



Automatic Parallelisation of Rust Programs at Compile Time

Author

Michael Oultram
1428105

Degree

MSci in Computer Science
11th February 2018

Supervisor

Dr Ian Batten

Institution

School of Computer Science
University of Birmingham

Contents

1	Preamble	3
1.1	Abstract	3
1.2	Acknowledgements	3
2	Introduction	4
2.1	Rust Language Features	4
3	Design	5
3.1	Dependency Analysis	5
3.2	Scheduler	5
3.3	Reconstructor	5
4	Implementation	6
4.1	Analysis Stage	6
4.2	Modification Stage	10
5	Evaluation	11
6	Discussion	12
7	Conclusion	13
8	Appendix	14
8.1	Submission File Structure	14
8.2	Running the Code	14
8.3	References	15

Chapter 1

Preamble

1.1 Abstract

TODO: Write Abstract

All software for this project can be found at <https://github.com/MichaelOultram/FYP>

1.2 Acknowledgements

Rust Compiler (<https://github.com/rust-lang/rust>)

Serde (<https://serde.rs/>): used to convert rust objects into JSON and back again

Chapter 2

Introduction

TODO: Write introduction

2.1 Rust Language Features

Rust is similar to other programming languages such as C++ but it does has some specific features that may not be known to the reader. This section briefly explains features of the language that are used in later sections of the report. For more in depth explanations of these features, the reader should look at the language documentation (*The Rust Book* 2017).

In Rust, all variables have an ownership. Only one block can have access to that variable at a time. This is enforced at compile time.

```
1 fn main() {  
2     let a = 10;  
3     f(&a);  
4     g(a);  
5     // Cannot access a here anymore  
6 }  
7 fn f(a: &u32){} // f borrows a  
8 fn g(a: u32){} // g takes ownership of a
```

TODO: finish and rewrite

Chapter 3

Design

3.1 Dependency Analysis

3.2 Scheduler

3.3 Reconstructor

DRAFT

Chapter 4

Implementation

4.1 Analysis Stage

Here is some text above

```
syntax_extension.rs
1 use syntax::codemap::Span;
2 use syntax::ptr::P;
3 use syntax::ast::{self, Stmt, Block, Item, ItemKind, Ident};
4 use syntax::ext::base::{MultiItemModifier, ExtCtxt, Annotatable};
5 use syntax::print::pprust;
6 use std::ops::Deref;
7
8 use serde_json;
9
10 use AutoParallelise;
11 use CompilerStage;
12 use dependency_analysis;
13 use shared_state::Function;
14 use scheduler::{self, ScheduleTree};
15
16 fn create_block(cx: &mut ExtCtxt, stmts: Vec<Stmt>) -> Block {
17     let block = quote_block!(cx, {});
18     Block {
19         stmts: stmts,
20         id: block.id,
21         rules: block.rules,
22         span: block.span,
23         recovered: block.recovered,
24     }
25 }
26
27 impl MultiItemModifier for AutoParallelise {
28     fn expand(&self, cx: &mut ExtCtxt, _span: Span, _meta_item: &ast::MetaItem, _item:
29     ↪ Annotatable) -> Vec<Annotatable> {
30         // Only make changes when on the Modification stage
31         if self.compiler_stage != CompilerStage::Modification {
32             return vec![_item];
33         }
34         let mut output = vec![];
35         // Unwrap item
36         if let Annotatable::Item(ref item) = _item {
37             // Find function name and the analysed function
38             let func_ident = item.ident;
39             let func_name = func_ident.name.to_string();
```

```

39         println!("\n\n{:?}", func_name); // Function Id
40
41         if let ItemKind::Fn(ref _fndecl, ref _unsafety, ref _constness, ref _abi, ref
↪ _generics, ref _block) = item.node {
42             println!("{:?}", _fndecl); // Function decl
43
44             // Find function from analysed stage
45             let mut maybe_analysed_function: Option<&Function> = None;
46             for func in &self.functions {
47                 if func.ident_name == func_name {
48                     maybe_analysed_function = Some(func);
49                 }
50             }
51             if let Some(analysed_function) = maybe_analysed_function {
52                 // Merge the dependency trees
53                 let mut base_deptree = dependency_analysis::analyse_block(&_block);
54                 dependency_analysis::merge_dependencies(&mut base_deptree,
↪ &analysed_function.encoded_deptree);
55
56                 println!("DEPTREE:");
57                 for node in &base_deptree {
58                     println!("{:?}", node);
59                 }
60
61                 // Produce a schedule
62                 let schedule = scheduler::create_schedule(&base_deptree);
63                 let schedule_json = match serde_json::to_string_pretty(&schedule) {
64                     Ok(obj) => obj,
65                     Err(why) => panic!("Unable to convert AutoParallelise to JSON:
↪ {:?", why),
66                 };
67                 println!("SCHEDULE:\n{}\n", schedule_json);
68
69                 // Convert schedule into multi-threaded code
70
71                 // Gather all the synclines and create them as variables
72                 let synclines = schedule.get_all_synclines();
73                 println!("Synclines:\n{:?}\n", synclines);
74                 let mut stmts = vec![];
75                 for ((to_a, to_b), (from_a, from_b)) in synclines {
76                     let line_name = format!("syncline_{}_{}_{}_{}", to_a, to_b,
↪ from_a, from_b);
77                     let sx = Ident::from_str(&format!("_{}_send", line_name));
78                     let rx = Ident::from_str(&format!("_{}_receive", line_name));
79                     let stmt = quote_stmt!(cx, let (lsx, lrx) =
↪ std::sync::mpsc::channel()).unwrap();
80                     stmts.push(stmt);
81                 }
82                 stmts.append(&mut spawn_from_schedule(cx, schedule.list()));
83
84                 let new_block = create_block(cx, stmts);
85

