



Automatic Parallelisation of Rust Programs at Compile Time

Author

Michael Ultram
1428105

Degree

MSci in Computer Science
28th March 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	Related Work	4
2.1.1	Parallelisation Models	4
2.1.2	Parallelisation Implementations	5
2.2	My Solution using Rust	6
2.2.1	Safety Features	6
2.2.2	Threads	7
2.2.3	Crates	7
3	Design	8
3.1	Analysis Stage	8
3.2	Modification Stage	9
4	Implementation	10
4.1	How to implement the design	10
4.2	Dependency Analysis	11
4.3	Modification Stage	11
4.3.1	Shared File	11
4.3.2	Scheduler	12
4.3.3	Reconstructor	12
5	Evaluation	13
5.1	Sequential Programs Parallelised	13
5.1.1	Simple Example	13
5.1.2	Password Cracker	14
5.2	Limitations/Improvements	15
6	Discussion	16
6.1	Achievements	16
6.2	Improvements	16
6.3	Approach Changes	16
7	Conclusion	17
8	Appendix	18
8.1	References	18
8.2	Submission File Structure	19
8.3	Running the Code	19
8.4	Rust Compiler Types	19

Chapter 1

Preamble

1.1 Abstract

TODO: Write Abstract

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

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

2.1 Related Work

Processors are being released with more cores. Sequential code cannot take advantage of these new cores. Writing parallel code is more difficult than writing sequential code. Existing sequential code would need to be rewritten to be parallelised. Ideally we want to gain the benefits of parallel code, whilst only having to write easy sequential code.

TODO: Mention that it is a well researched area One solution to this problem is a parallelising compiler. Research done in this field has mostly focused on the C/C++ language although other researches have had success using other languages. Some methods require manual annotation of the source code by the programmer to specify which parts of the program are parallelisable. Others have attempted to automatically detect these areas, but with an unsafe language like C++ it is challenging.

2.1.1 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.

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 state-ment 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.

Speculative parallelism

Zhong et al. (2008) shows that there is some parallelisable parts hidden in loops that affine scheduling and iteration space splicing 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.

2.1.2 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 2.1.1.

Eigenmann, Hoefflinger 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.

2.2 My Solution using Rust

Due to the limited time of the project, I am unlikely to contribute new ideas to the field. **TODO: Add link here.** For my project, I focused on the safe language rust. Rust has a unique way of managing memory such that only one thread can access the memory space at once. This is guaranteed at compile time, and should make the process of automatically detecting dependencies much easier. Rust also allows plugins into the compiler (nightly feature only as of writing) which gives modification access to the abstract syntax tree.

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. If the reader requires more in depth understanding than what is provided, then they should look at the language documentation, *The Rust Book* (2017).

2.2.1 Safety Features

“Rust is a systems programming language that runs blazingly fast, prevents segfaults, and guarantees thread safety” (*The Rust Programming Language* 2017). To get these safety properties, rust has some unique features. The biggest difference to other programming languages is how variable are handled.

Ownership

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
```

Listing 1: Borrowing and moving example

In this example, `a` is a local variable in the `main` method.

When `f` is called with parameter `a`, the function borrows that variable. This is similar to call-by-reference from other programming languages.

When `g` is called with parameter `a`, the variable is moved to `g`. This is unlike other programming languages as this is not call-by-value. Instead `g` takes ownership of `a`. When `g` is returned, the `main` method can no longer use `a`.

Mutability

Variables mutability is declared when the variable is declared. In rust, variables are immutable by default, but if specified they are mutable. When a variable is borrowed, it can either be immutably borrowed or mutably borrowed.

```
1 fn main() {
2     let a = 10;
3     let mut b = 20;
4     f(&a, &b);
5     g(&mut b);
6 }
7 fn f(a: &u32, b: &u32){} // f immutably borrows a and b
8 fn g(b: &mut u32){} // g mutably borrows b
```

Listing 2: Immutable and mutable borrowing

In the `main` method of this example, `a` is an immutable local variable and `b` is a mutable local variable. The `f` function borrows both `a` and `b` immutably. Even though `b` is declared as mutable, it cannot be changed inside `f`. Once `f` returns, `b` becomes mutable again inside the `main` method. The `g` function shows how `b` can be borrowed mutably.

Unsafe Blocks

The programmer may want can turn off some of rust's safety features by using an unsafe block. The most common use of an unsafe block is to modify a mutable static variable but it also allows de-referencing of a raw pointer and calling unsafe functions (i.e. an external c function). Using unsafe blocks may introduce race conditions as two threads could try to modify a global at the same time, and the rust language would not guarantee an order.

