveness analysis for instructions that use temporaries, so structural pattern matching is used to filter these out. Structural Pattern matching makes it easy to determine whether a temporary is used in the first or second operand or both. We therefore do not have to do a lot of conditional checks, because it is all given in the match case.

The first step is to create a mapping between labels and the closest instruction that uses temporaries. This information is used when performing control-flow, otherwise it is not known where the flow should go when a jump instruction is encountered.

Creating control flow is about chaining instructions, so that any instruction points to the next instruction in the flow. Every instruction is given a successor and predecessor attribute, though it was later realized that the predecessor was not needed. The successor attribute is not strictly necessary, as we would be able to obtain the same information during the iteration in another way, but it makes the subsequent steps a bit nicer. Both successor and predecessor would have been strictly necessary if we had based the iteration on following these pointers, but due to the use of match statements, it was judged that it was easier to use a for each loop for the iteration. Due to the scope implementation, it has been necessary to limit the analysis to each introduced scope and not just functions, which is the ideal solution. It is not possible to ensure the flow when registers are saved and restored along the way. Consider Figure 3.14 for an example of a flow graph. In panda, the scope is wrapped around the entire `while` construction. The flow is therefore as shown in the figure, apart from the fact that `body` is of course made up of multiple instructions.



Figure 3.14: Loop construction flow graph.

The limited flow is particularly expressed in an **if** statement. Optimally (but still locally) the flow would look like shown in Figure 3.15a, where, depending on the **if** condition, the flow continues in **then** or **else** and finally in the common final node. In this way, liveness can be calculated for the entire graph. Unfortunately, the case is as shown in Figure 3.15b. The flow is the same (shown with dotted lines), but in relation to the flow, we just take the direct route through the statement. **then** and **else** will be their own scopes, and therefore they will have their own local flow isolated from their surroundings. This is by no means desired, but it is a consequence of a choice made early in the process, because it standardized scopes and made it possible to declare variables in all scopes.

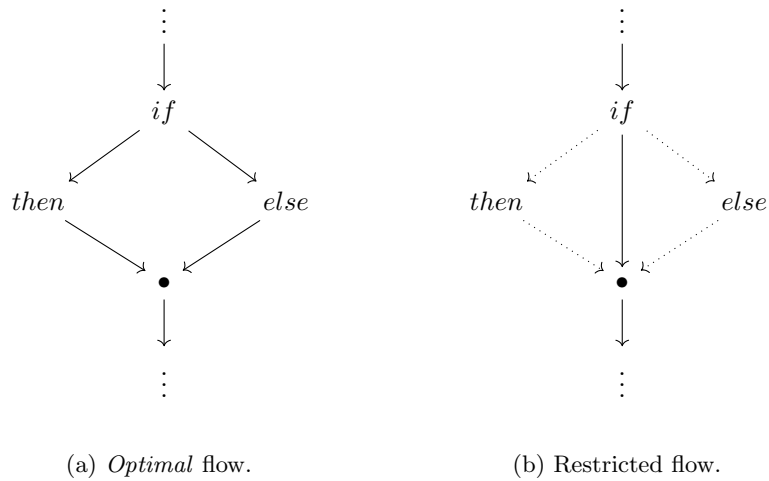


Figure 3.15: **if** construction flow graph.

Liveness analysis is about determining which temporaries are live at the same time. It is crucial in register allocation, because if two variables are not live at the same time, then they can share a register, otherwise they must necessarily claim their own register. We say a temporary is live if it holds a value that may be needed in the future, so this analysis is called *liveness* analysis. We obtain this information using the fixed point algorithm given in Appel 1997, p. 221. The implementation of this algorithm is given below:

```

1 def _liveness_analysis(self, code: list) -> None:
2     change = True
3
4     while (change):
5         change = False
6
7         for ins in reversed(code):
8             if self._do_set_calc(ins):
9                 change = True

```

Note that the instructions are parsed in reverse order. This is because the uses of a temporary come after its declaration. It's just a way to achieve fixed point faster. The algorithm works by preparing an *in* and *out* set for each instruction in the flow. The *in* set contains all the temporaries that are live for the instruction in question. *out* contains all the temporaries that must be taken care of for the next instruction. The algorithm keeps repeating the analysis of the flow until there are no changes to the *in* set. Worst-case run time of this algorithm is $O(N^4)$, *ibid.*, p. 224.

To give intuition for the calculation of *in* and *out* sets, consider Figure 3.16. In the figure, an arbitrary vertex from the flow is considered, where it is indicated how the *in* and *out* sets follow the flow. Two *out* edges are specified because it may be that there is a conditional jump.