```

```

86         // Convert function into use new_block
87         let new_func = ItemKind::Fn(_fndecl.clone(), *_unsafety, *_constness,
↪ *_abi, _generics.clone(), P(new_block));
88
89         let new_item = Item {
90             attrs: item.attrs.clone(),
91             id: item.id,
92             ident: item.ident,
93             node: new_func,
94             span: item.span,
95             tokens: item.tokens.clone(),
96             vis: item.vis.clone(),
97         };
98         // Prints the function
99         println!("converted_function:\n{}\n",
↪ pprust::item_to_string(&new_item));
100
101         let anno_item = Annotatable::Item(P(new_item));
102         output.push(anno_item);
103     } else {
104         panic!("{}", "was not found as an analysed function", func_name);
105     }
106 } else {
107     panic!("ItemKind was not FN");
108 }
109 } else {
110     panic!("Annotatable was not Item");
111 }
112
113 output
114     //vec![_item]
115 }
116 }
117
118 fn spawn_from_schedule<'a>(cx: &mut ExtCtxt, sch: &Vec<ScheduleTree<'a>>) -> Vec<Stmnt> {
119     let mut output = vec![];
120     let mut threads = vec![];
121
122     for i in 0..sch.len() {
123         if let ScheduleTree::SyncTo(ref stmtid1, ref stmtid2, _) = sch[i] {
124             let &(to_a, to_b) = stmtid1;
125             let &(from_a, from_b) = stmtid2;
126             let line_name = format!("syncline_{}_{}_{}_{}", to_a, to_b, from_a, from_b);
127             let sx = Ident::from_str(&format!("{}", "_send", line_name));
128             let prereq = quote_stmt!(cx, lsx.send()).unwrap();
129             output.push(prereq);
130         } else {
131             let mut thread_contents = vec![];
132             // Add prereqs and create a schedule for children
133             let ((lo, hi), mut children) = match sch[i] {
134                 ScheduleTree::Block(ref prereqs, ref spanning_tree, _) |
135                 ScheduleTree::Node(ref prereqs, ref spanning_tree) => {

```



```

136         println!("{:?} paths: {:?}", spanning_tree.node.get_stmtid(),
    ↪     spanning_tree.node.get_env());
137
138         // Add prereqs
139         let (from_a, from_b) = spanning_tree.node.get_stmtid();
140         for &(to_a, to_b) in prereqs {
141             let line_name = format!("syncline_{}_{}_{}_{}", to_a, to_b,
    ↪     from_a, from_b);
142             let rx = Ident::from_str(&format!("{}",_receive", line_name));
143             let prereq = quote_stmt!(cx, &rx.recv().unwrap()).unwrap();
144             thread_contents.push(prereq);
145         }
146
147         // Spawn children after node
148         let children = spawn_from_schedule(cx, &spanning_tree.children);
149
150         // Return node id
151         (spanning_tree.node.get_stmtid(), children)
152     }
153     ScheduleTree::SyncTo(_, _, _) => panic!("Unreachable case"),
154 };
155
156     // Add the current item
157     match sch[i] {
158         ScheduleTree::Node(_, ref spanning_tree) => {
159             // Copy statement of node
160             let stmt = spanning_tree.node.get_stmt().unwrap().deref().clone();
161             thread_contents.push(stmt);
162         }
163         ScheduleTree::Block(_, _, ref schedule) => {
164             // Add block to the schedule
165             let mut inner_block = spawn_from_schedule(cx, schedule.list());
166             let block = create_block(cx, inner_block);
167             let stmt = quote_stmt!(cx, &block).unwrap();
168             thread_contents.push(stmt);
169         }
170         ScheduleTree::SyncTo(_, _, _) => panic!("Unreachable case"),
171     };
172
173     // Add children
174     thread_contents.append(&mut children);
175
176     if i == sch.len() - 1 {
177         // Last uses the current thread
178         output.append(&mut thread_contents);
179     } else {
180         // All except the last is put into a concurrent thread
181         let thread_sname = format!("thread_{}_{}", lo, hi);
182         let thread_name = Ident::from_str(&thread_sname);
183         let thread_block = create_block(cx, thread_contents);
184         let thread_stmt = quote_stmt!(cx, let &thread_name =
    ↪     std::thread::spawn(move || &thread_block);).unwrap();

```

```
185         output.push(thread_stmt);
186         threads.push(thread_name);
187     }
188 }
189 }
190
191 // Join all threads
192 for thread in threads {
193     let thread_stmt = quote_stmt!(cx, thread.join().unwrap()).unwrap();
194     output.push(thread_stmt);
195 }
196
197 output
198 }
```

Listing 1: **TODO:** Remove this

Here is some text below 1

4.2 Modification Stage

Chapter 5

Evaluation

DRAFT

Chapter 6

Discussion

TODO: Write Discussion

DRAFT

Chapter 7

Conclusion

TODO: Write conclusion

DRAFT

Chapter 8

Appendix

8.1 Submission File Structure

TODO: Write section

8.2 Running the Code

Install rustup:

```
curl https://sh.rustup.rs -sSf | sh
```

Install rustc 1.25.0-nightly (0c6091fbd 2018-02-04) compiler:

```
rustup install nightly-2018-02-04
```

```
rustup default nightly-2018-02-04
```

Verify the version is correct using:

```
rustc --version
```

Create a new crate and write sequential code:

```
cargo init
```

Under [dependencies] in Cargo.toml, add:

```
auto_parallelise = { version = "0.1.0", git = "https://github.com/MichaelOultram/FYP/",  
branch = "release" }
```

At the top of your lib.rs or main.rs file, add:

```
#![feature(plugin)]  
#![plugin(auto_parallelise)]
```

At the top of every function, add:

```
#[autoparallelise]
```

Compile the code twice using (the first time will fail):

```
cargo build --release
```

Run the parallelised code using:

```
cargo run --release
```

8.3 References

The Rust Book (2017). URL: <https://doc.rust-lang.org/book/> (visited on 12/11/2017).