```
1 static mut global: u32 = 3;
2 fn main() {
3     let a = global;
4     inc_global();
5     unsafe {
6         global = 5;
7     }
8 }
9
10 unsafe fn inc_global() {
11     global += 1;
12 }
```

Listing 3: Immutable and mutable borrowing

2.2.2 Threads

Rust uses real threads. Due to memory design, only one thread can have access to a variable safely at once. Channels are used to communicate between threads and can move a variable from one thread to another. This variable must implement the `Send` trait.

2.2.3 Crates

Chapter 3

Design

This chapter covers the original design concepts of my parallelising compiler without realising the full internal details of the compiler. Once I began developing the compiler, I realised some mistakes in the design which I did not fully think through and some adaptations that I had to do due to the structure of the rust compiler. These changes are described in chapter 4. Also due to time constraints, not all the features described in this chapter were implemented.

The rust abstract syntax tree consists of three main types: Blocks, Statements and Expressions. A block contains a list of statements and a statement is a combination of expressions. A block can be represented as an expression (either directly, or inside a loop/if/etc.) which allows for infinite depth in the tree. Variables are a type of expression. My parallelising compiler will focus on statement level parallelisation within a given block.

Parallelising compilers described in the literature (section 2.1) are split into several stages. Similarly, the design of my compiler is in two main stages, the analysis stage and the modification stage. Both these stages contain several steps to achieve their goal. The analysis stage looks at each function in the source code and tries to find parts that could be parallelised. The modification stage takes the parts that can be parallelised and changes the source code so that they run in parallel.

3.1 Analysis Stage

Algorithm 1 Dependency Analysis Algorithm

Require: *block* as a list of Statements

```
1: envs = {}
2: deps = {}
3: block_env = []    ▷ Statement Environment of the Block when represented as a Statement
4: for stmt in block do
5:   stmt_env = Set of variables that stmt uses
6:   envs.push(stmt_env)
7:   stmt_deps = {}
8:   for i = length(stmt_env) - 1 to 0 do           ▷ Search backwards for vars in stmt_env
9:     dep_env = envs[i]
10:    for all var where (var is in stmt_env) and (var is in dep_env) do
11:      stmt_deps.push(i)
12:      stmt_env.remove(var)
13:   block_env.push(stmt_env)    ▷ Unmatched variables are part of the block's environment
```

Each statement is analysed individually to provide a list of variables that the statement accesses. This list describes the variable dependencies that the statement has, but it does not describe which statements must be executed before the current one for the program to remain correct. To get this information, the algorithm looks at each statement in turn. For each variable that

it requires, it looks at the statements before it in the block in reverse order for that variable to be its list of variable dependencies. The first statement found containing this variable as a dependency is added as a statement dependency. Any variables that were not found above the current statement must be defined outside the current block, and so are added as the current blocks dependency. This algorithm is shown in Algorithm 1.

3.2 Modification Stage

The dependency tree provided by the analysis stage shows what statements can be run in parallel. Some of these statements have multiple dependencies, all of which must be met before the statement is run. Each of these dependencies could be in a separate thread, and so some synchronisation technique is needed. The dependency tree is converted into a schedule tree so that we know which statements are run in which threads, and where/when synchronisation is required between threads.

Algorithm 2 Scheduling Algorithm

Require: *block* as a list of Statements

Require: *envs* as a list Statement Environments

Require: *deps* as a list Statement Dependencies