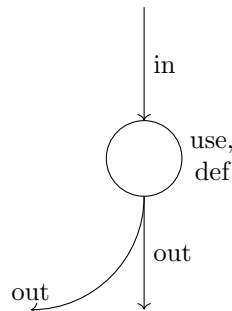


Figure 3.16: Abstract flow graph.

In and *out* sets must be calculated for every vertex in the flow. This is done by using the dataflow equations in Figure 3.17 – utilized in the `_do_set_calc` method. Also note how the definition of a temporary kills its liveness. This is expected behavior because if there are no uses of a definition, then we can simply remove those instructions. Executing useless instructions is useless.

$$in[n] = use[n] \cup (out[n] - def[n])$$

$$out[n] = \bigcup_{s \in succ[n]} in[s]$$

Figure 3.17: Dataflow equations for liveness analysis.

To take a concrete example, consider the register allocation process for the input program below.

```
1  int a = 5; int b = 7; print(a + b);
```

We skip control-flow because the phase does not provide any information to display. The liveness analysis produces the `in` sets below. We only consider instructions that use temporaries, and these instructions are extracted (left) and associated with their `in` set (right).

```
1  movq $1, r1:    {}
2  movq r1, r2:    {1}
3  movq $7, reg3:  {2}
4  movq r3, r4:    {2, 3}
5  movq r2, r5:    {2, 4}
6  addq r4, r5:    {4, 5}
7  push r5:        {5}
```

Based on all these `in` sets, an interference graph is created. The graph is made by adding an edge between any pair of temporaries that share an `in` set. The implementation in Panda makes use of a dictionary to create a mapping (temporary: adjacency set) to represent a graph. The literature suggests making an interference matrix, but it immediately seemed easier to implement the graph with adjacency sets, as the graph have to be strangulated and assembled again during coloring.

```
1  [{1: set(), 2: {3, 4}, 3: {2}, 4: {2, 5}, 5: {4}}]
```

The graph is drawn in Figure 3.18.

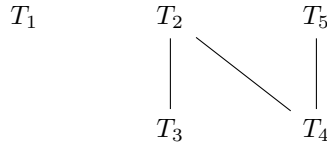


Figure 3.18: Interference graph.

Graph coloring is done using a simple heuristic, called Coloring by Simplification, in the following way:

1. Suppose the graph G has a vertex m with fewer than K neighbors, where K is the number of registers on the machine.
2. Let G' be the graph obtained by removing m from G .
3. When a new vertex m is added to the colored graph G' , it is assigned the lowest numbered color that is not assigned to any of its neighbors. If the neighbors of m have at most $K - 1$ colors among them, then a free color can always be found for m .

This leads to a stack-based algorithm for coloring: we repeatedly remove (and push on a stack) nodes of degree less than K . Each such simplification leads to more opportunity for simplification. If it is not the case that there is a node with degree less than K , then a random vertex is chosen in the graph. There are many heuristics that can be used, and one of them is the highest degree first heuristic because they are more likely to require a new color. The implementation in Panda just randomly pick a vertex with degree lower than K , and if there is no such vertex, then a random vertex with degree K or higher. The coloring implementation is provided below. $K = 10$ in Panda.

```

1 def _color_graph(self, graphs: list[defaultdict]) -> dict[int, int]:
2     colors = defaultdict(lambda: None)
3
4     for graph in graphs:
5         stack = self._take_graph_apart(graph)
6         self._assign_colors(stack, colors)
7
8     return colors

```

Running the algorithm on the input yields the color mapping (temporary: color) given below. The highest degree first heuristic would have colored the graph using only two colors.

```

1 {5: 1, 4: 2, 3: 1, 2: 3, 1: 1}

```

The graph in Figure 3.18 has been colored with the calculated colors in Figure 3.19.

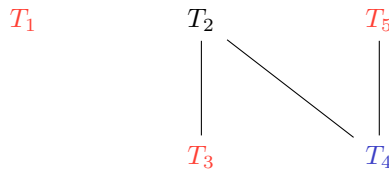


Figure 3.19: Colored interference graph.

3.7 Emit

The Emit phase is the last phase in the compiler's back end. It is in this phase that the internal representation introduced in the code generation is made into actual assembler instructions, which can then be compiled into machine code using `gcc`. Luckily it has been possible to develop an interface such that this class could be used for both the stack machine and register machine.

