

Computer Science Tripos

Part II Project Proposal Coversheet

Please fill in Part 1 of this form and attach it to the front of your Project Proposal.

Part 1

Name:	<input type="text" value="Martin Walls"/>	CRSID:	<input type="text" value="mrw64"/>
College:	<input type="text" value="Churchill"/>	Overseers: (Initials)	<input type="text" value="SH, NK"/>
Title of Project:	<input type="text" value="C to WebAssembly Compiler"/>		
Date of submission:	<input type="text" value="14 October 2022"/>	Will Human Participants be used?	<input type="text" value="No"/>
Project Originator:	<input type="text" value="Timothy M. Jones"/>		
Project Supervisor:	<input type="text" value="Timothy M. Jones"/>		
Directors of Studies:	<input type="text" value="Dr John K. Fawcett"/>		
Special Resource Sponsor:	<input type="text"/>		
Special Resource Sponsor:	<input type="text"/>		

Part 2

Overseer Signature 1: -----

Overseer Signature 2: -----

Overseers signatures to be obtained by Student Administration.

Overseers Notes:

Part 3

SA Date Received:	<input type="text"/>	SA Signature Approved:	<input type="text"/>
-------------------	----------------------	------------------------	----------------------

Part II Project Proposal: C to WebAssembly Compiler

Martin Walls

October 2022

Overview

With the web playing an ever-increasing role in how we interact with computers, applications are often expected to run in a web browser in the same way as a traditional native application. WebAssembly is a binary code format that runs in a stack-based virtual machine, supported by all major browsers. It aims to bring near-native performance to web applications, with applications for situations where JavaScript isn't performant enough, and for running programs originally written in languages other than JavaScript in a web browser.

I plan to implement a compiler from the C language to WebAssembly. C is a good candidate for this project because it is quite a low-level language, so I can focus on compiler optimisations rather than just implementing language features to make it work. Because C has manual memory management, I won't have to implement a garbage collector or other automatic memory management features. Initially I will provide support for the stack only, and if time allows I will implement `malloc` and `free` functionality to provide heap memory management.

I will compile a subset of the C language, to allow simple C programs to be run in a web browser. A minimal set of features to support will include arithmetic, control flow, variables, and functions (including recursion). I won't initially implement linking, so the compiler will only handle single-file programs. This includes not linking the C standard library, so I will provide simple implementations of some of the standard library myself, as necessary to provide common functionality such as `printf`.

I will use a lexer and parser generator to do the initial source code transformation into an abstract syntax tree. I will focus this project on transforming the abstract syntax tree into an intermediate representation—where optimisations can be done—and then generating the target WebAssembly code.

I plan to write the compiler in Rust, which is memory safe and performant, and has lexer/parser generators I can use.

To test and evaluate the compiler, I will write small benchmark programs that individually test each of the features and optimisations I add. For example, I will use the Fibonacci program to test recursion. I will also test it with Conway's Game of Life, as an example of a larger program, to test and evaluate the functionality of the compiler as a whole.

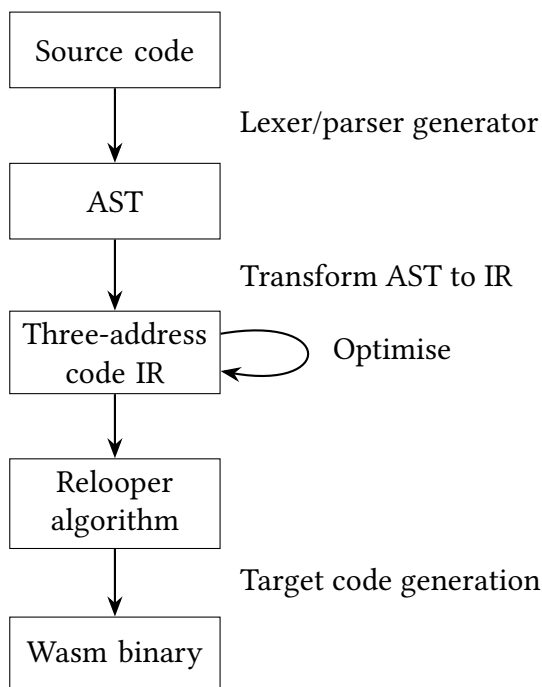
I will use a three-address code style of intermediate representation, because this lends itself to perform optimisations more easily. For example it's easier to see the control flow in three-address code

compared to a stack-based representation. To transform from abstract syntax tree to the intermediate representation, this will involve traversing the abstract syntax tree recursively, and applying a transformation depending on the type of node to three-address code.