```

1: schedule_trees = []
2: for i = 0 in length(block) - 1 do           ▷ Add all independent statements to separate trees
3:   if length(deps[i]) == 0 then
4:     schedule_trees.push(("Run", i, []))
5: while not all Statements from block are in schedule_trees do
6:   for stmtid = 0 to length(blocks) - 1 where block[stmtid] is not in schedule_trees do
7:     if all deps[stmtid] are in schedule_trees then
8:       dep_trees = find all ("Run", i, -) in schedule_trees for all i in deps[stmtid]
9:       sort descending dep_trees and deps[stmtid] by depth in schedule_trees
10:      for i = 1 to length(dep_trees) do
11:        (-, -, subtrees) = dep_trees[i]
12:        subtrees.push(("SyncTo", stmtid))
13:      (-, -, subtrees) = deps[0]
14:      subtrees.push(("SyncFrom", deps[stmtid][1 :], [("Run", stmtid, [])]))

```

The scheduling algorithm designed aims to run as much as possible in separate threads. It makes the naive assumption that threads have no overhead. To start of with, the algorithm looks for any statements with no dependencies. Each of these statements can be run in a separate thread. For all the remaining statements, the algorithm looks selects the set of statements that have all their dependencies in the schedule already. If the statement only has one dependency, then that statement is added directly after that dependency. If the statement has more than one, then the algorithm looks to see which dependency has the longest chain. The thought behind this is that this dependency should be the slowest, and so all the other dependencies should have been finished by this point. To make sure that there are no race conditions, a syncline is introduced from the other dependencies to just before the current statement. This algorithm is repeated until all the statements are in the schedule. Since there cannot be any cyclic dependencies due to the way that the dependency analysis algorithm was designed, this is guaranteed to terminate. This algorithm is shown in Algorithm 2.

Chapter 4

Implementation

This chapter looks at how the design was implemented in practice and the design decisions that had to be adapted due to unforeseen complexities. Each design change is justified with an example of why the original design fails, and alternatives that were considered.

4.1 How to implement the design

There were three choices on how I could implement the design: directly modifying the rust compiler source code and recompiling the compiler; using the rust compiler plugin system to modify the abstract syntax tree (AST) or writing a source to source translation from scratch. Modifying the rust compiler would give me the flexibility to change any part of the compiler that I needed but it would make seeing my individual contributions very difficult. Also, the compiler itself is very large and complex; it would take a while to compile from a clean state. Using a rust compiler plugin would give me less access, but I would only need to compile my plugin. The complexity of the compiler is still there with this option. Writing a source to source translation system from scratch would allow me to avoid touching the rust compiler and its complexity. In return, I would have to write code to extract the AST from a source file. I would have to model the ownership/borrowing information to detect when parallelisms properly. From all these choices, I decided to write a rust compiler plugin as it provides all the ownership information is already accessible. This option requires me to use the full AST; the other options had the possibility of using the less verbose high-level intermediate representation (HIR).

The rust compiler allows for plugins of different types. The two types of plugins of interest are Syntax Extension plugins and Linter plugins. Syntax Extension plugins are executed first in the rust compiler pipeline, and are generally used to convert macros into code. Linter plugins are run at a later stage, and are generally used to check code style to produce warnings (like unused variable). It has the complete AST of the code with all the macros expanded. All the information required about dependencies is accessible inside a linter plugin. However, once the rust compiler gets to the linter plugins, the AST can no longer be modified. Syntax extension plugins allow modification of the AST, but they do not have macros expanded. Some dependencies could be missed by trying to analyse the AST at this stage. The solution I decided on was to compile the program twice. On the first compile, the syntax extension does nothing and the linter plugin examines the expanded code. The dependency information gathered is saved into a file for the next compile. On the second compile, the syntax extension plugin reads the file to get all the dependency information. Any parallelisable parts are then modified to be run in parallel.

To use the rust plugin system, the source code will need to be annotated which provides the plugins with access to that element in the AST. Each function of the sequential source code should be annotated with `#[autoparallelise]`. This annotation does not provide any information about the parallelisability of the source code and is purely just a workaround for the rust plugin system.

4.2 Dependency Analysis

In the design section of the report, the dependency analysis algorithm takes in a block and examines it statement by statement to extract out the variables that are used. Actually doing this in the compiler plugin required a lot more effort than it seems on the face of it. The compiler calls a function in the plugin for each annotated function. To get the `Block` struct out of this requires expanding the `FnKind` enum. A `Block` contains a list of statements (`Stmt`). The `Stmt` struct contains a `StmtKind` enum. This describes what kind of statement it is, either a local variable binding, an unexpanded macro, or an expression ending with or without a semicolon. For the variable binding, a `Local` struct is given which contains a `Path` representing the variable name. Some variables will be assigned a value at creation and this is represented as an `Expr` struct. Unexpanded macros contain some other types which represent the arguments to the macro, but I did not end up going any deeper into this. Statements which are expression with or without a semicolon will give an `Expr` struct. Similar to the `Stmt` struct, the `Expr` struct contains a `ExprKind` enum which represents the 39 different types of expressions. Most of these options contain another `Expr`, but some contain other types like a `Pat` which is the second way a variable can be represented. Almost all of the different cases in `ExprKind` had to be dealt with to fully explore the tree and extract all the variables a statement uses (its environment).

Each statement is converted into our representation of the AST so that information can be stored about the dependencies. The original design used two types, a `expr` to represent a single statement and a `block` to represent a list of `exprs` and `blocks`. This design does not deal a `ExprKind` which contains a `Block` properly (i.e. a for loop). My first idea was to split the `Stmt` into a `expr` and a `block` and set the `block` to depend on the `expr`. This seemed to work for the dependency analysis stage, but not for the reconstruction stage. It was difficult to distinguish which `block` is the contents of the `ExprKind` if an external `block` was also dependent on the `expr`. I created a third type called `exprblock` which contains one `expr` and a list of `blocks` (as if-else statements will require more than one `block`). Once all the statements for a function are converted into our representation, and each statement had an environment, then the environments could be matched up as described in Algorithm 1. The statement ID stored is relative to the `block` of all of the dependencies.

TODO: Two environments so two lets with the same variable are not dependent. Originally only the dependency ids were going to be added to our representation of the AST, but I later realised that I should also store the environment. Later still I realised that just using one environment had a flaw which includes two environments containing which variables the statement requires and produces.

4.3 Modification Stage

4.3.1 Shared File

TODO: Shared file between stages. `StatementID` that is consistent between compiles. Tried merging dependencies then replacing dependencies

Once all functions have been analysed, the `DependencyTree` is converted into an `EncodedDependencyTree`. The code part of the converted statement is converted into a statement id. The statement ID is represented as a pair of numbers (`span.lo().0`, `span.hi().0`), which relate to the byte location of the source code. This will remain consistent between compile runs, whereas the `NodeID` does not. All `EncodedDependencyTrees` and function meta data is stored into a JSON file using `serde_json`.

The plugin detects the JSON file and loads it. This is stored as a shared state between different functions.

The first part of the dependency analysis is repeated for the syntax extension plugin this time. In later section of the design, we need access to the pure AST. There is no (easy) way for the AST to be store into a JSON file and recreated into structs that I could find. The reason that we use the linter plugin is so that we can see the dependencies hidden inside macros.

The dependencies are merged function by function from the shared state so that unexpanded macros get the missing dependencies. The dependencies of Statements that have the same StmtID are merged together.

4.3.2 Scheduler

Once the dependency analysis is complete, the scheduler takes the dependency tree and works out a schedule. All relative dependency ids are converted into StmtIDs. The idea around the scheduling algorithm is Maximum Spanning Trees. All statements that have no dependencies can be started at the very beginning. All remaining statements wait for all their dependencies to be put into the schedule. Once all the dependencies for a statement are added, this statement is selected as the next one to add to the schedule. As each statement requires all of its dependencies before it can be executed, it should be scheduled to run after the slowest dependency. This should minimise the amount of time that the statement has to wait for its dependencies. Synclines are created for all the remaining dependencies so that all dependencies are met. Each block gets its own schedule.

TODO: Add environment to syncline

4.3.3 Reconstructor

TODO: Reconstructor. Struggle of creating new statements with specific StatementIDs. Compiler bug that requires me to recompile for a third time **TODO:** Return types All parts of the reconstructor algorithm takes part in the second compile.

Chapter 5

Evaluation

TODO: Rewrite this entire paragraph. Evaluating the success of my solution is difficult. My implemented solution only works for some examples, and its focus is on parallelising as much as possible, not necessarily getting a faster program. Also, many research papers in the field compare their speedup to LLVM, but rust (like many new languages) uses LLVM itself. Making a pure comparison between the two is impossible, and so my result may be skewed.

5.1 Sequential Programs Parallelised

5.1.1 Simple Example