The mapping from the internal representation is almost 1:1, however there are small deviations such as the use of `%RSL`, which is a virtual register used to calculate links, and that register does not exist natively on the x86 architecture. In addition, it is necessary to handle a handful of `Meta` instructions, because it made sense to assemble these into composite instructions, and then let the Emit phase handle the actual code for this, since it is the same every time – these instructions could as well have been represented in the linear IR. The `Meta` instructions are:

```
1 self._enum_to_method_map = {
2     Meta.CALL_PRINTF: self._call_printf,
3     Meta.PROLOG: self._prolog,
4     Meta.EPILOG: self._epilog,
5     Meta.PRECALL: self._precall,
6     Meta.POSTRETURN: self._postreturn,
7     Meta.RET: self._ret
8 }
```

The `Meta` instructions is thus simply mapped using a dictionary in python, mapping the `Enum` to a function implemented within the `Emit` class.

The majority of these functions are functions used during procedure linkage, which is why it makes sense to declare exactly which registers are callee-save and caller-save registers:

```
1 self._callee_save_reg = ["rbx", "r12", "r13", "r14", "r15", "rbp"]
2 self._caller_save_reg = ["rcx", "rdx", "rsi", "rdi", "r8", "r9", "r10", "r11"]
```

It might also have been advantageous to have specified these fields in another class, so that variable offsets did not have to be hard-coded in the code generation, but could instead have been implemented as a function on the two lists. Adding extra callee-save registers and caller-save registers would then not require further intervention than adding them to the correct list.

`Meta.CALL_PRINTF` uses labels to ensure correct memory alignment, and we therefore want to reuse the `Labels` object instantiated during code generation. This is achieved using the property that the class is declared singleton, and Python will therefore just return the same instance as used in code generation, without manually having to pass that instance to `Emit`:

```
1 self._labels = Labels()
```

It is important that the counter within `Labels` is not reset, since this could, in the worst case, cause clashing labels, meaning that the program would be incorrect and in turn prevent `gcc` from being able to compile it, since labels are no longer ensured to be unique.

The code snippet given below is responsible for “dispatching” any ILOC instruction to an assembly instruction. Because the ILOC instructions are provided as a list, one can simply iterate through the list and translate the individual instruction as intended. The native function `map` is used though this is not the intended use of it, as `_dispatch` accumulates its result in a field called `_code`. However, it is nice to think of it as a mapping, why it is done that way. Also, `map` is a concise way to do iteration even if the resulting list just contains `None`.

```
1 def emit(self, iloc_ir: list[iloc.Instruction]) -> str:
2     self._program_prologue()
3     list(map(self._dispatch, iloc_ir))
4     self._code.append("\n")
5     return "\n".join(self._code)
```

The final result of `Emit` is a single string containing all the generated assembler code, including newlines. The output string can therefore just be inserted directly into any file using Python’s built-in I/O library.

A function called `_do_operand` takes care of handling targets. In order to create a common interface between the register machine and the stack machine, it has been chosen to simply expand the `match` statement, which can also be found in SCIL, so that it is possible to address several registers.

```
1 def _do_operand(self, operand: iloc.Operand) -> str:
2     match operand.target:
3         case iloc.Target(spec=T.IMI, val=val):
4             text = f"${val}"
5             ...
6         case iloc.Target(spec=T.REG, val=1):
7             text = "%rbx"
8         case iloc.Target(spec=T.REG, val=2):
9             text = "%rcx"
10            ...
11         case iloc.Target(spec=T.REG, val=11):
12             text = "%r15"
13         case _:
14             raise ValueError(f"Unkown target: {operand.target}")
```

However, the number of registers is a limited amount, and in this case the limit is 9 (in addition to two registers set aside for computation). The allocation phase may well risk passing on a program that requires 14 registers, which is why it becomes necessary to handle spills. Spills can fortunately be handled locally for the individual instructions, which is why it felt natural to handle this in the `Emit` phase as well.

Besides just handling spills, all the `move` instructions that operate on registers that have not been given a color (colored `None`) are removed. It can occur when, for example, assigning a value to a variable but never ever uses it afterwards. This is achieved with the following piece of code:

```

1 def _dispatch(self, instruction: iloc.Instruction) -> None:
2     if instruction.opcode == Op.MOVE:
3         for arg in instruction.args:
4             if arg.target.spec == T.REG and arg.target.val is None:
5                 return
6     ...

```

`return` in context of `_dispatch` means that we just ignore the instruction and don't generate code for it. This operation can be legally performed because, intuitively, why preserve something that is not used. Modern Compiler Implementation states the following:

“If there is no edge in the interference graph between the source and destination of a move instruction, then the move can be eliminated.”

— Appel 1997, p. 239

Since all nodes in the interference graph are colored, we can conclude that if a vertex has been colored `None`, then it not even part of the graph, thus, the `move` instruction can be safely eliminated.