DRAFT

Automatic Parallelisation of Rust Programs at Compile Time

Project Proposal

1 Problem

Kish (2002) estimated the end of Moore's Law of miniaturization within 6-8 years or earlier (based on their publication date) and as such, manufacturers have been increasing processors' core count to increase processor performance (Geer 2005). Writing parallelised programs to take advantage of these additional cores has some difficulty and often requires significant changes to the source code. Is it possible to automate these changes to convert sequential code into parallelised code? Previous attempts at solving this problem include D'Hollander, Zhang and Wang (1998) where they automated parallelisation of sequential FORTRAN code and Bas-karan, Ramanujam and Sadayappan (2010) where they automated conversion of sequential C into CUDA code. Both of these approaches use unsafe programming languages significant complexity to their solutions. Instead, can this problem be solved more easily with a safe programming language like rust (*The Rust Programming Language* 2017)?

2 Approach

My approach involves writing a plugin for the rust compiler. The rust compiler plugin system allows for different types of plugin, which are run at different stages of the compilation. A syntax extension plugin can manipulate the abstract syntax tree of any annotated function, one function at a time. A linter plugin can see the abstract syntax tree of every function without annotations but it cannot manipulate anything.

This problem requires read access to the entire abstract syntax tree, and then modify access to the parallelisable parts. Due to the order of execution, the syntax extension plugin would be executed first and then the linter plugin but I would require them to be executed in the opposite order. To solve this problem, I propose that the program is compiled twice.

On the first compilation, the syntax extension plugin would do nothing. The linter plugin would view the entire abstract data tree and analyse what each statement depends on and modifies. This would create a dependency tree, where any two statements that are independent can be run in parallel. The linter plugin would use this dependency tree to determine which parts should be parallelised, and save this information to a file.

On the second compilation, the syntax extension would be able to read the file the linter plugin created on the previous compilation. This lists all the changes required and the syntax plugin can apply those changes to the abstract syntax tree function by function. The linter plugin would also be able to view this file, and it could produce compiler warnings for any function that could be parallelised that is missing an annotation.

2.1 Requirements

I believe that I have access to all the required software and hardware for this project. Below is listed most of what I am planning to use:

- A computer to program on
- Source code of some rust programs to parallelise
- Rust compiler (and it's source code)
- A text editor to write code in (Atom)

2.2 Estimated Timeline

Estimated End Date	Milestone
30/10/2017	Create two plugins, a linter and a syntax extension, and get them to communicate via a file
13/11/2017	Get access to the abstract syntax tree inside the linter plugin
27/11/2017	In the linter plugin, analyse each statement to see what variables they depends on and what variables are modified (if any)
11/12/2017	Produce a dependency tree for the entire program based on this analysis
25/12/2017	Look for areas in the tree which do not depend on one another, these areas could be run in parallel
08/01/2018	Save these areas to the shared file used for communication. Make sure that the syntax extension plugin can read this file correctly
22/01/2018	In the syntax extension plugin, first try one statement in parallel to test it works and then try to run everything in parallel
29/01/2018	Run tests against rust programs for program correctness and calculate speedup/slowdown against the sequential version
12/02/2018	In the linter plugin, try to analyse the speed of each parallelisable part. Threads have an overhead so it may be faster to run in sequential for some cases. Record this into the shared file
26/03/2018	In the syntax extension plugin, take into consideration the time information and only run parts in parallel if it would be faster to run in parallel
05/03/2018	Run more tests against rust programs for speedup/slowdown

2.3 Possible Extensions

If I'm ahead of schedule with the coding section of this project, I could look into parallelising branches which have a slow condition (i.e. 'if' statements). Each branch would run isolated in separate thread using cloned data. When the condition is finally calculated the incorrect branches would be terminated. This kind of parallelisation is different from the project plan as some threads are "thrown away".

Another extension upon the previous extension would be to utilise the GPU in cases where the same isolatable task is repeated a large number of times but on different data such that the overhead of moving this task to the GPU is worthwhile. It is very unlikely that I would attempt this extension due to the amount of work required, and I'm unsure about how much real world code exists that could be automatically run on the GPU efficiently.

3 Evaluation

To evaluate whether the chosen approach is the best option to solve the problem, it is worth looking at the downsides to this approach, and look other possible approaches. I will also describe how I will evaluate the final solution to see whether it solves the original problem.

3.1 Disadvantages to the Chosen Approach

The first downside is that the target program would need to be compiled twice due to the limitation of rust compiler plugins as explained section 2. It might be possible to edit the rust compiler to make it so that the program needs to be compiled once, but this has its own downsides as explained in 3.2.1.

The second downside is installing the compiler plugins. To take an existing sequential rust program and compile it so that it runs in parallel would require a few steps. The compiler plugin crate would need to be imported, and every function would need to be annotated so that the syntax extension plugin has edit access to the abstract syntax tree for each function.

3.2 Other Approaches

Below are a few other approaches to the problem that I considered before reaching the chosen approach. Each alternative approach looks at the potential benefits over the chosen one and then explains why I did not choose that alternate approach.

3.2.1 Make a custom rust compiler

One of the downsides of the chosen approach is that the program has to be compiled twice. I may be able to bypass this limitation by making a fork of the rust compiler instead of using compiler plugins. It would then be possible to examine the entire abstract syntax tree and then edit it afterwards. However, this approach seems infeasible for me as the rust compiler code complexity would add too much time to the project.

