

Parslers: A staged selective parser combinator library for the Rust programming language

How to make fast and easy-to-use parser combinator libraries in Rust
faster and easier to use

Jordan Hall

Individual Project
Bachelors of Engineering – Computing
September 1, 2023

| | |
|-------------------------|---|
| INDIVIDUAL PROJECT | |
| Imperial College London | |
| | |
| Degree Programme: | Bachelors of Engineering – Computing |
| | |
| Identification number: | |
| Author: | Jordan Hall |
| Title: | Parslers: A staged selective parser combinator library for the Rust programming language How to make fast and easy-to-use parser combinator libraries in Rust faster and easier to use |
| Supervisor: | Jamie Willis |
| | |
| Abstract: | |
| | |
| Keywords: | Parser Combinator, Rust, Parslers, Parsing |
| Number of pages: | 58 |
| Language: | English |

CONTENTS

| | | |
|----------|--|-----------|
| 1 | Introduction | 5 |
| 2 | Background | 6 |
| 2.1 | Parsing in Computer Science | 6 |
| 2.2 | Formal Grammar | 6 |
| 2.2.1 | <i>Regular Grammar</i> | 7 |
| 2.2.2 | <i>Context-Free Grammar</i> | 7 |
| 2.3 | Tools for Writing Parsers | 8 |
| 2.3.1 | <i>Parser Generators</i> | 8 |
| 2.3.2 | <i>Handwritten Parsers</i> | 8 |
| 2.3.3 | <i>Parser Combinators</i> | 9 |
| 2.4 | The Rust Programming Language | 9 |
| 2.4.1 | <i>History</i> | 9 |
| 2.4.2 | <i>Rise to Popularity</i> | 9 |
| 2.4.3 | <i>Meta-programming using Procedural Macros</i> | 10 |
| 3 | Ethical Issues | 13 |
| 4 | Parslers | 14 |
| 4.1 | Requirements | 14 |
| 4.1.1 | <i>Expressive Power</i> | 15 |
| 4.1.2 | <i>Performance</i> | 15 |
| 4.1.3 | <i>Correctness</i> | 15 |
| 4.2 | Library Structure | 15 |
| 4.3 | Parser Definition Interface | 16 |
| 4.3.1 | <i>Fixing the structural problem</i> | 16 |
| 4.3.2 | <i>The Primitive Combinators</i> | 17 |
| 4.3.3 | <i>Auxiliary Combinators</i> | 18 |
| 4.3.4 | <i>The Parsler Trait</i> | 19 |
| 4.4 | The Reflect Trait | 22 |
| 4.4.1 | <i>Using Dynamic Dispatch</i> | 25 |
| 4.4.2 | <i>Returning to Static Dispatch</i> | 26 |
| 4.5 | Lazy Parsers | 26 |
| 4.6 | Named Parsers | 29 |
| 4.7 | Looping Parsers | 30 |
| 4.8 | Alternative Parser | 32 |
| 4.9 | Untyped AST | 32 |
| 4.10 | Backend Optimisations | 33 |
| 4.10.1 | <i>Lawful Reduction</i> | 34 |
| 4.10.2 | <i>Output Usage Analysis</i> | 34 |
| 4.10.3 | <i>Validation Parsers</i> | 35 |
| 4.10.4 | <i>Analysis of Recursive Parsers</i> | 35 |
| 4.11 | Code Generation | 36 |
| 4.11.1 | <i>Registering Functions</i> | 36 |
| 4.11.2 | <i>Parsing and Generating Rust Code with Syn and Quote</i> | 38 |
| 4.11.3 | <i>The type of Parsers</i> | 40 |
| 4.11.4 | <i>Recursively Compiling Parsers</i> | 41 |
| 5 | Evaluation | 43 |

| | | |
|----------|--|-----------|
| 5.1 | Performance | 43 |
| 5.1.1 | <i>Prior Discussion</i> | 43 |
| 5.1.2 | <i>Comparative JSON Benchmarks</i> | 43 |
| 5.1.3 | <i>Comparative Branflakes Benchmarks</i> | 44 |
| 5.1.4 | <i>Staging and Optimisation Benchmarks</i> | 45 |
| 5.1.5 | <i>Conclusion</i> | 46 |
| 5.2 | Expressive Power | 46 |
| 5.2.1 | <i>Sequence</i> | 47 |
| 5.2.2 | <i>Ordered Choice</i> | 47 |
| 5.2.3 | <i>Zero/One-or-More</i> | 48 |
| 5.2.4 | <i>Optional</i> | 48 |
| 5.2.5 | <i>And/Not Predicates</i> | 48 |
| 5.2.6 | <i>Conclusion</i> | 48 |
| 5.3 | Correctness | 48 |
| 5.4 | Usability | 49 |
| 5.4.1 | <i>Staging</i> | 49 |
| 5.4.2 | <i>The Reflect Trait</i> | 50 |
| 5.4.3 | <i>Conclusion</i> | 51 |
| 6 | Conclusions | 52 |
| 6.1 | Future Work | 52 |
| 6.1.1 | <i>The Reflect Trait</i> | 52 |
| 6.1.2 | <i>Performance Optimisations</i> | 52 |
| 6.1.3 | <i>Error Messages</i> | 54 |
| 6.1.4 | <i>Asynchronous Parsing</i> | 54 |
| 6.2 | Example of the Branflakes Parser | 55 |
| | References | 58 |

FOREWORD

ACKNOWLEDGMENTS

I'd like to thank all those who helped support me in this project, including (but not limited to):

Jamie Willis, for inspiring my curiosity in the field of parsing and especially parser combinators.

1 INTRODUCTION

In this project, I aim to implement a library for the Rust programming language. The aim of this library is to provide tools to help non-expert parser architects develop fast parsers with an easy-to-use interface. I will leverage the Rust language's powerful meta-programming tool – Procedural Macros – to implement a domain-specific language that is both easy to use and optimizes the built parser into fast, native Rust code at compile time.

The motivation is derived from the fact that generating parsers of structured data (e.g programming languages) often comes with a lot of complexity. Programmers have to concern themselves with speed, correctness and maintainability of their parsers, when parsing is oftentimes not the focal point of the language.

Although many parser generators and libraries exist, few are able to both optimize the resulting parser and provide a user-friendly interface for the programmer. Even fewer parser libraries and generators have structures similar to the original grammar. However, functional parsers (including parser combinators) oftentimes look similar to the grammar that they parse. Recent research into parser combinators show that if the structure of the parser is known at compile time, then they can be optimized by languages that offer staged compilation (such as Haskell, Scala 3, and Rust). This makes parser combinators a promising methodology for language designers.

| Languages | Automaton |
|------------------------|---|
| Regular | Finite state automaton |
| Context-free | Non-deterministic pushdown automaton |
| Context-sensitive | Linear-bounded non-deterministic Turing machine |
| Recursively enumerable | Turing machine |

2 BACKGROUND

The subject of parsers and parser libraries both have a rich history in the field of computer science. This section will introduce the background of parsing as a field of computer science, make the case for why further development into parsing speed is necessary, and lay out the groundwork for Parslers as the solution to writing fast and succinct parsers in the Rust programming language.

2.1 Parsing in Computer Science

In the field of computer science, parsing is a process used to analyse a string of symbols, usually text, in accordance with a set of rules known as a grammar. Early work looked to describe and understand natural languages (such as English or Spanish) in a way that made it amenable to analyse using mathematical principles. The resulting field of study, formal language theory (Jäger & Rogers 2012), laid the groundwork that made writing computer programs capable of analysing and generating said natural languages possible.

In modern times, parsing has become a core tool used by software engineers on a daily basis. Software engineers interact with compilers using a language format specification (the syntax of a language), which must be parsed into structured data for a compiler or interpreter to process. Different programs interact with one another using a serialisation/deserialisation format (such as JSON, XML) in order to communicate. Images are encoded into a shared structured format and must be parsed by image-viewing software in order to be displayed.

As parsing data is so ubiquitous in computer science, it follows that much research has been done on formalising the various methods computer scientists and programmers use to translate structured data into text, and vice versa.

2.2 Formal Grammar

For mathematicians to build upon each other's work, they must first agree on the method with which they communicate. Although linguists and etymologists have discussed the nature of languages for millennia prior, mathematicians, such as Noam Chomsky, defined the concrete notion of formal grammar as a set of rules and valid productions.

An example of such a grammar would be:

$$\begin{array}{l} \langle S \rangle ::= \langle S \rangle 'a' \\ \quad \quad | \quad 'b' \end{array}$$

Strings of symbols such as *ba*, *b*, *bbbaaa*, and *baba* are all valid strings within this grammar. However, *ab* would not be considered a valid string within this grammar, as there's no way to reduce this string down to *S*.

Chomsky formulated a hierarchy of formal grammar, which describe the power of all different grammar types and the type of automaton necessary to parse such a grammar (Allott, Lohndal & Rey 2021).

2.2.1 Regular Grammar

Regular grammars are the most basic form of formal grammar. They are grammars in which each production must always start with a terminal:

$$\begin{aligned}\langle S \rangle &::= 'a' \langle S \rangle \\ &| 'b' \langle S \rangle\end{aligned}$$

This is to say, a production such as $S \rightarrow aSa$ would not be a valid regular grammar. The most common use case of regular grammars is the regular expression language `Regex`. `Regex` provides fast text validation with simple expressions.

An example of a `Regex` expression is one used to validate dates in the form DD/MM/YYYY:

```
((([0-9]|[12][0-9]|3[01])([/])([0-9]{1,2}|10|12)([/])(\d{4}))|((([0-9]|[12][0-9]|30)([/]))  
→ ([0-9]{1,2}|11)([/])(\d{4}))|((([0-9]|1[0-9]|2[0-8])([/])(02)([/])(\d{4}))|((29)(\.|-|\/)(02)([/]))  
→ ([02468][048]00))|((29)([/])(02)([/])([13579][26]00))|((29)([/])(02)([/])([0-9][0-9][0][48]))  
→ |((29)([/])(02)([/])([0-9][0-9][2468][048]))|((29)([/])(02)([/])([0-9][0-9][13579][26]))))
```

2.2.2 Context-Free Grammar

Context-free grammar are written with production rules in the form:

$$\langle S \rangle ::= w$$

Where S is the non-terminal to be produced, and w is a string of terminals and non-terminals. What makes such a grammar context-free is that S can be produced anywhere that w is successfully matched. It does not depend on the surrounding context.