`handle_reg_spill` is called from `_dispatch`. `handle_reg_spill` takes care of inserting extra instructions in case a spill register needs adjustment. The method has been inserted below:

```

1 def _handle_reg_spill(self, instruction: iloc.Instruction) ->
  ↳ list[iloc.Instruction]:
2     if instructions := self._check_register_dependent_operations(instruction):
3         return instructions
4     instructions.extend(
5         self._check_if_spill_is_needed(instruction)
6     )
7     instructions.append(
8         self._make_reference_to_spilled_registers(instruction)
9     )
10    return instructions

```

The method handles the logic in three phases. Because of the horrible names, we list them in order. The first method call handles binary operations on spilled registers. The method ensures that the operands are in physical registers before performing the operation, and then puts the result back into the spilled register. The second method call assigns a stack place to a spilled register. The third method call manipulates the instructions in question such that `move` instructions refer correctly to the spilled registers on the stack.

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 are many tools available to solve this task, which is why the big challenge lies in choosing the tool 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 25 tests in 4.972s
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.

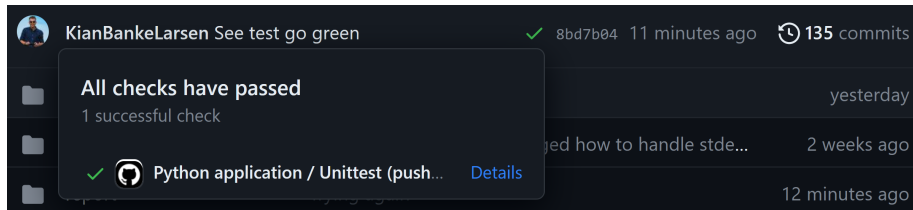


Figure 4.1: GitHub Actions for system testing.

That way, one knows 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.

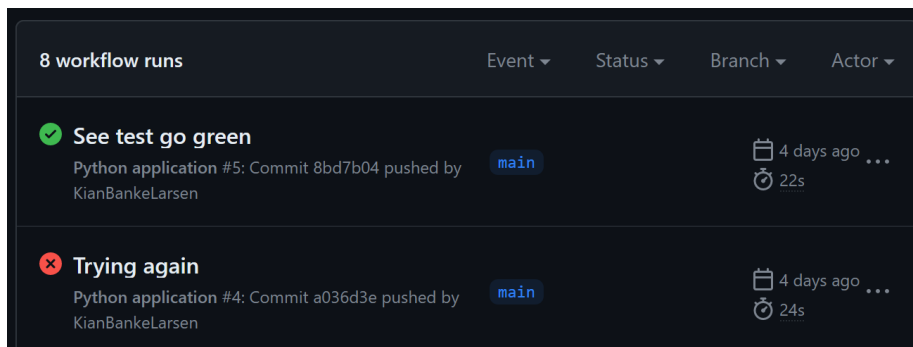


Figure 4.2: Workflow history.

A workflow build consists of the stages stated in Figure 4.3. Note that, in this example, the stack machine test passed, but the register version did not. Upon build failure, an email is sent to whom introduced the issue (me).

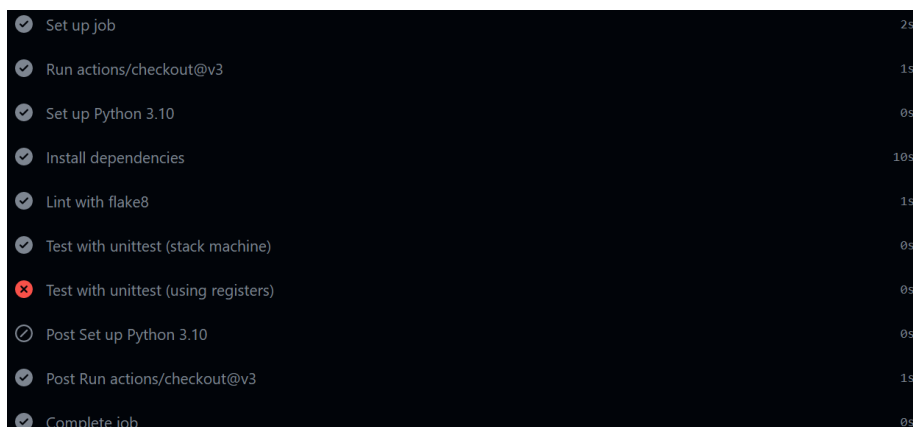


Figure 4.3: Workflow run. Red indicates failure.

The list of tests is too long to mention each one, instead three tests that caused particularly many problems have been chosen. The tests are representative because they have a relatively high complexity, as the input programs utilize recursion, loops and nested scopes. For this to work properly, it requires that the constructs and static link pointer works correctly. In addition, it must be possible to tear down several stack frames at once, since every scope introduces a new stack frame. Failing to remove a stack frame properly causes problems in the `postreturn` and `epilog` phase, and in worst case segmentation fault.