3.2.2 Use a different language from rust

There are many programs out there, written in many different languages. So why rust? Rust is not in the top 25 of most popular programming languages, according to StackOverflow (2017), so the number of existing programs may smaller than a more popular language. Rust is not unpopular however making it the most loved language and the 10th language on the most want to lean list (StackOverflow 2017). If there are not enough existing programs to justify specialising the tool just for rust (which there probably is), the number programs should increase exponentially in the future. The main reason I chose rust is because it is a safe language with guaranteed memory safety and threads without race conditions which would help significantly with the task.

3.2.3 Manual annotation of the parts of the program that are parallelisable

Instead of analysing the program for parallelisable parts, let the programmer use annotations to show which functions can be run in parallel, similar to OpenMP developed for C++ (Dagum and Menon 1998). This would be significantly simpler than the chosen approach as the analysis and dependency tree steps of the compilation could essentially be removed. The programmer is responsible for telling the compiler what can and cannot be parallelised. If the programmer accidentally makes a function parallel that should not be run in parallel, then the compiler would either need to spot this error (using the analysis step that was removed) or it would create an incorrect program. This approach would also not allow for automatic translation from sequential to parallel code and so it is not that much better than using threads manually.

3.2.4 Just run the sequential code

All of this analysis at compile time (and the fact that it would compile twice in the chosen approach) would increase compile times and running code in multiple threads has some runtime overhead. Is it worth the development time of automating the conversion of sequential code into parallelisable code? This is a valid argument if the number of sequential programs that could benefit from this software is small enough such that it would be faster and easier to redevelop these programs. It must also be true that developing future programs as parallel programs would have the same difficulty as developing it as a sequential program for this argument to hold. Since developing parallel programs is more difficult than developing sequential programs, and there are probably enough programs that could benefit from this software, it is worth the development time to create this piece of software.

3.3 Measuring Project Success

I can look at existing programs written in rust to see if there is any real world impact by measuring the speedup of the parallelised program vs the original sequential code. I can also calculate what percentage of the static sequential code has been parallelised, and what percentage of the runtime is using the parallel code. From these statistics, I can evaluate how well my solution solves the problem.

4 References

- Baskaran, Muthu, Jj Ramanujam and P Sadayappan (2010). “Automatic C-to-CUDA code generation for affine programs”. In: *Compiler Construction*. Springer, pp. 244–263.
- Dagum, Leonardo and Ramesh Menon (1998). “OpenMP: an industry standard API for shared-memory programming”. In: *IEEE computational science and engineering* 5.1, pp. 46–55.
- D’Hollander, Erik H, Fubo Zhang and Qi Wang (1998). “The fortran parallel transformer and its programming environment”. In: *Information sciences* 106.3-4, pp. 293–317.
- Geer, David (2005). “Chip makers turn to multicore processors”. In: *Computer* 38.5, pp. 11–13.
- Kish, Laszlo B (2002). “End of Moore’s law: thermal (noise) death of integration in micro and nano electronics”. In: *Physics Letters A* 305.3, pp. 144–149.
- StackOverflow (2017). *Stack Overflow Developer Survey 2017*. URL: <https://insights.stackoverflow.com/survey/2017> (visited on 09/10/2017).
- The Rust Programming Language* (2017). URL: <https://www.rust-lang.org/en-US/> (visited on 09/10/2017).

Automatic Parallelisation of Rust Programs at Compile Time

Michael Oultram
Student ID: 1428105

Dr Ian Batten
Project Supervisor

Abstract—Processors have been gaining more multi-core performance which sequential code cannot take advantage of. Many solutions exist to this problem including automatic parallelisation at the binary level, automatic parallelisation at the compiler and manually parallelising using annotations. The potential benefits of solving this problem are increased performance for existing programs, as well as making development of new software easier (as programmers do not need to worry about writing parallelised code).

This paper focuses on automatically converting sequential source code into a parallelised program. The literature is explored and presenting into two main areas: theoretical models of automatic parallelisation and existing real-world parallelising compilers. The author presents a design for a new parallelising compiler for the Rust programming language. To demonstrate the design, three sequential elements of Rust programs are manually converted by this paper: function calls, for-loops and branches.

I. INTRODUCTION

Kish (2002) estimated the end of Moore’s Law of miniaturization within 6-8 years or earlier (based on their publication date) and as such, manufacturers have been increasing processors’ core count to increase processor performance (Geer 2005). Writing parallelised programs to take advantage of these additional cores has some difficulty and often requires significant changes to the source code. One solution to this problem, which is the focus of this paper, is to automatically transform sequential source code into parallelised code. This solution, if achieved, would allow for existing uncompiled sequential programs to take advantage of the new hardware. Developing new programs would be easier for programmers as they can just write the code sequentially and let the compiler parallelise it.

Section II examines the literature to finds that most existing solutions focus on unsafe languages such as C++ and FORTRAN. Since a significant proportion of existing programs are written in these languages, it makes sense for other authors to focus on these languages. There is a downsides to focusing on these unsafe languages, dependency analysis becomes much more difficult. To avoid this downside, the design this paper introduces in section IV uses the safe programming language Rust.

“Rust is a systems programming language that runs blazingly fast, prevents segfaults, and guarantees thread safety” (*The Rust Programming Language* 2017). This brief introduction

to some of Rust’s features will explain the elements of the language necessary for the reader to understand for later sections of this paper. For further understanding of the language, it is recommended that the reader looks at *The Rust Book* (2017).