It is worth noting that most computer programming languages are described using context-free grammar because they are typically easier to generate efficient parsers for. Furthermore, context-free grammars allow programmers to quickly create mental models that can be applied when writing syntax for the language.

Although context-free grammars are the most widely-used category of grammar for everyday use cases, they are to some extent limited. A good example is parsing a number. A grammar for a floating-point number can be described using the following grammar:

$$\begin{aligned}\langle Float \rangle &::= '-' \langle UnsignedFloat \rangle \\ &| \langle UnsignedFloat \rangle \\ \langle UnsignedFloat \rangle &::= \langle UInt \rangle '.' \langle UInt \rangle \\ \langle UInt \rangle &::= \langle Digit \rangle \langle UInt \rangle \\ &| \langle Digit \rangle \\ \langle Digit \rangle &::= '0'..'9'\end{aligned}$$

Although this grammar correctly describes all real numbers, a context-free grammar cannot be used to easily describe the real numbers describable by a fixed-size floating point format. It is possible by enumerating all values that fit within this criteria, however this is not useful for generating efficient parsers. For this reason, parsing is generally performed before a semantic check of the parsed data. This way, grammar for programming languages may remain context-free for simplicity.

2.3 Tools for Writing Parsers

Almost every program that deals with text in some way likely needs to convert it into structured data to be useful. To fulfil this need many tools and libraries have been developed to ease the time, speed, and correctness of parsers.

2.3.1 Parser Generators

Parser generators offer an end-to-end solution for writing parsers, the most famous of which are ANTLR and Bison. They both target multiple languages, and work by taking in formatted grammar and output autogenerated code in the target language. Parser generators work well for languages that do not have strong meta-programming capabilities, and they oftentimes use grammar formats that look similar to the original grammar. ANTLR uses Extended Backus-Naur form, and Bison uses BNF. Parser generators that use produce LR parsers, like Bison, are able to optimise the underlying parser during generation.

However, parser generators typically produce concrete syntax trees after parsing. This makes it difficult to maintain the resulting parser because changes to the grammar will result in the language designer having to also change the code that converts the concrete syntax tree into the abstract syntax tree, suitable for semantic analysis.

2.3.2 Handwritten Parsers

Handwritten parsers are generally written from scratch in the target language. The most common of which are top-down parsers.

Languages with hand-written parsers include:

- Java (OpenJDK)
- Golang
- Swift
- TypeScript
- JavaScript
- C/C++ (both GCC and Clang)

There are several reasons a language designer would write their own hand-written parsers, but Golang (abbreviated *Go*) provides an interesting case-study. Go originally used a YACC-based parser, a parser generator similar to Bison and ANTLR mentioned above. In November 2015, Go switched to a hand-written parser (Griesemer & Fitzpatrick 2015), which yielded a 18% parsing speedup, and a 3% compilation speedup for the compiler itself. Because hand-written parsers give the language designer increased control in how lexing/parsing takes place, oftentimes resulting in faster parsers. Handwritten parsers can also generate the abstract syntax tree directly during parsing.

Handwritten parsers are also not limited to context-free grammars, as opposed to many parsing libraries and parser generators.

Unfortunately, handwritten parsers suffer from maintainability issues. Because lexing, parsing, error messaging, and AST generation must all be programmed by the language designer, the parser now becomes a significant portion of the codebase where it otherwise wouldn't be. Furthermore, writing

a correct and optimised handwritten parser requires expert knowledge on how to architect such parsers.

2.3.3 Parser Combinators

Parser combinators work by composing many higher-order functions together and then applying that function to an input stream to generate an output AST (Hutton & Meijer 1996). The parsing code can be integrated with the rest of the codebase, much like a hand-written parser. They can produce the target AST directly during parsing. Furthermore, parser combinators look very similar to the original grammar in many implementations, which make them maintainable.

Until recently, parser combinators were considered slow in comparison to handwritten parsers and parser generators. However, more recent research has shown that parser combinators can be optimized in a separate step from the rest of the program, at compile time, by removing the monad (Willis, Wu & Pickering 2020). This allows parser combinators to look like the grammar when writing them, but are optimized and then compiled down to native code in the target language.

2.4 The Rust Programming Language

2.4.1 History

The Rust programming language was developed by Graydon Hoare at Mozilla in 2006. Mozilla then began officially sponsoring the development of the language in 2009. The first stable version of the Rust programming language was released in 2015. In 2021, the Rust Foundation was established, taking over development of the Rust programming from Mozilla.

The initial motivation came from a common problem in modern software engineering: concurrency. As modern computer processing unit improvements focus less and less on raw speed and more towards very high core counts, concurrent programming is seen as the way to improve the performance of computer programs. However, this comes at a cost in the mental overhead involved in programming. The difficulty stems from the fact that modern CPUs use the Total-Store-Ordering concurrent memory model. This means that changes to memory locations in RAM made by one core may be stored in a core-specific write-back buffer, and not reflected in RAM (and thereby other cores) for an unpredictable period of time. This is known as a data race. When writing concurrent software in traditional programming languages, a great deal of effort must be made to ensure that potentially concurrent memory access is atomic, and other cores are notified in a predictable time frame.

A motivating example of this was the previous Mozilla Firefox CSS engine. CSS styling is a top-down process, wherein the children of the parent object can be styled independently top-down (and thus in parallel). There have been attempts to parallelize the previous C++ engine to leverage this, however none were successful as the codebase was too large, and the mental workload for maintaining the concurrent invariants turned out to be overwhelming.

To solve this problem, Rust utilizes a “borrowing” memory model, whereby objects can be “owned” by a scope, and then borrowed immutably many times, or borrowed mutably only once, but never both at the same time. This solves the problem of data races by ensuring that no two cores can read and write to a memory location at the same time. Rust then builds an abstraction over these rules by providing mutex and atomic variable implementations, as well as smart pointers that automatically destruct objects when they are no longer in use.

2.4.2 Rise to Popularity

Rust has become incredibly popular since release, being top of Stack Overflow Developer Survey’s “Most Loved Programming Language” every year from 2016-2022 (the time of writing). Many at-

tribute this popularity to its ease of use and speed.

2.4.3 Meta-programming using Procedural Macros

The Rust programming language also provides a powerful tool for assisting development. The Procedural Macro is a Rust program that is run during compile time, and transforms a sequence of parsed Rust tokens to another sequence of Rust tokens, which are then compiled alongside the rest of the Rust program. The Procedural Macro is a core part of the language as it provides a tool for library developers to make using said library easier to develop with, saving time.

Take, for example, Rocket. Rocket is a framework for developing web servers in Rust, and provides many macros to improve the ease of development.

An example of a simple web-server that returns “Hello World” when called.

```
1  #[macro_use]
2  extern crate rocket;
3
4  #[get("/")]
5  fn hello() -> &'static str {
6      "Hello, world!"
7  }
8
9  #[launch]
10 fn rocket() -> _ {
11     rocket::build()
12         .mount("/", routes![hello])
13 }
```

The `#[get("/")]` and `#[launch]` tags are macros that tell the compiler to parse the resulting functions into a stream of tokens, and pass them into a Rust function to transform them. We can see how the rocket procedural macros expand the little code above into the code below:

```
1  #[macro_use]
2  extern crate rocket;
3  fn hello() -> &'static str {
4      "Hello, world!"
5  }
6
7  #[doc(hidden)]
8  #[allow(non_camel_case_types)]
9  /// Rocket code generated proxy structure.
10 struct hello {}
11 /// Rocket code generated proxy static conversion implementations.
12 impl hello {
13     #[allow(non_snake_case, unreachable_patterns, unreachable_code)]
14     fn into_info(self) -> ::rocket::route::StaticInfo {
15         fn monomorphized_function<'__r>(<
16             __req: &'__r ::rocket::request::Request<'_,
17             __data: ::rocket::data::Data<'__r,>
```

```

18 ) -> ::rocket::route::BoxFuture<'__r> {
19     ::std::boxed::Box::pin(async move {
20         let ___responder = hello();
21         ::rocket::route::Outcome::from(__req, ___responder)
22     })
23 }
24 ::rocket::route::StaticInfo {
25     name: "hello",
26     method: ::rocket::http::Method::Get,
27     uri: "/",
28     handler: monomorphized_function,
29     format: ::std::option::Option::None,
30     rank: ::std::option::Option::None,
31     sentinels: <[_]>::into_vec(
32         #[rustc_box]
33         ::alloc::boxed::Box::new([
34             {
35                 #[allow(unused_imports)]
36                 use ::rocket::sentinel::resolution::{DefaultSentinel as _, Resolve};
37                 ::rocket::sentinel::Sentry {
38                     type_id: std::any::TypeId::of::<&'_ str>(),
39                     type_name: std::any::type_name::<&'_ str>(),
40                     parent: None,
41                     location: ("src/main.rs", 5u32, 15u32),
42                     specialized: Resolve::<&'_ str>::SPECIALIZED,
43                     abort: Resolve::<&'_ str>::abort,
44                 }
45             },
46             {
47                 #[allow(unused_imports)]
48                 use ::rocket::sentinel::resolution::{DefaultSentinel as _, Resolve};
49                 ::rocket::sentinel::Sentry {
50                     type_id: std::any::TypeId::of::<str>(),
51                     type_name: std::any::type_name::<str>(),
52                     parent: None.or(Some(std::any::TypeId::of::<&'_ str>())),
53                     location: ("src/main.rs", 5u32, 24u32),
54                     specialized: Resolve::<str>::SPECIALIZED,
55                     abort: Resolve::<str>::abort,
56                 }
57             },
58         ]),
59     ),
60 }
61 }
62 #[doc(hidden)]
63 pub fn into_route(self) -> ::rocket::Route {
64     self.into_info().into()
65 }
66 }
67

```

As you can see, macros do a lot of work behind the scenes for developers, and allow Rust to have libraries that are incredibly easy to use.

3 ETHICAL ISSUES

This project is aimed to help people make efficient, correct, and maintainable parsers. Because of this scope, there are very few ethical considerations that must be taken into account. That being said, it is important in the conduct of this project to point out that the resulting library may be used in production, or on projects where correctness is of great importance. This puts the responsibility of correctness onto me, the creator of the library. I must also be honest and clear in all places where the library may be found that the aim of the project was not to create safety-critical parsers, and there may be issues that I have not yet identified.