The test below is borrowed from SCIL. The test is named `summers.panda`.

```

1  int sum_recurse(int n) {
2      if(n == 1) {return 1;} else {
3          return n + sum_recurse(n - 1);
4      }
5  }
6  int sum_loop(int n) {
7      int sum, i;
8      i = 1; sum = 0;
9      while(i <= n) {
10         sum = sum + i;
11         i = i + 1;
12     }
13     return sum;
14 }
15
16 print(sum_recurse(9)); print(sum_loop(9));
17 print(sum_recurse(42)); print(sum_loop(42));

```

The test below does not calculate anything meaningful or difficult, but it checks that the functionality implemented in Figure 3.12 works as expected.

```

1  int i = 5;
2  int main(){
3      if(1){
4          if(1){
5              if(1){
6                  i = i + 1;
7                  return i;
8              }
9          }
10     }
11 }
12 print(main());

```

The goal has been to test all legal semantics thoroughly. This includes assignment to variables declared outside the current scope. Assignment to formal parameters, recursion, return from static nested scopes, comments, sugared declaration with assignment, binary operations (comparisons and arithmetic) and constructs like `while`, `for`, `if` and `print`.

A test that is particularly critical for the register machine is right-unbalanced trees,^[1] or complete trees, because these cause a greater register need. An example of such tree can be found in Figure 4.4. The figure is part of the syntax tree for the input program given below:

```
1 print(1+(1+(1+(1+(1+(1+(1+(1+(1+(1+(1+(1+(1+(1+(1+((
```

The reason is, in this case, that the left child is always visited first, after which the expression in the right child must be evaluated. This is clear in Figure 3.13. It is therefore necessary to preserve all left children until there are no more right children. In the case of the stack machine, this is not something one thinks too much about, as the stack just grows as expected. With register allocation, attention is required, because there must be a place to store values, since these must be in a register, regardless of whether this register is spilled onto the stack, it is still a “register”.

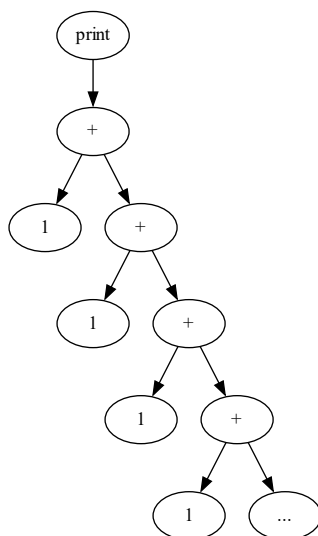


Figure 4.4: Unbalanced expression.

To be precise, the input program uses 27 temporary registers as a consequence of the code generation, which after preparation of the interference graph and coloring can be stored in 14 physical registers. The Emit phase only provides 9 registers for storage, meaning that 5 registers are located on the stack.

[1] This is tested with the test cases `/testing/test-cases/unbalanced_right_addition.panda` and `/testing/test-cases/unbalanced_right_subtraction.panda`.