```
simple-example/main.rs
1  #![feature(plugin)]
2  #![plugin(auto_parallelise())]
3
4  #[autoparallelise]
5  fn main() {
6      let mut a = 4;
7      let mut b = 3;
8      a += 1;
9      b += 1;
10     println!("{}", a, b);
11     println!("End of program");
12 }
```

Listing 4: A simple example program

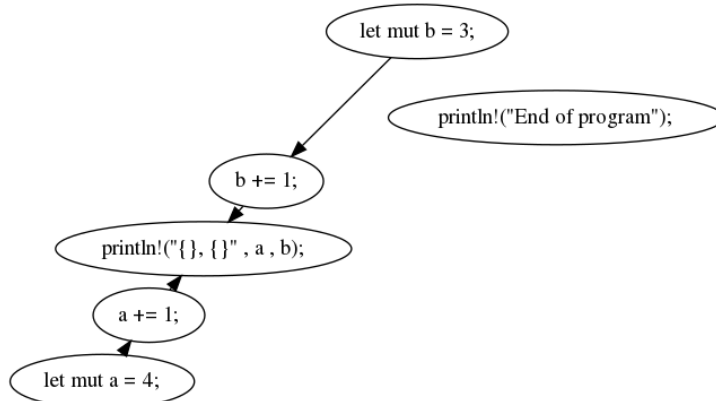


Figure 5.1: Dependency analysis graph of the simple example

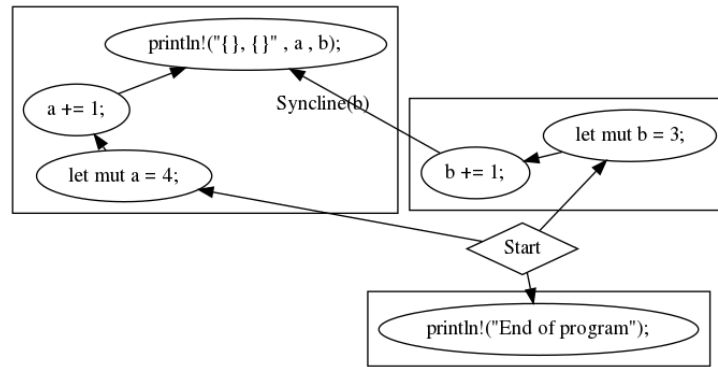


Figure 5.2: Schedule of the simple example

5.1.2 Password Cracker