4 PARSLERS

This chapter introduces Parslers, the library developed for this thesis. Parslers is a Rust parser-combinator framework developed to empower non-expert parser developers to generate blazingly fast parsers while only using high-level abstractions. This is achieved by using Rust's staged compilation to take high-level parsers and compile them into native Rust code, all at the compile step.

Section Requirements lays out the requirements that Parslers would have to meet to become a viable parser combinator library competitive with the rest of the monadic parser combinator libraries already available in Rust.

Section Library Structure shows how the user is expected to interact with Parslers. This section also details how Parslers makes use of Rust's staged compilation pipeline to perform optimisations on user-generated parsers and generate native Rust code.

Section Parser Definition Interface follows by showing how Parslers is different from other monadic parser combinator libraries by removing the monad entirely, and how Parslers uses the selective combinators to regain some of the expressive power lost with this decision. It also introduces the primitive and auxiliary combinators. Finally, the Parsler trait is introduced, one of the core traits that the library users interact with in order to build their parsers.

Section The Reflect Trait starts by introducing the key challenge in the staged-compilation pipeline - how to ensure that runtime functions and values survive the build phase into the generated native code. It also discusses the challenges involved in the chosen solution, and how Parslers overcomes these challenges to achieve runtime reflection of functions - a novel approach to runtime reflection.

With the internals of the library laid out, sections Lazy Parsers, Named Parsers, Looping Parsers, and Alternative Parser discuss the parser combinators that allow for complex structure and improved efficiency. In particular, Lazy Parsers discusses the `LazyParser`, which allows Parslers to overcome the problem of infinite recursion when parsers are defined structurally instead of the commonly used higher-order function approach. Named Parsers introduces the `NamedParser`, which is used to prevent parsers from being duplicated during code generation, which is essential for recursive parsers. Looping Parsers tackles the problem of efficiency arising from the fact that Rust does not yet support tail-call optimisation. This section then goes on to solve this problem by introducing a novel parser, the `LoopParser`, which allows transforming recursive parsers into iterative parsers.

Now that the user-facing interface has been properly introduced, section Backend Optimisations discusses how fixing the structure of the parser by removing the monad can be utilised to perform static analysis on the users' parsers at compile-time. This section will lay out the current optimisations implemented by Parslers, including lawful reduction, and usage analysis optimisations. It will then finish by discussing the possible extensions to static analysis that can be implemented in the future in order to further improve Parslers' performance characteristics.

Finally, section Code Generation details how these high-level parsers are transformed into native Rust code. It also discusses how functions and values used in the high-level parser definition survive into the generated code.

4.1 Requirements

In order for Parslers to become a viable parser combinator library suitable to replace the current existing parsing libraries available to the Rust ecosystem, a number of requirements must be set and

adhered to. This section discusses those requirements, which will be the guiding decisions used in the design-process of the library, and each will be evaluated in turn at the end of this report.

4.1.1 Expressive Power

Expressive power means the range of grammars that can be parsed by Parslers. Detailing this is key, as it will help define which parsing problems Parslers can be applied to. Because Parslers is a recursive-descent parser combinator library, it makes sense to require that Parslers can parse any PEG grammar. The evaluation section will show that Parslers can parse a wider range of grammars than PEG.

4.1.2 Performance

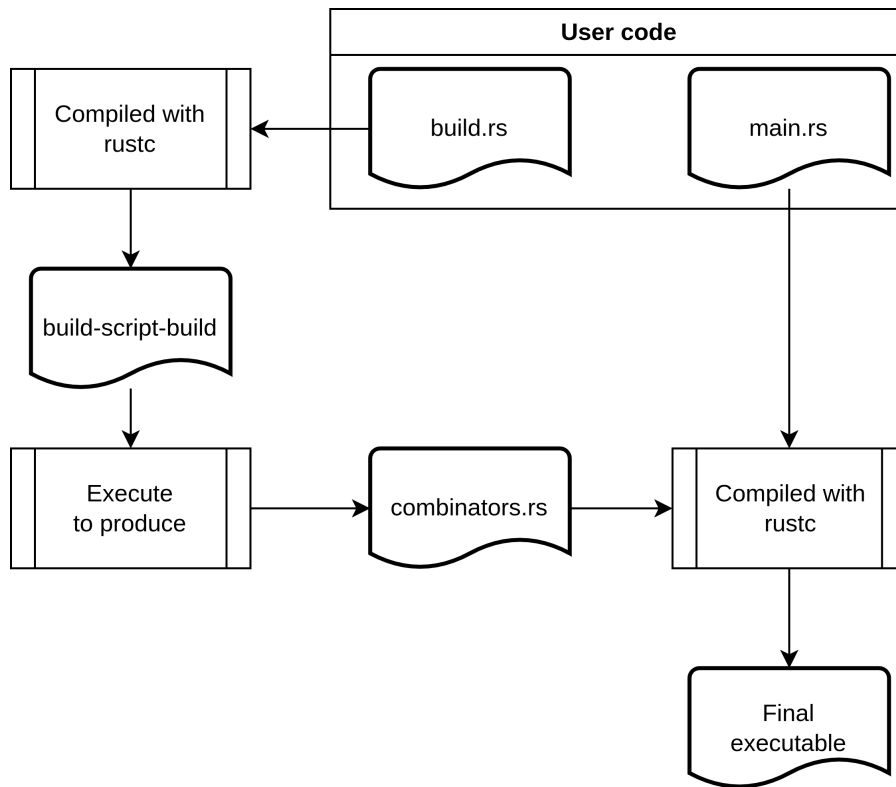
Parslers' main concern as a parsing library is using static analysis to optimise user-written parsers at compile time. It, therefore, follows that Parslers should have performance characteristics better than other parsing libraries available in the Rust ecosystem. However, it is acceptable for Parslers to be out-performed by hand-written parsers for specific grammars, as these parsers are embedded with the full possible domain-specific knowledge available. Parslers will, however, attempt to achieve performance characteristics similar to hand-written parsers where possible.

4.1.3 Correctness

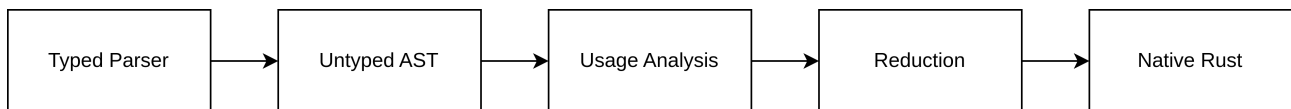
Although parsers are often tested on a large corpus of correct and incorrect example text, correctness is still one of the top priorities of Parslers. As per this requirement, Parslers is required to be type-safe in the user-facing aspects of the library. Furthermore, the restrictions on the types of parsers should prevent ill-formed parsers from being created at compile time. There should be no runtime checks on the correctness of the parsers, purely relying on the strong type-checking Rust provides at compile time.

4.2 Library Structure

Parslers takes a staged-compilation approach to parser generation. The user writes their parsers as typed Rust in the `build.rs` file that is compiled and run before the main binary is compiled. Parslers leverages this stage of compilation to perform optimisations on the user-written parsers before translating the parsers into native rust functions in a separate file. The user then includes this file as an import to their codebase, where it is compiled with the rest of the binary.



Because Rust will type-check the user’s parser for us, we are then able to convert to an untyped AST to perform optimisations on it. Because the optimisations performed on the AST do not change the type of the produced values, they may be compiled into Rust code without needing to embed Rust’s type checker into the Parslers library. The structure of Parslers is as follows:



Usage analysis and structural reduction are the two key optimisation techniques which make Parslers’ generated parsers blazingly fast. We will detail these optimisations later in the report.

4.3 Parser Definition Interface

Now that Parslers’ staged compilation has been explained, it is time to move on to how the library presents the type of the parsers users construct. This section starts off by introducing the monadic parser, and how representing the parser as a monad destroys opportunities for static analysis to be performed on said parsers. It then goes on to discuss how the monad plays a large role in the expressive power of monadic parsers, before introducing a combinator that helps regain some of this expressive power - the selective combinator. After this, the primitive combinators are introduced, which form the building blocks used to build parsers and auxiliary combinators. This section then finishes off by introducing the `Parslers` trait, which all parsers and combinators must implement.

4.3.1 Fixing the structural problem

The most commonly-used parser combinator libraries in Rust cannot currently perform static analysis on parsers created with them. This section will discuss why the root of this problem lies in the use of a parser as a monad, and how removing the monad can allow for static analysis in Parslers. Then, this section will discuss how removing the monad leaves us without the full expressive power it provides, and how introducing the selective parser combinator helps to retrieve

some of this power.

What separates Parslers from the monadic parser combinator libraries is that it lacks the ability to generate arbitrary parsing structure at runtime. That is to say, there is no valid combinator which, as the result of previous input, generates a new parser to be executed on the remaining input. Because Parslers requires that the structure of parsers is fixed, this allows for static analysis to be performed at compile time.

In essence, a combinator in parser combinators is usually expressed as a monad. Monadic parsers typically have two key functions:

- `pure :: a -> Parser a`
- `flatMap :: (Parser a) -> (a -> Parser b) -> (Parser b)`.

`pure` simply lifts the value `a` into a parser which produces the value `a` when applied to an input string. `pure` does not pose a problem to static analysis, because its structure is known at compile-time. The problem lies with the `flatMap` method (also known as `bind`). `a -> Parser b` can generate an arbitrary parser based on the value of `a`. If the structure of the parser is to be known at compile-time, then `flatMap` will not be available to us.

The fixed structure does come at a cost of expressive power. If the user wishes to write an integer parser, they are constrained to the types Rust provides (`u8`, `u16`, `u32`, `u64`, `usize`, `i8`, `i16`, `i32`, `i64`, `isize`). If the user would like to constrain the parsed values to fit inside of a `u8`, then the parser would have to change behaviour based on the parsed value. This would typically be done using the `flatMap` operator, yielding either the parsed value as a `u8`, or an error if the value is too large.