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 erase
2     && python3.10 -m coverage run -a main.py -td
3     && python3.10 -m coverage run -a main.py -tsd
4     && python3.10 -m coverage report
5
6 Name                               Stmts  Miss  Cover
7 -----
8 main.py                             24     2    92%
9 src/compiler.py                      70     3    96%
10 src/dataclass/AST.py                 125     0   100%
11 src/dataclass/iloc.py                 22     0   100%
12 src/dataclass/symbol.py               36     2    94%
13 src/enums/code_generation_enum.py     38     0   100%
14 src/enums/symbols_enum.py             5     0   100%
15 src/phase/allocator.py                213     5    98%
16 src/phase/code_generation_base.py     58     3    95%
17 src/phase/code_generation_register.py 262     9    97%
18 src/phase/code_generation_stack.py    199     5    97%
19 src/phase/emit.py                     199    17    91%
20 src/phase/lexer.py                    44     8    82%
21 src/phase/parser.py                   101     2    98%
22 src/phase/parsetab.py                  18     0   100%
23 src/phase/symbol_collection.py         117     2    98%
24 src/phase/syntactic_desugaring.py     65     0   100%
25 src/printer/ast_printer.py            141     3    98%
26 src/printer/generic_printer.py         17     0   100%
27 src/printer/symbol_printer.py          40     0   100%
28 src/utils/error.py                     5     0   100%
29 src/utils/interfacing_parser.py        1     0   100%
30 src/utils/label_generator.py           7     0   100%
31 src/utils/x86_instruction_enum_dict.py 2     0   100%
32 testing/test.py                       75     0   100%
33 -----
34 TOTAL                               1884    61    97%
35 Wrote HTML report to htmlcov/index.html
```

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 see which lines have been executed and which are missing, as shown in the figure below:

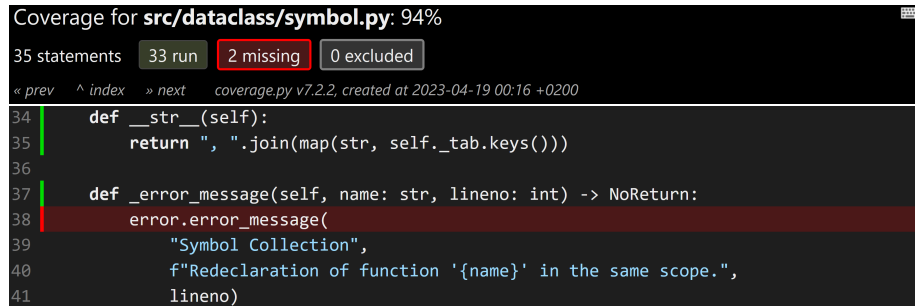


Figure 4.5: Coverage HTML report.

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.

Performance Comparison

It is time to compare the performance of the stack machine and the register machine. Only the three tests indicated in Figure 5.1 have been made, as it is limited how many tests can be made without composite types. The tests are:

- `sum_loop` from `test-cases/summers.panda` called with actual parameter 300000000.
- `fib` from `test-cases/fibonacci_classic.panda` called with actual parameter 39.
- `for-loop` performing 500000000 iterations of adding two numbers.

All test results are calculated as an average of 10 measurements. The tests have been performed on an Intel x86 Core i9-8950HK Coffee Lake Processor, and timed using the native tool `time` in Ubuntu (based on `real` time).

We did not expect much performance gain by using the register machine versus the stack machine, since the register allocation is rather limited. One might even have expected slower executions because values move between more physical registers than necessary (peep-hole optimization required). A remarkable performance gain is observed for all tests when using the register machine.

The reason for the increased performance of the register machine is, by hypothesis, due to the use of move instructions despite push and pop instructions. Based on a scientific paper written by Agner Fog from DTU, Fog, Agner 2022, pp. 314–315, the hypothesis is not without authority. Figure 5.1 presents a table comparing reciprocal throughput. Without getting too technical and analytical about the compiled programs, we achieve approximately 50% better performance by replacing push and pop instructions with move instructions. It agrees very well with the performance gain we see in `sum_loop`. It is expected that Fibonacci is dominated by the complexity associated with handling stack frames, and therefore we do not achieve the same performance gain. The result in the `for-loop` test is crazy, but the primary work is also optimized from a reciprocal throughput of almost 1 to around 0.25. The tests are satisfactory, and we conclude that the hypothesis holds.

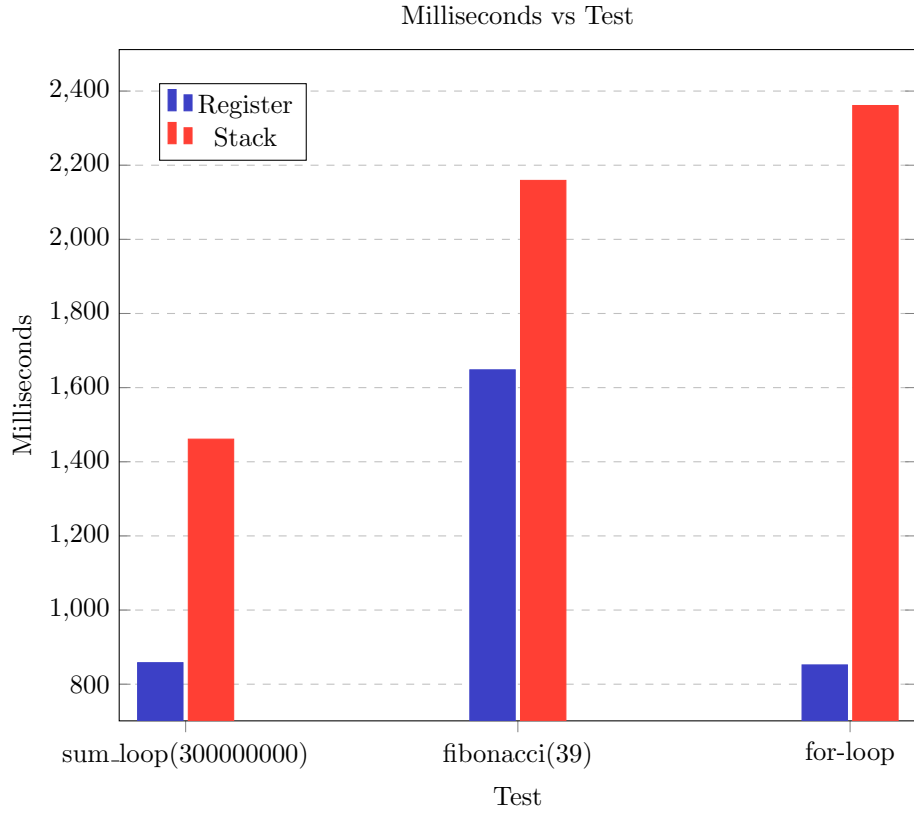


Figure 5.1: Stack Machine vs Register Machine.

Instruction	Operands	Reciprocal throughput
MOV	r,i	0.25
MOV	r32/64,r32/64	0.25
MOV	r32/64, m	0.5
PUSH	r	1
PUSH	i	1
PUSH	m	1
PUSH	stack pointer	1
POP	r	0.5
POP	stack pointer	3
POP	m	1

Table 5.1: Coffee Lake integer instructions.

This chapter is devoted to reflecting on the learning through the project, as well as describing some issues that could be interesting to work on if the project had had a longer duration.

6.1 Retrospect

There are two fundamental principles that must be observed when building a compiler. The first principle is inviolable:

“The compiler must preserve the meaning of the input program.”

— Copper and Torczon 2022, p. 6

The second principle has a more practical aspect:

“The compiler must discernibly improve the input program.”