```

password-cracker/main.rs
27  #[autoparallelise]
28  fn hash(word: &String) -> String {
29      let mut hash_word = word.clone();
30      // Hash word using Sha256
31      for i in (0..1000).rev() {
32          let mut hasher = Sha256::new();
33          hasher.input_str(&hash_word);
34          hash_word = hasher.result_str();
35      }
36      hash_word
37  }
  
```

Listing 5: Hash function of the password cracker program

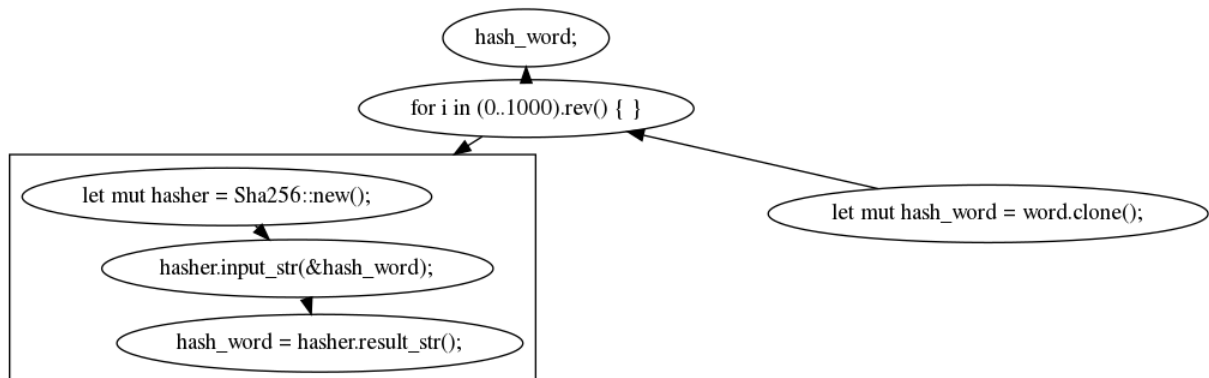


Figure 5.3: Dependency analysis graph of the hash function from the password cracker program

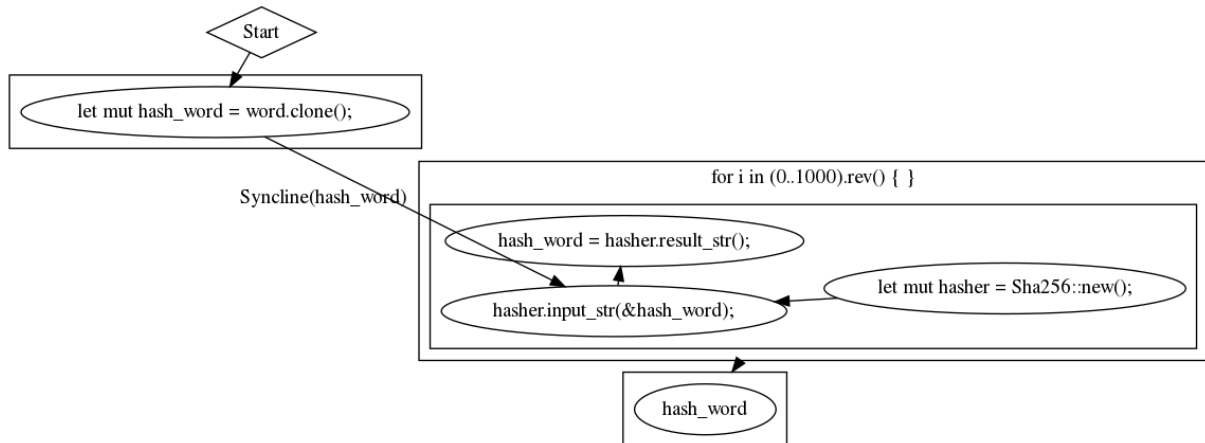


Figure 5.4: Schedule of the hash function from the password cracker program