In C++ (or other unsafe languages), a variable can be accessed and modified if it is in scope. In Rust, each variable is immutable by default and if a variable is passed to another function, then that function takes ‘ownership’ of the variable. The code below would not compile in Rust, but similar code would compile in C++. When `f` is called, it takes ownership of the variable `i`. This essentially moves it out of scope, and so it cannot be moved into `g`.

```
fn main() {  
    let mut i = vec![1,2,3];  
    f(i);  
    g(i);  
}
```

To make the code above compile, we have two options. The first option is to let `f` ‘borrow’ the variable instead of moving the ownership of `i` to `f`. In the example below, variable `i` is ‘mutably borrowed’ by `f`, and then ‘immutably borrowed’ by `g`. The syntax of Rust makes it explicit that `f` may modify this variable and `g` cannot. This information expressed is available to the compiler.

```
fn main() {  
    let mut i = vec![1,2,3];  
    f(&mut i);  
    g(&i);  
}
```

The second option is to clone `i` instead. This would keep the types of `f` and `g` the same as the first example. In the example below, we make a copy of `i` before moving it to `f`. Any changes that `f` makes to `i` will not be transferred to `g`, unlike the previous example.

```
fn main() {  
    let mut i = vec![1,2,3];  
    let j = i.clone();  
    f(i);  
    g(j);  
}
```

Note: `i` must be cloned before moving `i` to `f`, or else we would not have access to `i` to clone it.

II. RELATED WORK

A. Parallelisation Models

In this section, we look at theoretical models of automatic parallelism. The static parallelism subsection shows related work where the schedule is fixed and calculated at ‘compile’ time. It is shown how rearranging loop iterations and optimising memory access patterns for multiple threads can increase performance. The speculative parallelism subsection shows related work where the schedule is more flexible. This kind of parallelism tries to run dependent tasks in parallel and detecting when there is a conflict. When a conflict occurs, some parallel thread is ‘undone’ and rerun.

1) *Static parallelism*: Feautrier (1992) describes one model of a parallel program as a set of operations Ω on an initial store, and a partial ordering binary relation Γ also known as a dependency tree. It is shown that this basic model of a parallel is equivalent to affine scheduling, where Ω and Γ are described as linear inequalities. Finding a solution where these linear inequalities hold produces a schedule for the program where dependent statements are executed in order. There are some programs where no affine schedule exists. Bondhugula et al. (2008) uses the affine scheduling model on perfectly, and imperfectly nested loops. They describe the transformations needed to minimise the communication between threads, further increasing the performance of the parallelised code.

An alternative method to affine scheduling is iteration space slicing introduced by Pugh and Rosser (1997). “Iteration space slicing takes dependency information as input to find all statement instances from a given loop nest which must be executed to produce the correct result”. Pugh and Rosser (1997) shows how this information can be used to transform loops on example programs to produce a real world speedup. Beletskaya et al. (2011) shows that iteration space slicing extracts more coarse-grained parallelism than affine scheduling.

2) *Speculative parallelism*: Zhong et al. (2008) shows that there is some parallelisable parts hidden in loops that affine scheduling and iteration space slicing cannot find. They propose a method that runs future loop iterations in parallel with past loop iterations. If a future loop iteration accesses some shared memory space, and then a past iteration modifies that same location, the future loop iteration is ‘undone’ and restarted. It is shown that this method increases the amount of the program that is parallelised.

Prabhu, Ramalingam and Vaswani (2010) introduce two new language constructs for C# to allow the programmer to manually specify areas of the program that can be speculatively parallelised. Yiapanis, Brown and Luján (2016) designs a parallelising compiler which can automatically take advantage of speculative parallelism.

B. Parallelisation Implementations

In this section we look at: parallelising compilers which focus on parallelising FORTRAN programs; OpenMP which is an model for shared memory programming and parallelising compilers which convert sequential CPU code into parallelised GPU code. Some of these parallelising compilers are based off of models described in subsection II-A.

Eigenmann, Hoeflinger and D. Padua (1998) manually parallelises the PERFECT benchmarks for FORTRAN which are compared with the original versions to calculate the potential speedup of an automatic parallelising compiler. D’Hollander, Zhang and Wang (1998) developed a FORTRAN transformer which reconstructs code using GOTO statements so that more parallelisms can be detected. It performs dependency analysis and automatically parallelised loops by splitting the task into jobs. These jobs can be split between networked machines to run more jobs concurrently. Rauchwerger and D. A. Padua (1999) introduce a new language construct for FORTRAN programs which allows for run-time speculative parallelism on for loops. Their implementation parallelises some parts of the PERFECT benchmarks which existing parallelising compilers of the time could not find.

Quiñones et al. (2005) introduce the Mitosis compiler which combines speculation with iteration space slicing. There is always only one non-speculative thread which is seen as the base execution; all other threads are speculative. The Mitosis compiler computes the probability of two iterations conflicting. If this probability is low, and there is a spare thread unit, then the loop iteration is executed in parallel. The non-speculative thread detects any conflicts as it is the only thread that can commit results.

Dagum and Menon (1998) introduces a programming interface for shared memory multiprocessors called OpenMP targeted at FORTRAN, C and C++. The programmer annotates the elements of the program that are parallelisable, which the compiler recognises and performs the optimisation. OpenMP is compared to alternative parallel programming models. Kim et al. (2000) introduces the ICP-PFC compiler for FORTRAN which uses the OpenMP model. All loops in the source code are analysed by calculated a dependency matrix. The compiler automatically adds the relevant OpenMP annotations to the loop. Lam (2011) extends OpenMP using machine learning to automate the parallelisation. The system is trained using a set containing programs already parallelised using OpenMP. The knowledge learned is applied to sequential programs to produce parallelised programs.

A CPUs architecture is typically optimised for latency whereas a GPUs architecture is typically optimised for throughput. This can make GPUs perform much better than CPUs for a certain type of task. Baskaran, Ramanujam and Sadayappan (2010) uses the affine transformation model to convert sequential C code into parallelised CUDA code. For loops are tiled for efficient execution on the GPU.