— *ibid.*, p. 7

It has been a core value throughout the project to make a well-defined compiler implementation that does not do anything inappropriate. Furthermore, the goal has been to create a clear, consistent, easily understandable grammar that a user would be able to pick up with ease. It is assessed that the first principle has been complied with. In this project, no special effort has been made in optimization apart from the preparation of an implementation that uses registers, which is why the second principle must be trivially respected. An example of an optimization that could challenge this would be instruction scheduling, since instructions are literally rescheduled after the invariant based linear IR has been generated.

The approach to the project has been academically based on literature, in order to gain an overview of new as well as old knowledge. Although this has given a confident approach to the project, it must be recognized that it has dulled the process unnecessarily. When coding a project such as a compiler that consists of relatively many lines of code, and many hours of debugging, it is necessary to

purposefully choose a theme and then learn as needed, otherwise many hours are lost on something that is not directly applicable. The hours are of course not wasted, but not beneficial for the project either.

It has always been desired that it should be possible to declare new variables in any scope, which is why it was decided to let `body` (syntactic category from parsing) be what constituted a scope. The choice later made it clear, during code generation phase, that this imposed too many constraints on how the stack should be set up, since a simple generalized procedure linkage was wanted. This resulted in suddenly having to handle all scopes as procedure calls, though it had nothing to with a procedure call. One also encounter problems when doing local register allocation, because flow is only considered locally for a scope (normally local to functions). This consequence limits the analysis to the individual scope, making the content of a `then` block, of an `if` statement, a black box for the surrounding scope. Thus, the analysis of control flow becomes less effective, which contributes to overall poorer performance of the register machine compiler. The easiest solution, and most likely most effective, would be to redo the scope implementation. Global register allocation would be the hard option, though it would not solve the underlying problem besides patching it.

6.2 Further Development

The first thing to solve is the scope implementation. This is possibly where the greatest performance is to be gained, based on the compiler's current functionality. Being able to analyze flow for entire functions will make it easier to optimize the use of registers. Especially in the case of using variables outside the current scope, the variables will then be marked as escaping, which causes them to live on the stack, which is annoying when wanting to use registers. At present, loops cannot really benefit from registers, although `sum_loop` still managed to achieve significantly better performance.

Another challenge is that there is no `push imm64` instruction, which can cause the error: `Error: operand type mismatch for 'push'`. The error occurs, for example, when a number greater than 32 bit two's complement is printed, as the number to be printed will be pushed directly to the stack. This is annoying because the registers and the reserved stack space is 64 bit, but we can efficiently only use 32 bit at worst. A similar error apply to the register machine, as `movq` cannot handle immediate source operands larger than 32 bit two's complement. Workarounds are available, but they enforce a more advanced implementation than prepared in this project.

At present, all spilled registers are stored as temporary saves on the stack, but it might be desirable to move some spilled registers to the local data area and parameter area, so that the space already allocated on the stack is used efficiently.

Finally, it could be pretty cool with composite data types like arrays to make the language Turing Complete. At the present time, it is not possible to calculate anything significant, as one can only work with integers and Boolean values (0 or 1 integer). Composite data types will only further the potential of the already implemented programming constructs.

6.3 Single Point Of Contact

As a final rounding and summary of the project, GitHub Pages have been used as an addendum to the report, hoping to provide a brief user-friendly introduction to the code and the report. A screenshot of the website can be found in Figure 6.1.

The website has been created with Jekyll, which is a tool that can automatically transform plain text into a static website. The tool only requires a small YAML file (`/_config.yml`), listed below, and a markdown file for the page's content.