To transform from the intermediate representation to WebAssembly, I will need to convert the three-address code representation into a stack-based format, since WebAssembly is stack-based. This stack-based format will have a direct correspondence to WebAssembly instructions, so the final step of the compiler will be writing out the list of program instructions to a WebAssembly binary file.

C allows unstructured control flow (e.g. goto), whereas WebAssembly only supports structured control flow. Therefore I will need a step in the compiler to transform unstructured to structured control flow. One algorithm to do this is the Relooper algorithm, which was originally implemented as part of Emscripten, a LLVM to JavaScript compiler¹.

Compiler pipeline overview



Starting point

I don't have any experience in writing compilers beyond the Part IB Compiler Construction course. I haven't previously used any lexer or parser generator libraries. I've briefly looked at Rust over the summer, but haven't written anything other than simple programs in it.

I have briefly looked up the instruction set for WebAssembly and have written a single-function program that does basic arithmetic, in WebAssembly text format. I used wat2wasm to convert this to a WebAssembly binary and ran the function using JavaScript.

¹<https://github.com/emscripten-core/emscripten/blob/main/docs/paper.pdf>

I have briefly researched lexer and parser generators to see what's out there and to help decide on which language to write my compiler in, but I haven't used them before.

Success criteria

The project will be a success if:

- The program generates an abstract syntax tree from C source code.
- The program transforms the abstract syntax tree into an intermediate representation.
- The program uses the Relooper algorithm to transform unstructured to structured control flow.
- The program generates WebAssembly binary code from the intermediate representation.
- The compiler generates binary code that produces the same output as the source program.

Optimisations

First I will implement some simple optimisations, before adding some more complicated ones.

One of the simple optimisations I will implement is peephole optimisation, which is where we look at short sections of code and match them against patterns we know can be optimised, then replacing them with the optimised version. For example, redundant operations can be removed, such as writing to the same variable twice in a row (ignoring the first value written), or a stack push followed immediately by a pop. Null operations (operations that have no effect, such as adding zero) can also be removed.

Constant folding is another quite simple optimisation that performs some arithmetic at compile time already, if possible. For example, the statement $x = 3 + 4$ can be replaced by $x = 7$ at compile time; there is no need for the addition operation to be done at runtime.

These optimisations will be run in several passes, because doing one optimisation may then allow another optimisation to be done that wasn't previously available. The optimisation passes will run until no further changes are made.

The stack-based peephole optimisations (such as removing pushes directly followed by a pop) will be done once the three-address code representation has been transformed into the stack-based format in the final stage.

A more complicated optimisation to add will be tail-call optimisation, which removes unnecessary stack frames when a function call is the last statement of a function.

Other harder optimisations are left as extensions to the project.

Extensions

Extensions to this project will be further optimisations. These optimisations are more complicated and will involve more analysis of the code.

One optimisation would be dead-code elimination, which looks through the code for any variables that are written to but never read. Code that writes to these variables is removed, saving processing power and space.

Another optimisation would be unreachable-code elimination, where we perform analysis to find blocks of code that can never be executed, and removing them. This will involve control flow analysis to determine the possible routes the program can take.

Evaluation

To test and evaluate the compiler, I will use it to compile a variety of different programs. Some of these will be small programs I will write to specifically test the features and optimisations of the compiler individually. I will also write a larger test program to evaluate the compiler as a whole.

In addition, I will use some pre-existing benchmark programs to give a wider range of tests. For example, cBench is a set of programs for benchmarking optimisations, which I could choose appropriate programs from. The source for cBench is no longer available online, but my supervisor is able to give me a copy of them.

For each of these, I will verify that the generated WebAssembly code produces the same output as the source program when run.

To evaluate the impact of the optimisations, I will run the compiler once with optimisations enabled and once with them disabled, on the same set of programs. I will then benchmark the performance of the output program to identify the impact of the optimisations on the program's running time, and I will also compare the size of the two programs to assess the impact on storage space.