III. PROBLEM DETAILS

The literature focuses on unsafe languages such as C/C++ and FORTRAN. As a result, most of their methods revolve around understanding conflicts between statements/loop iterations of complex code. Since memory in Rust has ‘ownership’, this dependency information is more readily available.

The Rust compiler allows for plugins of different types but there are two types that are of interest to this problem. Written in the order of execution, they are:

- Syntax Extension: can modify the abstract syntax tree of any annotated function.
- Early Lint Pass: can see abstract syntax tree of each uncompiled function, with macros expanded, without annotations, but it cannot edit them.

Before we can automate the parallelisation of Rust programs, we must look at some example programs to fully understand the problem. The examples were written sequentially and then manually parallelised. The pattern of parallelisation could be applied to different examples, and a design is proposed to automate this in section IV. In these examples, threads are used as if they have no overhead, to simplify the explanation of how a sequential program could be parallelised. In reality this is not the case and so in section IV `thread::spawn` is replaced with something more optimal.

A. Parallel Function Optimisations

In some cases, the function arguments are known before the result of the function is required. When the program is sequential, the program must wait for the result of that function call. If the program was parallelised, this function could be started in the background as soon as the arguments are decided. This would allow the result to be ready for when it is requested or at least closer to ready than the pure sequential version.

Example: Algorithm 1 is a program that calculates Fibonacci numbers in a very inefficient way. The main method calculates `i` on the second line, and since it is not mutable it cannot be changed. The `slow_method` is executed before the `fibonacci` function, even though the two functions are independent. Just running the `fibonacci` function earlier would not improve performance as the `slow_method` would now have to wait for `fibonacci` function to finish. Ideally, `fibonacci` should start as soon as `i` is calculated, and at the same time `slow_method` should be executed.

Algorithm 2 shows one way of converting Algorithm 1 into a more parallelised version. These changes allow for `fibonacci_parallel` to be called as soon as `i` is decided. `fibonacci_parallel` does not block, and will return almost immediately with a `JoinHandle`. The inner thread is executed in the background, allowing `slow_method` to be executed at the same time. Inside the

Algorithm 1 Sequential Fibonacci Function

```
fn main() {  
    // Set i = first program argument  
    let program_args = std::env::args();  
    let i = program_args[0].trim().parse::<u32>();  
    slow_method();  
    println!("{}", fibonacci(i));  
}  
  
fn fibonacci(n: u32) -> u64 {  
    match n {  
        0 => 0,  
        1, 2 => 1,  
        _ => fibonacci(n-1) + fibonacci(n-2),  
    }  
}
```

Algorithm 2 Parallel Fibonacci Function

```
fn main() {  
    // Set i = first program argument  
    let program_args = std::env::args();  
    let i = program_args[0].trim().parse::<u32>();  
    let fib = fibonacci_parallel(i);  
    slow_method();  
    println!("{}", fib.join());  
}  
  
fn fibonacci_parallel(n: u32) -> JoinHandle<u64> {  
    thread::spawn(move || {  
        match n {  
            0 => 0,  
            1, 2 => 1,  
            _ => {  
                let n1 = fibonacci_parallel(n-1);  
                let n2 = fibonacci_parallel(n-2);  
                n1.join() + n2.join();  
            },  
        }  
    })  
}  
  
fn fibonacci(n: u32) -> u64 {  
    let fib = fibonacci_parallel(n);  
    fib.join()  
}
```

background thread `fibonacci_parallel` is called again which spawns another thread. These changes also allow for `n1` and `n2` to be executed in parallel. Calling `join()` on a `JoinHandle` will block until the thread is terminated, and returns the result of the thread. The `fibonacci` function is changed so that any external code that uses this function will still work (although sequentially).

B. For-Loop Optimisations

Loop iterations that are independent of each other could be run at the same time to increase performance. It is not always clear how many iteration of a loop may occur, especially while-loops. For-loops are more explicit about how many iterations will occur, but they can still terminate early with a `break` or a `return` statement.

Example: Algorithm 3 is a real world example with a two level for-loop combined with an if-statement which returns the password for the first valid hash. The inner for-loop is

not parallelisable as each iteration reads and modifies the `hash_word` variable. The outer for-loop is parallelisable, but if we were to naively put the contents into separate threads, as shown in Algorithm 4, then we may end up with the wrong result. In the sequential version, the first password in the list to match the hash would be returned but in the naive parallel version, the password returned depends on the order the threads are run. In this example program at most only one iteration could possibly enter the if-branch, terminating the for-loop early, but in other examples there could be more than one. Algorithm 5 is the final complete parallelisation of Algorithm 3. Each word is hashed in it's own thread, compared to `hash_password` and returned to be stored in the initial thread. The function does not return a result until all previous iteration threads have also returned.

Algorithm 3 Sequential Password Cracker

```
fn crack_password(dictionary: &Vec<String>,
    hash_password: String)
    -> Option<String> {
    for word in dictionary {
        // Hash word using Sha256
        let mut hash_word = word;
        for _ in 0..40 {
            let mut sha = Sha256::new();
            sha.input_str(word);
            hash_word = sha.result_str();
        }
        // Check if hash matches
        if hash_password == hash_word {
            return Some(word.clone());
        }
    }
    // No hash matched
    None
}
```

C. Branch Optimisations

In the previous optimisations, all the code that is run in parallel would have also been run in sequential normally. This optimisation more speculative than the other optimisations described as some threads are going to be disregarded. Each side of the branch for an if-statement could be run in parallel whilst waiting for the condition. When the condition is finally calculated, the result of the correct branch is kept. Algorithm 7 shows this optimisation applied to Algorithm 6. Ideally, when the condition is calculated and the incorrect branch is still being calculated in a thread, then this thread should be cancelled. This is not easily achievable in Rust so this optimisation is left as a possible extension of the design.