Because `flatMap` is not available to Parslers, a selective combinator must be introduced. Previous work on selective parser combinators (Willis, Wu & Pickering 2020) produced the `branch(b, l, r)` combinator. In essence, it is a parser with the type:

```
branch :: Parser (Either a b) -> Parser (a -> c) -> Parser (b -> c) -> Parser c
```

`branch` executes parser `b`, and based on the result of that, will execute parser `l` or `r`. With this combinator, it would not be too difficult to generate a parser which reads a string of digits, attempts to parse it into a `u8` (returning `Either<u8, ()>`), and then either return it in parser `l` or throw an error in parser `r`.

By using `branch`, some of the expressive power of the monad can be restored to the fixed structure of Parslers. Further work has been done on selective parser combinators that implements what's known as a register (Willis, Wu & Schrijvers 2022), a mutable value that can be threaded throughout the parser. Using a register with the selective `branch` combinator restores the most of the expressiveness of the monad while maintaining the fixed structure of the parser.

4.3.2 The Primitive Combinators

Parser combinator libraries are designed with the approach that complex parsers can be represented by composing higher-order functions together. The leaf nodes of these higher-order functions are primitive combinators, which perform a very basic parsing task such as reading a character. Function-based parser combinator libraries (such as `nom`) allow the users to create these primitive combinators manually, and provide frequently used primitive combinators. However, because

Parslers represents parser combinators structurally, the set of primitive combinators is fixed and must encapsulate everything needed for the user to write their parsers.

The list of available primitive combinators in Parslers is:

| | | |
|--------------|--------------|----------------------|
| Pure(a) | Ap(p, q) | Branch(b, l, r) |
| Satisfy(f) | Map(p, f) | LazyParser(f) |
| Attempt(p) | Recognise(p) | NamedParser(p, name) |
| Or(p, q) | Look(p) | Loop(p0, pN) |
| Then(p, q) | NegLook(p) | |
| Before(p, q) | Empty | |

Pure(a) simply returns the value a, consuming no input.

Satisfy(f) takes a function f in the form `FnOnce(char) -> bool`. It will apply f to a consumed character, and return that character upon success. Will fail if f returns false and backtracks.

Attempt(p) will perform parser p, backtracking if p fails.

Or(p, q) will attempt to perform parser p, and then parser q if p fails. Does not backtrack.

Then(p, q) will perform parser p, then parser q, returning only the result of parser q.

Before(p, q) will perform parser p, then parser q, returning only the result of parser p.

Ap(p, q) takes parsers p and q where p returns a function to be applied to q. It will perform both in sequence, and then apply the results.

Map(p, f) is equivalent to Ap(Pure(f), p), but added for brevity.

Recognise(p) performs p, and then returns the string consumed by p.

Look(p) performs the parser p, and then backtracks on the consumed characters.

NegLook(p) performs the parser p, and fails if p succeeds.

Empty always fails.

Branch(b, l, r) is described in the section Fixing the structural problem

LazyParser(f) is described in the section Lazy Parsers

NamedParser(p, name) is described in the section Named Parsers

Loop(p0, pN) is described in the section Looping Parsers

4.3.3 Auxiliary Combinators

The primitive combinators introduced are enough to form any higher-order combinator that the user could use within the context of Parslers. However, many non-primitive combinators are frequently used in the construction of parsers. For this reason, Parslers comes pre-packaged with many frequently-used higher-order combinators to reduce the workload of the user.

The list of pre-packaged auxiliary combinators in Parslers are:

| | | |
|-------------|-------------------|------------|
| some_rev(p) | match_char(c) | opt(p) |
| many_ref(p) | one_of(cs) | ws(p) |
| some(p) | filtered_by(p, f) | dynamic(p) |
| many(p) | ident(f) | |
| some_map(p) | tag(s) | |
| many_map(p) | not(c) | |

some_rev(p) performs p at least once and collects the results into a `Vec` in reverse order. This is the recursive version of some(p).

many_rev(p) performs p zero or more times and collects the results into a `Vec` in reverse order.

This is the recursive version of `many(p)`.

`some(p)` does the same as `some_rev(p)`, but the results are added to the `Vec` in parsed order. This uses the `LoopParser` instead of recursion.

`many(p)` does the same as `many_rev(p)`, but the results are added to the `Vec` in parsed order. Like `some(p)`, this uses the `LoopParser` instead of recursion.

`some_map(p)` does the same as `some(p)`, but the results are collected into a `HashMap` instead of a `Vec`.

`many_map(p)` does the same as `many(p)`, but the results are collected into a `HashMap` instead of a `Vec`.

`match_char(c)` is equivalent to `Satisfy(|ch| ch == c)`. It is added for brevity.

`one_of(cs)` takes an iterator that produces at least one `char`, and produces a parser which attempts each character in sequence.

`filtered_by(p, f)` takes a parser with output type `O`, and function of the type `FnOnce(O) -> bool`. It succeeds if `f` returns true on the output of `p`, failing otherwise.

`ident(f)` matches a valid identifier that is not a keyword (as denoted by the function `f` returning false).

`tag(s)` takes `s: &str` and parses this string character by character. Returns the original string upon success.

`not(c)` does the same as `match_char(c)`, but fails if the character is matched.

`opt(p)` makes `p` optional, and is the same as `Or(Then(p, pure(())), pure(()))`.

`ws(p)` consumes trailing whitespace after parsing `p`.

`dynamic(p)` converts `p` into a heap-allocated dynamically dispatched version of `p`. This is useful for recursion.

4.3.4 The Parsler Trait

The previous sections introduce that the primitive and auxiliary parsers have been introduced which show the building blocks of complex parsers written in this library, and how Parslers does not use the monadic parser. This section wraps up the user interface by introducing the trait that all of these combinators implement, the `Parsler` trait. It provides the trait methods that allow parsers to be tested and compiled into an untyped AST. Despite its importance to the library, the definition is simple:

```
1 trait Parsler {
2     type Output;
3
4     fn parse(&self, input: &mut Chars) -> Result<Self::Output, String>;
5     fn compile(&self, info: &mut CompileContext) -> ast::Parser;
6
7     fn ap<B, F, P>(self, p: P) -> Ap<B, F, Self, P>
8     where
9         Self: Sized,
10         P: Parsler<Output = F>,
11         F: FnOnce(Self::Output) -> B,
12     {
13         Ap(self, p)
14     }
15
16     fn map<B, F>(self, f: F) -> Map<B, Self, F>
17     where
18         Self: Sized,
19         F: FnOnce(Self::Output) -> B + Clone + Reflect,
```

```

20 {
21     Map(self, f)
22 }
23
24 fn then<P>(self, p: P) -> Then<Self, P>
25 where
26     Self: Sized,
27     P: Parsler,
28 {
29     Then(self, p)
30 }
31
32 fn or<P>(self, p: P) -> Or<Self, P>
33 where
34     Self: Sized,
35     P: Parsler<Output = Self::Output>,
36 {
37     Or(self, p)
38 }
39
40 fn before<P>(self, p: P) -> Before<Self, P>
41 where
42     Self: Sized,
43     P: Parsler,
44 {
45     Before(self, p)
46 }
47
48 fn attempt(self) -> Attempt<Self>
49 where
50     Self: Sized,
51 {
52     Attempt(self)
53 }
54
55 fn look(self) -> Look<Self>
56 where
57     Self: Sized,
58 {
59     Look(self)
60 }
61
62 fn neg_look(self) -> NegLook<Self>
63 where
64     Self: Sized,
65 {
66     NegLook(self)
67 }
68 }

```

Every parser combinator in Parslers must implement the Parsler trait. It is built up of three components:

1. **type Output**; represents what the parsers produce. This is vital for making sure that combined parsers are well-formed upon creation.
2. **fn parse(&self, input: &mut Chars) -> Result<Self::Output, String>**; must be implemented for each combinator separately, and is used as a debug parser before code generation.
3. **fn compile(&self, info: &mut CompileContext) -> ast::Parser**; must also be implemented for each combinator, and describes how the parser can be compiled into the backend abstract syntax tree for generating native code.
4. `ap`; `map`; `then`; `or`; `before`; `attempt`; `look`; `negLook` merely act as chained combinators, and are implemented for all cases. They allow users to chain parsers together instead of relying on deeply nested structs, which can be hard to read.

Notice that for the `map` function, it takes a function `F: FnOnce(Self::Output) -> B + Clone + Reflect`. This is a function that can only be called once, can be cloned, and implements a `Reflect` trait. I will discuss this trait in the next section, but it is what allows us to achieve runtime reflection.

Let's take a look at the simplest implementor of Parsler, the `Pure` combinator:

```
1  #[derive(Clone, Debug)]
2  struct Pure<A: Reflect>(A);
3
4  impl<A: Clone + Reflect> Parsler for Pure<A> {
5      type Output = A;
6      fn parse(&self, _input: &mut Chars) -> Result<Self::Output, String> {
7          let val = self.0.clone();
8          Ok(val)
9      }
10
11     fn compile(&self, _info: &mut CompileContext) -> ast::Parser {
12         ast::Parser::Pure(ast::PureVal::Val(self.0.to_string()))
13     }
14 }
15
16 fn pure<A: Reflect>(a: A) -> Pure<A> {
17     Pure(a)
18 }
```

All `Pure` does as a combinator is return a value. It does not consume any input, and always succeeds. Notice how `A` is also required to implement `Reflect`, in order to retrieve the definition of the value at compile time.

Notice that the `parse()` function takes a mutable reference to a `Chars` struct. `Chars` is an iterator that produces characters from a string. Because Parslers does not implicitly backtrack, we only need

to pass around a mutable reference to this iterator. Therefore, the `parse()` function does not need to return the remaining unparsed string.

For combinators that enable backtracking, like `Attempt`, `Chars` can be cloned. This produces a new iterator at the same starting point as the source iterator. We use this functionality in the `Attempt` operator as follows:

```
1 impl<P1: Parsler> Parsler for Attempt<P1> {
2     type Output = P1::Output;
3     fn parse(&self, input: &mut Chars) -> Result<Self::Output, String> {
4         let copied_input = input.clone();
5         let a = self.0.parse(input);
6         if a.is_err() {
7             *input = copied_input;
8         }
9         a
10    }
11    // ...
12 }
```

So `Attempt` will clone `Chars` before parsing it to the inner parser. If the inner parser fails, it will then backtrack by replacing the original input with the cloned input. In this way, we obtain backtracking functionality.