```
1 theme: jekyll-theme-slate
2 title: Panda Compiler
3 description: This is a bachelor project, Southern University of Denmark
4 show_downloads: true
```

The website is simply a nice presentation of the `README.md` file. In addition, it is possible to download the project directly here, or follow the link to GitHub.

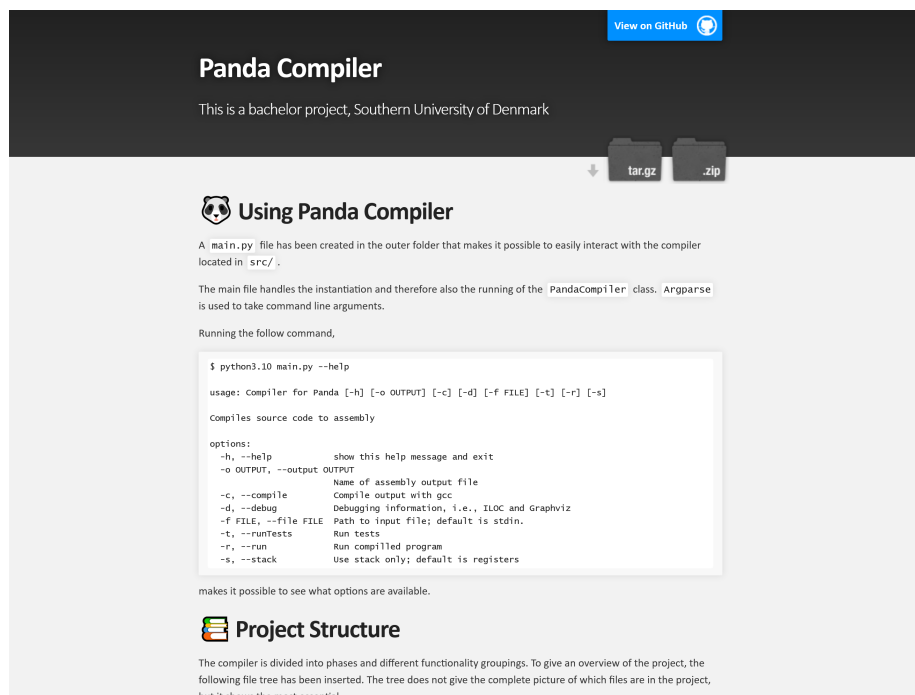


Figure 6.1: Automatically generated Jekyll Slate website hosted on GitHub Pages: <https://kianbankelarsen.github.io/Compiler/>.

GitHub Pages takes care of deployment when a YAML is specified. Subsequently, the website will be automatically deployed on any push to GitHub – at the same time as the workflow for testing runs.

Conclusion

In this bachelor project, a compiler has been successfully implemented, which can be used both as a stack machine and a register machine. The compiler consists of a scanner, parser, symbol collector, desugar, code generator, allocator and a code emitter. The compiler consists of approximately 3500 lines of Python code.

The implemented language is named Panda, and resembles a subset of C. The implemented types are Boolean and Integer (actually both are integers). Constructions such as if-else statements, while-loops and for-loops have also been implemented. 64 bit are used for all types of values. It had been well regarded with composite types, such as arrays, but unfortunately it was not achieved in the time frame.

The performance of the stack machine and the register machine have been compared, and it turns out that the register machine runs remarkably faster, which is a very satisfactory result. Depending on the given test, a performance gain of up to a factor of three is observed.

Correctness of the compiler has been tested using 25 test cases and all tests run successfully. Furthermore, a workflow has been set up to carry out these tests via GitHub Actions, which ensures that all code pushed to production is tested. In the event of an error, the person responsible will receive a notification accordingly. Although the 25 tests do not provoke an error, this does not mean that the compiler does not contain any errors, just that they have not been discovered yet. The 25 tests achieve a code coverage of 97%, which is sufficient.

Every project has a README. In addition to providing a light introduction to the project, a static website has been generated on the file, from which the code is publicly available: <https://kianbankelarsen.github.io/Compiler/>

The implementation of scopes is not optimal, which is why it could be interesting to refactor, given that there was more time to work on the project. Only when this is done will it be possible to fully utilize the register allocation, and in that way, hopefully, achieve an even greater performance gain.

Bibliography

- Appel, Andrew W. (1997). *Modern Compiler Implementation in C*. English. Cambridge University Press. ISBN: 978-0-52-158390-9.
- 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.
- Fog, Agner (2022). *4. Instruction tables*.
https://www.agner.org/optimize/instruction_tables.pdf.
Online; accessed 28. May 2023.