```

password-cracker/main.rs
39 #[autoparallelise]
40 fn main() {
41     let now = Instant::now();
42
43     let dictionary: Vec<String> = load_dictionary();
44     let password_hash =
45     ⇨ format!("0954229bd82060f9d55ccc310b315ea831b9ba8faee1b76f66a19cc71140dfd7");
46
47     for id in 0..dictionary.len() {
48         let word = dictionary[id].clone();
49         let hash_word = hash(&word);
50         if hash_word == password_hash {
51             println!("Password is {}", word);
52         }
53     }
54
55     let elapsed = now.elapsed();
56     let sec = (elapsed.as_secs() as f64) + (elapsed.subsec_nanos() as f64 / 1000_000_000.0);
57     println!("Seconds: {}", sec);
58 }

```

Listing 6: Main method of the password cracker program

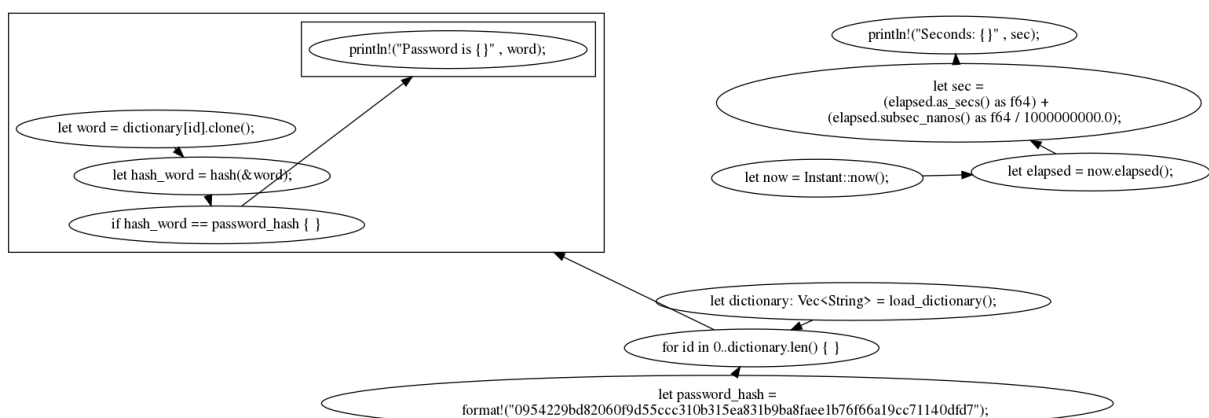


Figure 5.5: Dependency analysis graph of the main method from the password cracker program

5.2 Limitations/Improvements

My implementation to the problem does not work for all sequential programs.

Chapter 6

Discussion

6.1 Achievements

TODO: Deals with inter-weaved programs pretty well, as long as they use features implemented

6.2 Improvements

TODO: Deal with unsafe, non Send traits, mutable/immutable, etc.

TODO: Immutable variables should be cloned, instead of moved/synchronised across threads. Increases memory usage, but would increase program speed

TODO: Manage number of threads being used

TODO: Combine threads together when it is not worth it to parallelise.

6.3 Approach Changes

TODO: Expression level parallelisms

TODO: Work with HIR instead so I don't have to deal with as much syntax sugar

Chapter 7

Conclusion

The problem that I wanted to tackle was automatically converting sequential source code into a parallelised program. This is a hard problem which is well established in the field of computer science. A literature review was included as part of this report, but due to the size of the field, it is no where near complete. It was very unlikely that I could further the research in this field with the amount of time I had.

The solution I chose was to use a compiler plugin which had the downside that I had to deal with all of the syntax sugar. This dramatically increased the number of cases I had to deal with, and not all of them were dealt with correctly.

TODO: Write conclusion

Chapter 8

Appendix

8.1 References

- Baskaran, Muthu, Jj Ramanujam and P Sadayappan (2010). “Automatic C-to-CUDA code generation for affine programs”. In: *Compiler Construction*, pp. 244–263.
- Beletska, 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.
- 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.
- 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.

8.2 Submission File Structure

TODO: Write section

8.3 Running the Code

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

```
curl https://sh.rustup.rs -sSf | sh -s -- -y --default-toolchain nightly-2018-03-05
```

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="*", git="https://github.com/MichaelOultram/Auto_Parallelise/"}
```

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 once to run the analysis stage:

```
cargo build --release
```