4.4 The Reflect Trait

Parslers is a staged selective parser combinator library. The staged aspect of Parslers is derived from the fact that Parslers compiles user-created parsers during the build stage of Rust's compilation pipeline. The native code that Parslers generates is placed into a module for the user to include in their Rust programs. This poses an interesting architectural problem, as functions are not inspectable during the runtime of Rust programs.

If Parslers is to achieve staged compilation, it must also place functions and values defined by users in the `build.rs` file into the generated Rust code. This is equivalent to runtime reflection to the `build.rs` file. Although Rust does not support this directly, work has been done by other libraries needing runtime reflection. An example of this is the game engine `Bevy` (Anderson 2020), which uses the helper library `bevy_reflect` in order to achieve runtime reflection of values. However, there is no equivalent library in Rust that allows for runtime reflection of functions. In order to make staged compilation viable for Parslers, it must also support runtime reflection of functions, and allow function definitions to survive into the generated native Rust code. This section details how Parslers achieves this.

The `Reflect` trait is how we are able to achieve runtime-reflection:

```
1 trait Reflect {
2     fn name(&self) -> &'static str {
3         std::any::type_name::<Self>()
4     }
5
6     fn reflect(&self) -> String;
```

```
7 }
```

The purpose of this trait is to have the definition of the value available at runtime. The `to_string()` method gives us the value definition as a string, which can then be inserted into the generated code during the native code compilation step. It is already defined for all primitives and compound types provided in Rust. It's implementation for `char` is rather simple:

```
1 impl Reflect for char {  
2     fn reflect(&self) -> String {  
3         format!("{}", self)  
4     }  
5 }
```

However, it gets a bit more complicated for functions. This is because functions are black boxes in Rust, and can not be inspected at compile time or runtime. To make matters worse, you cannot implement traits over specific functions as of yet. To get around this, Parslers comes with an attribute macro `#[reflect]`. This macro replaces the function with a struct, and then implements `FnOnce` for that struct. The semantics of the function are the same before and after this transformation, but now that it is a struct, we can implement the `Reflect` trait.

To the user of this macro, they merely add the attribute macro above their function like so:

```
1 #[reflect]  
2 fn is_a(c: char) -> bool {  
3     c == 'a'  
4 }
```

But at compile time, the compiler will expand the macro into the following code:

```
1 #[allow(non_camel_case_types)]  
2 #[allow(incorrect_ident_case)]  
3 struct is_a;  
4  
5 impl FnOnce<char> for is_a {  
6     type Output = bool;  
7     extern "rust-call" fn call_once(self, (c,): (char,)) -> bool {  
8         c == 'a'  
9     }  
10 }  
11 impl Reflect for is_a {  
12     fn reflect(&self) -> String {  
13         "fn is_a(c : char) -> bool { c == 'a' }".to_owned()  
14     }  
15 }
```


Notice that because we implement `FnOnce` for `is_a`, it can still be called as a function in the same way as before:

```
1 println!("{}", is_a('b')); // false
```

However now we are able to read the body of the function at runtime, thanks to the automatic implementation of `Reflect`.

There is one caveat to automatically implementing the `Reflect` trait for functions. The user may want to implement a function that returns a function, to be applied to another value. An example is the parser `match_char('a').map(append).ap(pure(vec![]))`. This parser matches the character `'a'`, converts it into a function that appends `'a'` to a `Vec`, and then applies it to an empty vec. In other words, it parses a singleton list of the character `'a'`. In Rust, the most natural way to implement `append` is a function like the following:

```
1 fn append<A: 'static>(a: A) -> impl FnOnce(Vec<A>) -> Vec<A> {
2     move |mut v| {
3         v.push(a);
4         v
5     }
6 }
```

The `impl FnOnce(Vec<A>) -> Vec<A>` signifies that we will return a type that implements a function in the form `Vec<A> -> Vec<A>`. This is because Rust does not allow naming the concrete type of closures. The Rust compiler will figure out what the concrete type of the closure is internally, and use that to resolve the `impl` return type. However, if we look at the macro-expansion for the `FnOnce` trait implemented on this function by `#[reflect]` we see the following:

```
1 impl<A: 'static> FnOnce<(A,)> for append {
2     type Output = impl FnOnce(Vec<A>) -> Vec<A>; // impl in the associated type
3     extern "rust-call" fn call_once(self, (a,): (A,)) -> Self::Output {
4         move |mut v| {
5             v.push(a);
6             v
7         }
8     }
9 }
```

Unfortunately, Rust does not yet allow `impl` occurring in the associated type of a trait implementation. There are two options for resolving this issue:

1. Use dynamic dispatch and heap-allocate the return function.
2. Use the unstable `impl_trait_in_assoc_type`.

4.4.1 Using Dynamic Dispatch

With option 1, the user would have to reimplement the `append()` function to look like this:

```
1 fn append<A: 'static>(a: A) -> Box<dyn FnOnce(Vec<A>) -> Vec<A>> {
2     Box::new(move |mut v| {
3         v.push(a);
4         v
5     })
6 }
```

Because dynamic dispatch is used on a heap-allocated value, we have transformed the returned opaque type into a concrete type. The type `Box<dyn FnOnce(Vec<A>) -> Vec<A>>` can be used without requiring the use of a new unstable language feature. This is not without cost however, as Rust does not yet support whole program devirtualisation, the implementation of the returned function cannot be optimised by the Rust compiler. It would not be acceptable for Parslers to force users to use dynamic dispatch, as it would violate the performance guarantee. In order to account for this, the `##[reflect]` macro also comes with the ability to unbox a value in the code-generation step. The way to activate this feature is to write `##[reflect(unbox)]` instead. After macro-expansion with this feature, `##[reflect]` generates:

```
1 impl<A: 'static> FnOnce<(A,)> for append {
2     type Output = Box<dyn FnOnce(Vec<A>) -> Vec<A>>;
3     extern "rust-call" fn call_once(self, (a,): (A,)) -> Self::Output {
4         Box::new(move |mut v| {
5             v.push(a);
6             v
7         })
8     }
9 }
10 impl Reflect for append {
11     fn name(&self) -> &'static str {
12         "append"
13     }
14     fn reflect(&self) -> String {
15         "fn append<A:'static>(a : A,) -> impl FnOnce(Vec<A>) -> Vec<A> {
16             move | mut v | { v.push(a) ; v }
17         }".to_owned()
18     }
19 }
```

Notice that in the implementation of `FnOnce`, the associated type `Output` still remains `Box<dyn FnOnce(Vec<A>) -> Vec<A>>`. However, the implementation of `Reflect::reflect` the generated code has removed the `Box<dyn ...>` type and replaced it with `impl ...`, and removed the box in the function body. This means that, with the feature enabled, the generated Rust code will use static dispatch instead of dynamic dispatch. This resolves Parslers' performance requirement.

4.4.2 Returning to Static Dispatch

Although the `#[reflect(unbox)]` macro works for many cases, it is not perfect. For example, it is completely feasible for the user to write a parser like:

```
1  #[reflect(unbox)]
2  fn f1(c: char) -> Box<dyn FnOnce() -> char> {
3      Box::new(move || c)
4  }
5  #[reflect]
6  fn f2(f: Box<dyn FnOnce() -> char>) -> char {
7      f()
8  }
9  fn parser() -> impl Parsler<Output = char> {
10     pure('a').map(f1).map(f2)
11 }
```

As the user writes this parser, the Rust type system sees that `f1` produces a boxed function, and `f2` consumes a boxed function. However, the generated code of `f1` will generate an unboxed closure. This will cause the parser to fail its type check.

The alternative to unboxing the type at the code-generation step is to use the compiler feature `impl_trait_in_assoc_type`. This feature expands Rust's type inference capabilities to resolve opaque types from their occurrences elsewhere in the implementation. With this compiler feature enabled, we can return unboxed closures:

```
1  #[reflect]
2  pub fn append<A: 'static>(a: A) -> impl FnOnce(Vec<A>) -> Vec<A> {
3      move |mut v| {
4          v.push(a);
5          v
6      }
7  }
8  // Macro expanded code:
9  pub struct append;
10 impl<A: 'static> FnOnce<(A),> for append {
11     type Output = impl FnOnce(Vec<A>) -> Vec<A>; // Opaque type declared
12     extern "rust-call" fn call_once(self, (a,)): (A,)) -> Self::Output {
13         // Opaque type resolved
14         move |mut v| {
15             v.push(a);
16             v
17         }
18     }
19 }
```

4.5 Lazy Parsers

As discussed in the section Requirements, one of Parslers main goals is ensuring that users cannot write a malformed parser at the compile-time of the build.rs file. The section Parser Definition

Interface introduces building parsers as typed structures instead of higher-order functions. This is for the purpose of translating the parser into native Rust code, which would not be possible by composing higher-order functions. However, representing parsers as typed structures poses a new problem of representing recursion within these structures. This section introduces one of the parser combinators used to achieve recursion in parsers, which will be expanded upon by Named Parsers.

If the user attempts to create the following parser:

$$\langle \text{some_a} \rangle ::= 'a' \langle \text{some_a} \rangle$$
$$| 'a'$$

That is, to parse at least one 'a' character, then a parser like this would seem reasonable:

```
1 fn some_a() -> impl Parsler<Output = ()> {
2     match_char('a').then(opt(some_a()))
3 }
```

However, the Rust compiler will throw the error:

```
1 warning: function cannot return without recursing
2   --> parsers-benchmark/build.rs:153:1
3     |
4 153 | fn some_a() -> impl Parsler<Output = ()> {
5     | ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^ cannot return without recursing
6 154 |     match_char('a').then(opt(some_a()))
7     |                               ----- recursive call site
8     |
9     = help: a 'loop' may express intention better if this is on purpose
10    = note: '#[warn(unconditional_recursion)]' on by default
11
12 error[E0720]: cannot resolve opaque type
13   --> parsers-benchmark/build.rs:153:16
14     |
15 153 | fn some_a() -> impl Parsler<Output = ()> {
16     |           ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^ recursive opaque type
17 154 |     match_char('a').then(opt(some_a()))
```

There are two problems here. Firstly, as the Rust compiler correctly points out, the function recurses infinitely. This is because our parsers are purely structural, and so calling this function will cause it to recurse infinitely. We are also posed with the problem that Rust's type system cannot handle infinite recursion, which is the second error message. Without some help, the user will not be able to write recursive parsers.