Work Plan

1	14th - 28th Oct	<p>Preparatory research, set up project environment, including toolchain for running compiled WebAssembly. I will research the WebAssembly instruction set.</p> <p>I will also write test C programs for Fibonacci and Conway's Game of Life. To help with my WebAssembly research, I will implement the same Fibonacci program in WebAssembly by hand.</p> <p>Milestone deliverable: <i>I will write a short LaTeX document explaining the WebAssembly instruction set, from the research I do.</i></p> <p><i>C programs of Fibonacci and Conway's Game of Life, and a WebAssembly implementation of Fibonacci.</i></p>
---	-----------------	---

2	28th Oct - 11th Nov	<p>Lexer and parser generator implementation.</p> <p>This will involve writing the inputs to the lexer and parser generators to describe the grammar of the source code and the different types of tokens.</p> <p>Milestone deliverable: <i>Lexer and parser generator inputs. The compiler will be able to generate an abstract syntax tree (AST) representation from a source program.</i></p>
3	11th - 25th Nov	<p>Implementation of transforming the AST into the intermediate representation. This will require defining the intermediate code to generate for each type of node in the AST.</p> <p>Milestone deliverable: <i>The compiler will be able to generate an intermediate representation version from a source program.</i></p>
4	25th Nov - 9th Dec	<p>Researching and implementing the Relooper algorithm.</p> <p>Milestone deliverable: <i>The compiler will be able to transform unstructured control flow into structured control flow using the Relooper algorithm. I will also write a short LaTeX document describing the algorithm.</i></p>
5	9th - 23rd Dec	<p>Implementation of target code generation from intermediate representation.</p> <p>For each type of instruction in the intermediate representation, I will need to define the transformation that generates WebAssembly from it.</p> <p>Milestone deliverable: <i>The compiler will be able to generate target code for a source program. The generated WebAssembly will be able to be run in a web browser.</i></p>
Two weeks off over Christmas		
6	6th - 20th Jan	<p>(I'll be more busy during the first week of this with some extracurricular events before term.)</p> <p>Slack time to finish main implementation if necessary. Implement some peephole optimisations (how many I do here depends on how much of the slack time I need).</p> <p>Milestone deliverable: <i>The basic compiler pipeline will be complete. Some peephole optimisations will be implemented.</i></p>
7	20th Jan - 3rd Feb	<p>Write progress report.</p> <p>Continue implementing optimisations, in particular implementing tail-call optimisation.</p> <p>Milestone deliverable: <i>Completed progress report. (Deadline 03/02)</i></p>
8	3rd - 17th Feb	<p>(I'll be more busy here with extra-curricular events.)</p> <p>Slack time to finish main optimisations if necessary. If time allows, work on extension optimisations.</p>

		<i>Milestone deliverable: The compiler will be able to generate target code with optimisations applied. Evidence to show the impact of the optimisations.</i>
9	17th Feb - 3rd Mar	Evaluate the compiled WebAssembly using a variety of programs (as described above), including correctness and impact of optimisations. Write these evaluations into a draft evaluation chapter. <i>Milestone deliverable: Draft evaluation chapter.</i>
10	3rd - 17th Mar	Write introduction and preparation chapters. <i>Milestone deliverable: Introduction and preparation chapters.</i>
11	17th - 31st Mar	Write implementation chapter. <i>Milestone deliverable: Implementation chapter.</i>
12	31st Mar - 14th Apr	Write conclusions chapter and finish evaluations chapter. <i>Milestone deliverable: Evaluations and conclusions chapter. First draft of complete dissertation.</i>
13	14th - 28th Apr	Adjust dissertation based on feedback. <i>Milestone deliverable: Finished dissertation.</i>
14	28th Apr - 12 May	Slack time in two weeks up to formal deadline, to make any final changes. <i>Milestone deliverable: Final dissertation submitted. (Deadline 12/05)</i>

Resources

I will primarily use my own laptop for development. I accept full responsibility for this machine and I have made contingency plans to protect myself against hardware and/or software failure.

My laptop specifications are:

- Lenovo IdeaPad S540
- CPU: AMD Ryzen 7 3750H
- 8GB RAM
- 2TB SSD
- OS: Fedora 35

I will use Git for version control and will regularly push to an online Git repository on GitHub. I will clone this repository to the MCS and regularly update the clone, so that if my machine fails I can immediately continue work on the MCS.