Normally you would compile the code again with the same command to apply the modifications but due to a bug in the rust nightly compiler (#46489), this doesn't work. Instead you must pipe stdout into a file:

```
cargo build --release > parallel_code.rs
```

Normally you would just run the parallelised code but due to the bug, you will need to create a new project and copy parallel_code.rs along with any imports. Then you can:

```
cargo run --release
```

8.4 Rust Compiler Types

```
rust/src/libsyntax/visit.rs
34 pub enum FnKind<'a> {
35     /// fn foo() or extern "Abi" fn foo()
36     ItemFn(Ident, Unsafety, Spanned<Constness>, Abi, &'a Visibility, &'a Block),
37
38     /// fn foo(&self)
39     Method(Ident, &'a MethodSig, Option<&'a Visibility>, &'a Block),
```

```

40
41     /// |x, y| body
42     Closure(&'a Expr),
43 }

```

Listing 7: FnKind enum

```

rust/src/libsyntax/ast.rs
489 pub struct Block {
490     /// Statements in a block
491     pub stmts: Vec<Stmt>,
492     pub id: NodeId,
493     /// Distinguishes between `unsafe { ... }` and `{ ... }`
494     pub rules: BlockCheckMode,
495     pub span: Span,
496     pub recovered: bool,
497 }

```

Listing 8: Block struct

```

rust/src/libsyntax/ast.rs
781 pub struct Stmt {
782     pub id: NodeId,
783     pub node: StmtKind,
784     pub span: Span,
785 }

```

Listing 9: Stmt struct

```

rust/src/libsyntax/ast.rs
815 pub enum StmtKind {
816     /// A local (let) binding.
817     Local(P<Local>),
818
819     /// An item definition.
820     Item(P<Item>),
821
822     /// Expr without trailing semi-colon.
823     Expr(P<Expr>),
824     /// Expr with a trailing semi-colon.
825     Semi(P<Expr>),
826     /// Macro.
827     Mac(P<(Mac, MacStmtStyle, ThinVec<Attribute>)>),
828 }

```

Listing 10: StmtKind enum

```

rust/src/libsyntax/ast.rs
899 pub struct Expr {
900     pub id: NodeId,
901     pub node: ExprKind,
902     pub span: Span,
903     pub attrs: ThinVec<Attribute>
904 }

```

Listing 11: Expr struct

```

rust/src/libsyntax/ast.rs
987 pub enum ExprKind {
988     /// A `box x` expression.
989     Box(P<Expr>),
990     /// First expr is the place; second expr is the value.
991     InPlace(P<Expr>, P<Expr>),
992     /// An array `[a, b, c, d]`
993     Array(Vec<P<Expr>>),
994     /// A function call
995     ///
996     /// The first field resolves to the function itself,
997     /// and the second field is the list of arguments.
998     /// This also represents calling the constructor of