In order to fix the first problem, infinite recursion, we introduce the `LazyParser`. It has a simple definition:

```

1 pub struct LazyParser<P, F>(std::cell::OnceCell<P>, F)
2 where
3     P: Parsler,
4     F: FnOnce() -> P + Clone;

```

`LazyParser` a structure containing two values, a `OnceCell`, and a function, `F`. The `OnceCell` is a value that can be created as uninitialised, and then initialised with a function at a later time. `F` is a function that produces a parser. This way, we have a way of delaying initializing a parser until we need to use it. As we can see in its implementation of the `Parsler` trait:

```

1 impl<P: Parsler, F: FnOnce() -> P + Clone> Parsler for LazyParser<P, F> {
2     type Output = P::Output;
3
4     fn parse(&self, input: &mut Chars) -> Result<Self::Output, String> {
5         self.0.get_or_init(self.1.clone()).parse(input)
6     }
7
8     fn compile(&self, info: &mut CompileContext) -> ast::Parser {
9         self.0.get_or_init(self.1.clone()).compile(info)
10    }
11 }

```

We can simply initialise the parser as of when we need it using the `get_or_init()` function. This solves the first problem, as we lazily initialize our parsers. However, there is still the problem of infinite type recursion, because the full type of the underlying parser behind `LazyParser` is still visible in `P`. To solve this, we use Rust's polymorphism:

```

1 pub fn dynamic<P>(p: P) -> Rc<dyn Parsler<Output = P::Output>>
2 where
3     P: Parsler + 'static,
4 {
5     Rc::new(p)
6 }

```

`Rc<dyn Parsler<Output = P::Output>>` is a reference-counted smart pointer to an object that implements `Parsler` and whose output is the same as the input parser. This maintains the type information we need while erasing the true type of the parser. Using these two concepts, we can rewrite the user's original parser into:

```

1 fn some_a() -> impl Parsler<Output = ()> {
2     match_char('a').then(opt(dynamic(LazyParser::new(some_a))))
3 }

```

This is enough to get our parser to compile. However, there is one more step needed to achieve recursive parsing, which will be the topic of the next section. Also note that because Parslers uses staged compilation, using polymorphism and smart pointers has no effect on the performance of the resulting parser. These will be erased in the final parser into functions and called by name.

4.6 Named Parsers

The Lazy Parsers section introduces the first tool that can be used to generate recursive parsers by masking the concrete type of the parser behind dynamic dispatch. However, one more Looking back to the implementation of `LazyParser`, notice that both the `parse()` and `compile()` functions call the respective methods on the wrapped parser. For `parse()`, this is fine because the underlying parser may fail, and not recurse further. If the user writes a left-recursive parser, for example:

```
1 <some_a> ::= <some_a> "a"
```

Then the parser will recurse infinitely. However, this is expected of all top-down parsers, and so is not a problem. However, the `compile()` function will never fail if the parser is recursive, and so we are posed with the problem of identifying parsers that have already been compiled. In the original Staged Selective Parser Combinators paper (Willis, Wu & Pickering 2020), the author makes use of Haskell's **StableNames** to identify parsers that are equal. As Rust has no equivalence of **StableNames**, naming parsers manually is the necessary fallback. And so a new parser is introduced, the `NamedParser`:

```
1 pub struct NamedParser<P: Parsler>(pub P, pub String);
```

The `String` acts as an identifier that can be used to catch otherwise infinitely recursive calls to `compile()`:

```
1 impl<P: Parsler> Parsler for NamedParser<P> {
2     type Output = P::Output;
3     //...
4     fn compile(&self, info: &mut CompileContext) -> ast::Parser {
5         let hash = hash(self);
6         let name = format!("{}", self.1, hash);
7         match info.register_parser(&name) {
8             NamedParserStatus::Registered => ast::Parser::Ident(name),
9             NamedParserStatus::Unregistered => {
10                 let p = self.0.compile(info);
11                 info.insert_parser::<P>(&name, p.clone());
12                 ast::Parser::Ident(name)
13             }
14         }
15     }
16 }
```

Notice that the parser is also hashed to extend the identifier. We do this because manually naming parsers poses the risk of the user giving two different parsers the same name. Hashing the parsers

adds an extra safeguard against this. The call to `info.register_parser()` registers the parser's name, and then returns if it was previously registered. If it was previously registered, then we just reference the parser's name. If it has not been seen before, then we compile the wrapped parser, insert it into the `CompileContext`, and then return the parser's name. This way, we catch potential infinite recursion in the compilation of self-referential parsers.

Because naming a parser often implies recursion, the library contains a helper function `name` :

```
1 pub fn name<P: Parsler + 'static>(  
2     name: &str,  
3     f: impl FnOnce() -> P + Clone + 'static,  
4 ) -> Rc<dyn Parsler<Output = P::Output>> {  
5     Rc::new(NamedParser(LazyParser::new(f), name.to_owned()))  
6 }
```

Finally, the recursive parser from before can be re-written as:

```
1 fn some_a() -> impl Parsler<Output = ()> + Clone {  
2     name("some_a", || match_char('a').then(opt(some_a())))  
3 }
```

4.7 Looping Parsers

It is common for a parser combinator library to implement the `some` and `many` combinators. These are simple combinators to implement:

```
1 pub fn some_rev<P>(p: P) -> Rc<dyn Parsler<Output = Vec<P::Output>>>  
2 where  
3     P: Parsler + Clone + 'static,  
4     P::Output: Clone + 'static + Reflect,  
5 {  
6     name("some_rev", || p.clone().map(append).ap(many_rev(p)))  
7 }  
8  
9 pub fn many_rev<P>(p: P) -> Rc<dyn Parsler<Output = Vec<P::Output>>>  
10 where  
11     P: Parsler + Clone + 'static,  
12     P::Output: Clone + 'static + Reflect,  
13 {  
14     name("many_rev", || some_rev(p).or(pure(vec![])))  
15 }
```

As can be seen above, they both rely on mutual recursion. Many functional programming languages support tail-call optimisation, which identify and replace recursion with imperative loops (Muchnick et al. 1997). Unfortunately, Rust does not officially support tail-call optimisation yet. As a result,

these two parsers will mutually recurse, which will incur the costs associated with recursion. Furthermore, `Vec` (equivalent to `ArrayList` in Java) grows from the end, making appending to the beginning of the vector expensive. Therefore, `some_rev()` and `many_rev()` will give the resulting parsed values in reverse-order from their parse order, and will require the user to reverse the list afterwards. This makes having a looping combinator necessary for writing performant parsers.

In order to accommodate for the need of loops within the parser, we introduce the `Loop` parser:

```
1 pub struct Loop<F, P0, PN>(P0, PN)
2 where
3     F: FnOnce(P0::Output) -> P0::Output,
4     P0: Parsler,
5     PN: Parsler<Output = F>;
```

This is a combinator that takes two parsers:

`P0` : A parser which produces the initial value.

`PN` : A parser which will run in a loop, and produces a function that transforms the output of previous iterations.

This becomes more clear when looking at the implementation of the parser:

```
1 impl<F, P0, PN> Parsler for Loop<F, P0, PN>
2 where
3     F: FnOnce(P0::Output) -> P0::Output,
4     P0: Parsler,
5     PN: Parsler<Output = F>,
6 {
7     type Output = P0::Output;
8
9     fn parse(&self, input: &mut Chars) -> Result<Self::Output, String> {
10         let mut v = self.0.parse(input)?;
11         while let Ok(f) = self.1.parse(input) {
12             v = f(v);
13         }
14         Ok(v)
15     }
16     // ...
17 }
```

Implementing the `some()` and `many()` combinators then becomes simple:

```
1 pub fn some<P>(p: P) -> Rc<dyn Parsler<Output = Vec<P::Output>>>
2 where
3     P: Parsler + Clone + 'static,
4     P::Output: Clone + 'static,
```



```

5 {
6     dynamic(Loop(p.clone().map(singleton), pure(append).ap(p)))
7 }
8
9 pub fn many<P>(p: P) -> Rc<dyn Parsler<Output = Vec<P::Output>>>
10 where
11     P: Parsler + Clone + 'static,
12     P::Output: Clone + 'static + Reflect,
13 {
14     dynamic(Loop(pure(vec![]), Map(p, append)))
15 }

```

The only difference now is that the recursion has been removed from the implementation of the combinators. Note that the generality of the `Loop` operator allows writing any parser. For example, we can write a parser that counts the occurrences of `'a'` in a string:

```

1 fn count_as() -> impl Parsler<Output = usize> {
2     Loop(pure(0), match_char('a').then(pure(inc)))
3 }

```

4.8 Alternative Parser

4.9 Untyped AST

Parslers' backend requires optimisations on the user-written parser's structure during the runtime of the `build.rs` module. This is done by conditionally matching on the structure of the parse tree, and applying substitutions that would improve the performance of the parser. As Rust follows the functional programming paradigm, the best way to perform this structural analysis would be via pattern matching. Unfortunately, there are several things preventing pattern matching from being used in this context.

Recall that, in its current form, the user writes parsers by combining typed structs together. This approach offers the great benefit that the user cannot write malformed parsers at compile time. Looking at the implementation of the `Or` combinator provides a concrete example:

```

1 pub struct Or<P1: Parsler, P2: Parsler<Output = P1::Output>>(pub P1, pub P2);

```

As we can see, there is a compile-time assertion in the definition of `Or` that the outputs of `P1` and `P2` (that is, the two alternative parsers) are the same. This would prevent a user from writing a parser that could produce two different types if `P1` succeeded or if `P2` succeeded.

However, the downside to this compile-time guarantee is that we lose the ability to perform structural pattern-matching on the parsers, as pattern matching can only be performed on enums. Furthermore, combinators such as `LazyParser` and `dyn Parsler` are not inspectable at runtime for pattern matching, as the internal structures are opaque. Furthermore, using typed enums is not possible as Rust does not currently support higher-kinded types, which would enable generic types to be defined in the individual enumerations. Although using Rust's Generic Associated Types could

allow for this, experiments performed in the writing of Parslers found that this method generates too much complexity on the side of the user.