IV. DESIGN OVERVIEW

A typical parallelising compiler, such as those described in section II, have a few stages. First the compiler looks at each statement of the source code, and performs dependency analysis to calculate the critical path. Any independent statements can be run in parallel. A scheduling algorithm calculates which order the statements should be executed, grouping statements that are dependent on each other into

Algorithm 4 Naive Parallel Password Cracker

```
fn crack_password(dictionary: &Vec<String>,
    hash_password: String)
    -> Option<String> {
    // Create a communication channel
    let (tx, rx) = mpsc::channel();
    // Start a thread for each dictionary entry
    for i in 0..dictionary.len() {
        let word = dictionary[i].clone();
        let tx = tx.clone();
        thread::spawn(move || {
            let result = {
                // Hash word using Sha256
                let mut hash_word = word;
                for _ in 0..40 {
                    let mut sha = Sha256::new();
                    sha.input_str(word);
                    hash_word = sha.result_str();
                }
                // Send result via channel
                if hash_password == hash_word {
                    Some(word)
                } else {
                    None
                }
            };
            tx.send(result);
        });
    }
    // Receive up to dictionary.len() results
    for _ in 0..dictionary.len() {
        if let Some(result) = rx.receive() {
            return result;
        }
    }
    // No hash matched
    None
}
```

the same task. A compile time performance metric is used on each task to estimate the potential speedup of the parallel version over the sequential version. Due to the overhead of threads, this potential speedup would not actually be achieved. The parallelising compiler takes this into account and will only use the parallel version if it predicts it will really be faster.

This papers design takes elements from a typical parallelising compiler, dividing the task into two main stages, an analysis stage and a modification stage. The analysis stage is run by a linter plugin and the modification stage is run by a syntax extension plugin. Analysis stage must come before the modification stage, so compiling must be done twice (once for each stage). When the plugin is loaded, it determines which stage it is by looking for a shared JSON file. If this file does not exist, then it is the analysis stage, otherwise the file's content is loaded and the stage is updated to be the modification stage.

A. Analysis stage

On the first compilation, the syntax extension plugin would do nothing. The linter plugin would view the entire abstract data tree and analyse what each statement depends on and modifies. This would create a dependency tree, where

Algorithm 5 Parallel Password Cracker

```
fn crack_password(dictionary: &Vec<String>,
    hash_password: String)
    -> Option<String> {
    let (tx, rx) = mpsc::channel();
    for i in 0..dictionary.len() {
        let word = dictionary[i].clone();
        let tx = tx.clone();
        thread::spawn(move || {
            let result = {
                // Hash word using Sha256
                let mut hash_word = word;
                for _ in 0..40 {
                    let mut sha = Sha256::new();
                    sha.input_str(word);
                    hash_word = sha.result_str();
                }
                // Check if hash matches
                if hash_password == hash_word {
                    Some(word)
                } else {
                    None
                }
            };
            tx.send((i, result));
        });
    }
    // Receive all the results
    let mut results = vec![None; list.len()];
    let mut verified_upto = -1;
    for _ in 0..results.len() {
        // Receive result and store in location
        let (i, result) = rx.receive();
        results[i] = Some(result);
        // Check for final result
        for i in 0..results.len() {
            if let Some(result) = results[i] {
                if let Some(word) = result {
                    return word;
                }
            } else {
                // Have not received i result yet
                break;
            }
        }
    }
    // No hash matched
    None
}
```

Algorithm 6 Sequential Slow If

```
fn f(a: u32, b: u32) -> u32 {
    if slow_condition(a, b) {
        (a * (b - a)) + 5
    } else {
        ((a * b) + 19) ^ 2
    }
}
```

any two statements that are independent can be run in parallel. The linter plugin would use this dependency tree to determine which parts should be parallelised, and save this information to a file. Some parts cannot be parallelised.

1) *Side Effects*: A pure function will return the same output if it is given the same input. If the function can return something different when given the same input, then it is not pure. The element of the function that changes the output is known as a side effect. The main side effect found in

Algorithm 7 Parallel Slow If

```
fn f(a: u32, b: u32) -> u32 {
    let true_branch = {
        let a = a.clone();
        let b = b.clone();
        thread::spawn(move || (a * (b - a)) + 5)
    };
    let false_branch = {
        let a = a.clone();
        let b = b.clone();
        thread::spawn(move || ((a * b) + 19) ^ 2)
    };
    if slow_condition(a, b) {
        true_branch.join()
    } else {
        false_branch.join()
    }
}
```

programs is IO. For example, a function takes a filename as input and outputs the contents of that file. If the contents of that file changes, the function would outputs something different even though its parameters are the same. Printing to standard out is also considered a side effect. Side effects generally cannot be parallelised, as the order of execution must be kept the same (we don't want to print in the wrong order). Additionally, everything that occurs before the side effect in the sequential version should also be run before the side effect in the parallel version.

B. Modification stage

On the second compilation, the syntax extension reads the file the linter plugin created on the previous compilation. This lists all the changes required and the syntax plugin can apply those changes to the abstract syntax tree function by function. The modifications applied would be similar to the sample programs shown in section III. Those modifications assumed that threads have no overhead and were used as a method of describing how a sequential program could be run in parallel. In the real world, these parallelised sample programs produce an unreasonable amount of threads and most of the threads are just waiting for other threads to return. This synchronisation can cause the parallel version to perform a lot worse than sequential code. The following subsections describe some of the problems with the parallelised sample programs, and provides some solutions.