```

```

999      /// tuple-like ADTs such as tuple structs and enum variants.
1000      Call(P<Expr>, Vec<P<Expr>>),
1001      /// A method call (`x.foo::<'static, Bar, Baz>(a, b, c, d)`)
1002      ///
1003      /// The `PathSegment` represents the method name and its generic arguments
1004      /// (within the angle brackets).
1005      /// The first element of the vector of `Expr`s is the expression that evaluates
1006      /// to the object on which the method is being called on (the receiver),
1007      /// and the remaining elements are the rest of the arguments.
1008      /// Thus, `x.foo::<Bar, Baz>(a, b, c, d)` is represented as
1009      /// `ExprKind::MethodCall(PathSegment { foo, [Bar, Baz] }, [x, a, b, c, d])`.
1010      MethodCall(PathSegment, Vec<P<Expr>>),
1011      /// A tuple `(a, b, c, d)`
1012      Tup(Vec<P<Expr>>),
1013      /// A binary operation (For example: `a + b`, `a * b`)
1014      Binary(BinOp, P<Expr>, P<Expr>),
1015      /// A unary operation (For example: `!x`, `*x`)
1016      Unary(UnOp, P<Expr>),
1017      /// A literal (For example: `1`, `"foo"`)
1018      Lit(P<Lit>),
1019      /// A cast (`foo as f64`)
1020      Cast(P<Expr>, P<Ty>),
1021      Type(P<Expr>, P<Ty>),
1022      /// An `if` block, with an optional else block
1023      ///
1024      /// `if expr { block } else { expr }`
1025      If(P<Expr>, P<Block>, Option<P<Expr>>),
1026      /// An `if let` expression with an optional else block
1027      ///
1028      /// `if let pat = expr { block } else { expr }`
1029      ///
1030      /// This is desugared to a `match` expression.
1031      IfLet(P<Pat>, P<Expr>, P<Block>, Option<P<Expr>>),
1032      /// A while loop, with an optional label
1033      ///
1034      /// `label: while expr { block }`
1035      While(P<Expr>, P<Block>, Option<SpannedIdent>),
1036      /// A while-let loop, with an optional label
1037      ///
1038      /// `label: while let pat = expr { block }`
1039      ///
1040      /// This is desugared to a combination of `loop` and `match` expressions.
1041      WhileLet(P<Pat>, P<Expr>, P<Block>, Option<SpannedIdent>),
1042      /// A for loop, with an optional label
1043      ///
1044      /// `label: for pat in expr { block }`
1045      ///
1046      /// This is desugared to a combination of `loop` and `match` expressions.
1047      ForLoop(P<Pat>, P<Expr>, P<Block>, Option<SpannedIdent>),
1048      /// Conditionless loop (can be exited with break, continue, or return)
1049      ///
1050      /// `label: loop { block }`
1051      Loop(P<Block>, Option<SpannedIdent>),
1052      /// A `match` block.
1053      Match(P<Expr>, Vec<Arm>),
1054      /// A closure (for example, `move |a, b, c| a + b + c`)
1055      ///
1056      /// The final span is the span of the argument block `|...|`
1057      Closure(CaptureBy, P<FnDecl>, P<Expr>, Span),
1058      /// A block `{ ... }`
1059      Block(P<Block>),
1060      /// A catch block `catch { ... }`
1061      Catch(P<Block>),
1062
1063      /// An assignment (`a = foo()`)
1064      Assign(P<Expr>, P<Expr>),
1065      /// An assignment with an operator
1066      ///
1067      /// For example, `a += 1`.
1068      AssignOp(BinOp, P<Expr>, P<Expr>),
1069      /// Access of a named struct field (`obj.foo`)
1070      Field(P<Expr>, SpannedIdent),
1071      /// Access of an unnamed field of a struct or tuple-struct

```

```

1072    ///
1073    /// For example, `foo.0`.
1074    TupField(P<Expr>, Spanned<usize>),
1075    /// An indexing operation (`foo[2]`)
1076    Index(P<Expr>, P<Expr>),
1077    /// A range (`1..2`, `1..`, `..2`, `1...2`, `1...`, `...2`)
1078    Range(Option<P<Expr>>, Option<P<Expr>>, RangeLimits),
1079
1080    /// Variable reference, possibly containing `::` and/or type
1081    /// parameters, e.g. `foo::bar::baz`.
1082    ///
1083    /// Optionally "qualified",
1084    /// E.g. `::SomeType`.
1085    Path(Option<QSelf>, Path),
1086
1087    /// A referencing operation (`&a` or `&mut a`)
1088    AddrOf(Mutability, P<Expr>),
1089    /// A `break`, with an optional label to break, and an optional expression
1090    Break(Option<SpannedIdent>, Option<P<Expr>>),
1091    /// A `continue`, with an optional label
1092    Continue(Option<SpannedIdent>),
1093    /// A `return`, with an optional value to be returned
1094    Ret(Option<P<Expr>>),
1095
1096    /// Output of the `asm!()` macro
1097    InlineAsm(P<InlineAsm>),
1098
1099    /// A macro invocation; pre-expansion
1100    Mac(Mac),
1101
1102    /// A struct literal expression.
1103    ///
1104    /// For example, `Foo {x: 1, y: 2}`, or
1105    /// `Foo {x: 1, .. base}`, where `base` is the `Option<Expr>`.
1106    Struct(Path, Vec<Field>, Option<P<Expr>>),
1107
1108    /// An array literal constructed from one repeated element.
1109    ///
1110    /// For example, `[1; 5]`. The first expression is the element
1111    /// to be repeated; the second is the number of times to repeat it.
1112    Repeat(P<Expr>, P<Expr>),
1113
1114    /// No-op: used solely so we can pretty-print faithfully
1115    Paren(P<Expr>),
1116
1117    /// `expr?`
1118    Try(P<Expr>),
1119
1120    /// A `yield`, with an optional value to be yielded
1121    Yield(Option<P<Expr>>),
1122 }

```

Listing 12: ExprKind enum

```

rust/src/libsyntax_pos/lib.rs
143 pub struct SpanData {
144     pub lo: BytePos,
145     pub hi: BytePos,
146     /// Information about where the macro came from, if this piece of
147     /// code was created by a macro expansion.
148     pub ctxt: SyntaxContext,
149 }

```

Listing 13: SpanData struct