The key insight to solving this issue is to restrict optimisations performed on parsers such that they are type-invariant. That is to say that no optimisations performed on a parser can change the output type of the parser. It will become clear that this is more of an observation of the current optimisations available rather than a limitation of the optimisations we can select from in later chapters. With this observation, we can begin to erase the types of the parsers at runtime and translate the parsers into an untyped AST using Rust's enums.

```
1 pub struct AnalysedParser {
2     pub output_used: bool,
3     pub returns_func: bool,
4     pub parser: Parser,
5 }
6 pub enum Parser {
7     Ident(String),
8     Pure(PureVal),
9     Satisfy(Func),
10    Try(Box<AnalysedParser>),
11    Look(Box<AnalysedParser>),
12    NegLook(Box<AnalysedParser>),
13    Ap(Box<AnalysedParser>, Box<AnalysedParser>),
14    Then(Box<AnalysedParser>, Box<AnalysedParser>),
15    Before(Box<AnalysedParser>, Box<AnalysedParser>),
16    Or(Box<AnalysedParser>, Box<AnalysedParser>),
17    Recognise(Box<AnalysedParser>),
18    Empty,
19    Branch(
20        Box<AnalysedParser>,
21        Box<AnalysedParser>,
22        Box<AnalysedParser>,
23    ),
24    Loop(Box<AnalysedParser>, Box<AnalysedParser>),
25 }
```

Translation from a typed structure to an untyped AST is the main purpose of the `Parsler::compile` method. Now that we have an untyped AST in the form of a Rust enum, we may perform structural optimisations on the parser. Note that this is bundled as a part of the backend of Parslers, and should never be a concern of the library user.

As this method requires optimisations on the `Parser` enum to be type-invariant, the library maintainer is now responsible for ensuring that this invariance propagates throughout all optimisation passes. However, this has not yet been a problem in the writing of Parslers.

4.10 Backend Optimisations

Optimisation plays a key part in Parslers' performance characteristics. In its current state, Parslers performs two major optimisations passes: Lawful Reduction, and Usage Analysis.

4.10.1 Lawful Reduction

Research done on parser combinators show that the combinators that form parsers are a type of algebra (Willis, Wu & Pickering 2020), and have equivalence laws. Below are the laws that apply to the combinators presented in Parslers:

$$\begin{array}{ll}
Ap(Pure(f), Pure(x)) = Pure(f(x)) & (1) \\
Ap(u, Pure(x)) = Ap(Pure(|f|f(x)), u) & (2) \\
Ap(u, Ap(v, w)) = Ap(Pure(compose), Ap(u, Ap(v, w))) & (3) \\
Or(Or(p, q), r) = Or(p, Or(q, r)) & (4) \\
Or(Empty, p) = Or(p, Empty) = p & (5) \\
Ap(Empty, p) = Empty & (6) \\
Or(Pure(x), p) = Pure(x) & (7) \\
Branch(Pure(Left(x)), p, q) = Ap(p, Pure(x)) & (8) \\
Branch(Pure(Right(y)), p, q) = Ap(q, Pure(y)) & (9) \\
Branch(Then(x, y), p, q) = Then(x, Branch(y, p, q)) & (10) \\
Try(NegLook(p)) = NegLook(p) & (11) \\
Look(Empty) = Empty & (12) \\
Look(Pure(x)) = Pure(x) & (13) \\
NegLook(Empty) = Pure(()) & (14) \\
NegLook(Pure(x)) = Empty & (15) \\
Look(Look(p)) = Look(p) & (16) \\
Or(Look(p), Look(q)) = Look(Or(Try(p), q)) & (17) \\
NegLook(NegLook(p)) = Look(p) & (18) \\
NegLook(Or(Try(p), q)) = Then(NegLook(p), NegLook(q)) & (19) \\
Or(NegLook(p), NegLook(q)) = NegLook(Then(Look(p), Look(q))) & (20)
\end{array}$$

The laws stated above allow Parslers to optimise user-written parsers to a great degree. Take, for example, the following parser:

```

1  #[reflect]
2  fn is_a(c: char) -> bool { c == 'a' }
3  #[reflect]
4  fn neg(b: bool) -> bool { !b }
5  pure('a').map(is_a).map(neg) // Always returns false

```

When written out as an AST, it becomes `Ap(neg, Ap(Pure(is_a), Pure('a')))`. Using the lawful reduction rules stated above, the AST can be reduced:

```

Ap(neg, Ap(Pure(is_a), Pure('a')))
    Ap(neg, Pure(is_a('a')))
        Pure(neg(is_a('a')))

```

Parslers performs the lawful reduction optimisation in a bottom-up recursive pass through the parser AST, applying substitutions where possible and then retrying if a match has been found (to ensure any new substitution found after applying a reduction is checked).

4.10.2 Output Usage Analysis

Usage analysis is the most powerful optimisation that Parslers has in its toolkit. It relies on the observation that certain combinators discard the returned value after parsing has completed. Take a look at the following parser:

```

1  #[reflect]
2  fn digit(c: char) -> bool { ('0'..'9').contains(&c) }
3  #[reflect]
4  fn parse_usize(s: String) -> usize { s.parse().unwrap() }
5  fn uint() -> impl Parsler<Output = usize> {
6      Recognise(some(Satisfy(digit))).map(parse_usize)
7  }

```

When broken down, `some(Satisfy(digit))` returns a `Vec<char>` of at least one, and then `Recognise(some(digit))` discards the output of `some(Satisfy(digit))`, and returns its consumed string. This string is then parsed into a `usize`. That `Vec<char>` requires allocation, reallocation (when the `Vec` grows), and the pushes to the `Vec` after each successful digit parsed. These clock cycles are wasted with no use. In some cases, the compiler recognises this and can discard wasted computation, but this is not perfect, as will be seen in the later chapter discussing performance benchmarks.

In order to account for this, Parslers incorporates usage analysis. This process performs a top-down pass through the untyped parser AST, and marks subtrees as unused if not used by their parents. During code-generation, if a parser subtree's output is unused, it will return the unit struct `()`, in Rust). Because this happens all the way down the parse tree, `some(Satisfy(digit))` does not return any value, and we will see that this can have a great effect in the next section.

The `Branch(b, l, r)` combinator is the only one which prevents usage analysis from propagating into the combinator subtrees. Because `l` and `r` are conditionally executed based upon the result of `b`, `b` must always be returned.

```

1  Then(p, q)      // p is discarded.
2  Before(p, q)    // q is discarded.
3  Recognise(p)    // p is discarded.
4  Branch(b, l, r) // b is always used.

```

4.10.3 Validation Parsers

An interesting side effect of the usage analysis optimisation can be found in validating parsers. If a user has a very large parser (for example, a JSON parser), but instead of wanting to parse text into structured data just validate that the text is a valid instance of the grammar. In this case, it suffices to simply augment the parser `p` into `p.then(pure())`. The usage analysis optimiser will see that the result of `p` is discarded, and deeply remove any computation and returned values. As can be seen in the benchmarking chapter, this provides a very large boost to performance.

4.10.4 Analysis of Recursive Parsers

In order for the benefits of analysis optimisation to be fully realised, it must be aggressive. Given a parser that is both recursive and unused (for example, a JSON validating parser), usage analysis must also identify recursion points and mark them to be revisited by the optimiser. Recall that recursive parsers are marked via a named hole. When it is identified that the result of this named parser is unused, the recursion point is renamed with a tag that marks it is unused, and the original parser is added to a set of parsers that are to be deeply optimised to reflect that the return value is

unused.

Optimisation on the renamed, unused parser may result in more named parsers being marked unused. Because of this, the optimisation is run repeatedly until there are no more named parser marked unused that do not already exist in the set of optimised parsers. The optimiser is guaranteed to terminate as there can only be a single unused variant of a parser for each of the named parsers, and so the loop is bounded.

4.11 Code Generation

The Untyped AST section introduces the intermediate representation of parsers as a Rust enum, which reflect the primitive combinators but strips the type safety that structs give us. Then, the Backend Optimisations details the optimisations that the untyped AST can take advantage of. This section concerns itself with how Parslers goes from the untyped AST back into native Rust code, which completes the loop.

To begin with, Parsing and Generating Rust Code with Syn and Quote will introduce the library that helps tokenise Rust code using a macro library called `quote`. Then, The type of Parsers will introduce how function definitions can be derived from user-generated parsers in the code-generation step.

4.11.1 Registering Functions

The section The Reflect Trait showed how Parslers achieves runtime reflection of function definitions using the `Reflect` trait and the `#[reflect]` macro. These functions must be added to the compile context so that they may be compiled into the native Rust code. Looking at the `CompileContext` struct, which was previously seen in the `Parsler::compile(&self, info: &mut CompileContext)` function:

```
1 pub struct CompileContext {
2     functions: HashMap<String, String>,
3     // ...
4 }
```

The `functions` field is a `HashMap<String, String>`. This stores the bodies of the functions, along with a generated name. The name is replaced and auto-generated as it is only the body of the function that needs to be unique, and Parslers is able to de-duplicate functions when they arise.

The `functions` field is only modified by three of the primitive combinators:

- `Pure`
- `Map`
- `Satisfy`

Now the usage in `Satisfy` and `Map` is not difficult, as both of these primitive combinators guarantee that they are constructed with a function type. Therefore, registering the function into the compile context is simple:

```

1  impl<B, P1, F> Parsler for Map<B, P1, F>
2  where
3      P1: Parsler,
4      F: FnOnce(P1::Output) -> B + Clone + Reflect,
5  {
6      type Output = B;
7      fn compile(&self, info: &mut CompileContext) -> ast::Parser {
8          let name = info.add_function(&self.1);
9          let p = self.0.compile(info);
10
11          ast::Parser::Ap(
12              Box::new(AnalysedParser::new(p)),
13              Box::new(AnalysedParser::new(ast::Parser::Pure(ast::PureVal::Func(
14                  ast::Func { name },
15                  )))),
16          )
17      }
18      // ...
19  }
20
21  impl<F: FnOnce(char) -> bool + Reflect + Clone> Parsler for Satisfy<F> {
22      type Output = char;
23
24      fn compile(&self, info: &mut CompileContext) -> ast::Parser {
25          let name = info.add_function(&self.0);
26
27          ast::Parser::Satisfy(ast::Func { name })
28      }
29      // ...
30  }

```

In these, as they both are restricted by `F: FnOnce(...) -> ...`, using `info.add_function(...)` is always safe. However, this is not the case with the `Pure` primitive combinator, because it does not restrict the type of value it is. This becomes clear when looking at `Pure`'s definition:

```

1  pub struct Pure<A: Reflect>(& A);

```

The fact that `Pure` can accept any value that implements `Reflect` means that a bit more logic is required to differentiate functions from values. To solve this, `Parslers` uses the `syn` crate to attempt to parse the reflected value held in `Pure` into a `ItemFn`. If `syn` can successfully parse the string returned by `Reflect::reflect(&self)` into a function, then it is safe to assume that it is a function, and can be registered the same way as in the implementations of the other two primitive combinators:

```

1  impl<A: Clone + Reflect> Parsler for Pure<A> {
2      type Output = A;
3

```

```

4  fn compile(&self, info: &mut CompileContext) -> ast::Parser {
5      if let Ok(_) = syn::parse_str::<syn::ItemFn>(&self.0.reflect()) {
6          let name = info.add_function(&self.0);
7          ast::Parser::Pure(ast::PureVal::Func(ast::Func { name }))
8      } else {
9          ast::Parser::Pure(ast::PureVal::Val(self.0.reflect()))
10     }
11 }
12 // ...
13 }

```

Now Parslers is able to register every function that can be found in the user-generated parsers.