1) *Infinite Thread Problem*: The initial solution for too many threads is to use a thread-pool so only a fixed number of threads are executed at the same time. This solution would not work in this case as if all the threads are waiting on a future task, which is also waiting for a thread to be free then we get stuck in a deadlock. By tweaking the design of the thread-pool slightly, we can have a fixed the number of tasks and prevent this deadlock.

As the task queue is the main cause of this deadlock, it is removed from this new proposed design which will be referred to as a no-queue thread-pool. Whenever a task is created, it looks for an available thread and that thread starts

executing the task in the background. If there is no thread available, then the current thread should execute the task. Also, if the the current thread is waiting for a task to return, then it could execute another task whilst waiting.

To modify the sample programs to use the no-queue thread-pool would be of minimal work; it would require a shared memory space (to gain access to the threads) and instead of calling `thread::spawn`, it would call another function.

2) *Performance Analysis*: Parallelising small tasks, or tasks that would require a lot of synchronisation when run in parallel actually run significantly slower than the sequential code. In these cases, the code should not be parallelised. There is a problem though, how does the compiler know if it is faster parallelised or sequential.

One possible solution is to run the sequential version in one thread, and the parallel version in the other threads. Whichever version finishes first is used, and the other version is cancelled. This approach leaves the performance approach until run-time, but it is a very inefficient solution as the parallel version now has fewer threads.

Another solution is to use a performance metric at compile time to estimate which version will be faster. If we assume each statement takes one unit of execution time, we can attempt to calculate how many time units the sequential and parallel version will take. Since we know that there is an overhead to the parallel version, we should add on some extra time units. The faster version would be the version that is compiled. This solution would not really work as described as the execution time may be dependent on how many iterations of a loop (which isn't necessarily known at compile time). It is also dependent on the number of cores the machine has, also not known at compile time.

The best solution is a combination of the previous two ideas. The compiler should create an function which estimates both the sequential and the parallel version execution time. This function can take in any run-time conditions i.e. the number of iterations of a loop or the number of cores the machine has. This function should be fast such that the overhead of calculating which version is faster should be negligible. At run-time this function predicts the faster version using the run-time variables, and runs only that version.

V. REFERENCES

- Baskaran, Muthu, Jj Ramanujam and P Sadayappan (2010). "Automatic C-to-CUDA code generation for affine programs". In: *Compiler Construction*, pp. 244–263.
- Beletskaya, Anna et al. (2011). "Coarse-grained loop parallelization: Iteration space slicing vs affine transformations". In: *Parallel Computing* 37.8, pp. 479–497.
- Bondhugula, Uday et al. (2008). "Automatic transformations for communication-minimized parallelization and locality optimization in the polyhedral model". In: *International Conference on Compiler Construction*, pp. 132–146.
- D'Hollander, Erik H, Fubo Zhang and Qi Wang (1998). "The fortran parallel transformer and its programming environment". In: *Information sciences* 106.3-4, pp. 293–317.
- Dagum, Leonardo and Ramesh Menon (1998). "OpenMP: an industry standard API for shared-memory programming". In: *IEEE computational science and engineering* 5.1, pp. 46–55.
- Eigenmann, Rudolf, Jay Hoeflinger and David Padua (1998). "On the automatic parallelization of the Perfect Benchmarks (R)". In: *IEEE Transactions on Parallel and Distributed Systems* 9.1, pp. 5–23.
- Feautrier, Paul (1992). "Some efficient solutions to the affine scheduling problem. I. One-dimensional time". In: *International journal of parallel programming* 21.5, pp. 313–347.
- Geer, David (2005). "Chip makers turn to multicore processors". In: *Computer* 38.5, pp. 11–13.
- Kim, Hong Soog et al. (2000). "ICU-PFC: An automatic parallelizing compiler". In: *Proceedings - 4th International Conference/Exhibition on High Performance Computing in the Asia-Pacific Region, HPC-Asia 2000*, pp. 243–246.
- Kish, Laszlo B (2002). "End of Moore's law: thermal (noise) death of integration in micro and nano electronics". In: *Physics Letters A* 305.3, pp. 144–149.
- Lam, Nam Quang (2011). "A Machine Learning and Compiler-based Approach to Automatically Parallelize Serial Programs Using OpenMP". In: *Master's Projects* 210.
- Prabhu, Prakash, Ganesan Ramalingam and Kapil Vaswani (2010). "Safe programmable speculative parallelism". In: *ACM Sigplan Notices* 45.6, pp. 50–61.
- Pugh, William and Evan Rosser (1997). "Iteration space slicing and its application to communication optimization". In: *Proceedings of the 11th international conference on Supercomputing*. ACM, pp. 221–228.
- Quiñones, Carlos García et al. (2005). "Mitosis compiler: an infrastructure for speculative threading based on pre-computation slices". In: *ACM Sigplan Notices* 40.6, pp. 269–279.
- Rauchwerger, Lawrence and David A. Padua (1999). "The LRPD test: Speculative run-time parallelization of loops with privatization and reduction parallelization". In: *IEEE Transactions on Parallel and Distributed Systems* 10.2, pp. 160–180.
- The Rust Book* (2017). URL: <https://doc.rust-lang.org/book/> (visited on 12/11/2017).
- The Rust Programming Language* (2017). URL: <https://www.rust-lang.org> (visited on 12/11/2017).
- Yiapanis, Paraskevas, Gavin Brown and Mikel Luján (2016). "Compiler-Driven Software Speculation for Thread-Level Parallelism". In: *ACM Transactions on Programming Languages and Systems (TOPLAS)* 38.2, pp. 1–45.
- Zhong, Hongtao et al. (2008). "Uncovering hidden loop level parallelism in sequential applications". In: *High Performance Computer Architecture, 2008. HPCA 2008. IEEE 14th International Symposium on*, pp. 290–301.