The `CompileContext::add_function(&mut self, func: &impl Reflect)` method is incredibly simple thanks to Rust's `Entry` API:

```

1  impl CompileContext {
2      pub fn add_function(&mut self, func: &impl Reflect) -> String {
3          let num_functions = self.functions.len();
4          self.functions
5              .entry(func.reflect())
6              .or_insert(format!("{}", num_functions))
7              .clone()
8      }
9      // ...
10 }

```

This method gets the current number of functions registered (so that each unique function has a unique name) and then attempts to retrieve the unique function name from `functions` map. If this entry does not already exist, it inserts a new name with the number of registered functions and returns the new name.

During the generation of the native Rust code, the keys and values of `functions` are iterated over, and the previous name of the function body is replaced with the generated unique name. Then, the `quote` library is used to insert each of these functions into the generated Rust module. These functions can then be referenced and called within the generated parser at run-time. Both the `syn` and `quote` libraries will be discussed in the section `Parsing and Generating Rust Code with Syn and Quote`.

4.11.2 Parsing and Generating Rust Code with Syn and Quote

The section `Registering Functions` discussed how functions can survive the `build.rs` file into the generated binary, and expanded further on the purpose of the section `The Reflect Trait`. This section concerns itself with two useful tools from the Rust ecosystem that allow Parslers to parse and generate Rust code, the `syn` and `quote` libraries respectively.

The `syn` library is used in several places in Parslers. This library is important to Rust library writers because of the way that Rust's staging works. In Rust, procedural macros are expanded after the parse stage of the compiler. This stage produces tokens in the form of a `TokenStream`, which hold

no semantic information about the code that has been parsed. Furthermore, `TokenStream` is not given in the form of a Rust AST, mainly because procedural macros allow library writers to create domain specific languages that do not necessarily conform to the Rust syntax standard (such as a HTML template in the `typed_html` library). However, in cases where the incoming `TokenStream` is fully compliant with the Rust syntax standards, just using `TokenStream` can be cumbersome. This is where the `syn` library comes in, which allows parsing `TokenStream`s into a Rust-compliant AST.

Although `syn` is used throughout the codebase of Parslers, the best example of its use is in the renaming of functions as described in Registering Functions:

```
1 let functions = context
2   .functions
3   .iter()
4   .map(|(body, name)| {
5       let mut body = syn::parse_str::<syn::ItemFn>(body).unwrap();
6
7       body.sig.ident = syn::parse_str::<syn::Ident>(name).unwrap();
8
9       // ...
10  })
11  .collect::<Vec<TokenStream>>();
```

In the above code, `functions` is iterated over (the keys and values are of type `String`). Then, the function body is parsed into a `syn::ItemFn`, which is a full function in the form:

```
1 fn func(x: i32) -> i32 {
2     a + 1
3 }
```

Then, the new name of the function is parsed into a valid identifier. Now it becomes easy to rename the function by replacing the function signature's identifier with the new function name. The hidden code behind `// ...` is the next topic of interest, which converts the typed AST that `syn` provides back into a `TokenStream`.

The final step in code generation involves converting the untyped AST back into native Rust code. Currently, Rust does not have good first-class support for generating Rust code in an ergonomic manner. This could be solved by using string templating, but this is not ergonomic and puts too much cognitive load on the library programmer. To solve this, a macro library has been developed called `quote`. This is a simple and minimalist library with one core feature, the `quote!{...}` macro. Using this macro is quite simple.

For example, if a library developer wanted to develop a macro that generates a function to print a string, but in a shorter way than the standard `println!(...)`. To create a `pl!(...)` macro, it is quite simple:


```

1  #[proc_macro]
2  pub fn pl(input: proc_macro::TokenStream) -> proc_macro::TokenStream {
3      let input = proc_macro2::TokenStream::from(input);
4      quote::quote!{
5          println!(#input)
6      }.into()
7  }

```

Here it can be seen on lines 4-6 that the quote macro is used to tokenize rust-like syntax. It also performs syntax analysis at the compile-time of the macro. This allows many of the potential syntax bugs that could arise in the string templating approach to code generation. In Parslers, the use of `quote!{...}` can be seen in the function generation:

```

1  let functions = context
2      .functions
3      .iter()
4      .map(|(body, name)| {
5          let mut body = syn::parse_str::<syn::ItemFn>(body).unwrap();
6
7          body.sig.ident = syn::parse_str::<syn::Ident>(name).unwrap();
8
9          quote! {
10              #[inline(always)]
11              #body
12          }
13      })
14      .collect::<Vec<TokenStream>>();

```

Notice that `#[inline(always)]` is used here. This is because functions used in Parslers are often-times incredibly small, and inlining these functions can always be justified. It may be useful in the future to add inlining hints to the `#[reflect]` attribute macro, however, inlining every function results in faster parsers as found in the benchmarks.

4.11.3 The type of Parsers

Now the tools used to parse and generate valid Rust syntax have been discussed in the previous section, it is time to move on to how function definitions can be derived from user-defined parsers.

This is important as user-defined parsers should be translated into functions that can be called on some object. It is sensible for the type of these functions to match the type of the `Parsler::parse()` function to ensure cross-compatibility. Therefore, the type of the output parsing function should be `fn(input: &mut std::str::Chars) -> Result<..., &'static str>`.

Determining what the `...` in the output result type should be considered carefully. This should require runtime reflection of the parsers' output types to correctly identify and insert it into the missing space. Luckily, the Rust language's powerful monomorphization can be used to circumvent adding even more runtime reflection into Parslers.

The Rust standard library comes pre-packaged with the `std::any` module. This is a module dedicated to compile-time reflection and makes use of monomorphization in order to achieve compile-time reflection instead of runtime reflection. This can then be used to retrieve the full type name of an object at compile-time.

`std::any::type_name<T>()` -> `&'static str` function does not take in any values, but instead a function generic that is used to retrieve the name of the type. When inserting a named parser into the `CompileContext`, `Parslers` uses an internal method `CompileContext::insert_parser<P: Parsler>(&mut self, name: &str, parser: ast::Parser)`. This method does not take the user-defined parser, but instead the untyped AST. However, it does take a generic parameter `P: Parsler`. This allows the type name of `P::Output` to be retrieved at compile-time, and be used to generate the true type of the output parser, filling in the `Result<..., &'static str>` result type signature in the generated parser function.

4.11.4 Recursively Compiling Parsers

The sections `Registering Functions`, `Parsing and Generating Rust Code with Syn and Quote`, and `The type of Parsers` have laid out the building blocks to converting user-defined parsers into native Rust code. However, there is one more step required before this process is complete - generating the parser bodies themselves.

The code generation is performed via a top-down post-order traversal of the untyped parser AST. It is worth noting that compilation is split into two different functions: `compile_used()`, and `compile_unused()`. This is to account for the fact that the generated output code looks very different for parser combinators whose outputs are used and unused. Looking at the implementation of some of the primitive combinators shows that thanks to Rust's pattern-matching syntax, this stage of compilation is easy to perform.

Firstly, see the `Satisfy` variant of the untyped parser AST:

```

1  Parser::Satisfy(ident) => {
2      let ident = syn::parse_str::<syn::Expr>(&ident.name).unwrap();
3      quote! {
4          let old_input = input.clone();
5          input.next()
6          .ok_or("Found EOF when character was expected")
7          .and_then(|c|
8              if (#ident)(c) {
9                  Ok(c)
10             } else {
11                 *input = old_input;
12                 Err("Expected a specific character")
13             }
14         )
15     }
16 }
```

Notice two things:

1. The code within the `quote!{...}` macro looks very close to the `Parsler::parse(...)` method in the implementation of `Parsler` for `Satisfy`. This is frequently the case and

expected considering that the output generated code is Rust code.

2. `input` is never defined in the generated code but is used. This is possible because the whole generated parser mutates `input`, and every sub-parser is nested in a new scope level. Therefore, it is always valid to assume that `input` exists and is correct within its context so long as backtracking parsers handle the logic correctly.

Now the code generation logic for the `Satisfy` variant has been introduced, it is worth looking at the `compile_unused()` function to see where differences typically lie:

```
1 Parser::Satisfy(ident) => {
2     let ident = syn::parse_str::<syn::Expr>(&ident.name).unwrap();
3     quote! {
4         let old_input = input.clone();
5         input.next()
6         .ok_or("Found EOF when character was expected")
7         .and_then(|c|
8             if (#ident)(c) {
9                 Ok(())
10            } else {
11                *input = old_input;
12                Err("expected a specific character")
13            }
14        )
15    }
16 }
```

Notice how `Ok(c)` has now become `Ok(())`. As mentioned in Output Usage Analysis, this is an optimisation that allows for significant performance improvements. This will be expanded upon in the evaluation section.

5 EVALUATION

This section concerns itself with the performance evaluation of Parslers as a general-purpose parsing library.

5.1 Performance

5.1.1 Prior Discussion

Evaluating the performance of any general-purpose library is a problematic pursuit. Especially when benchmarks must be comparative to other libraries. This is because a user could implement a program in many ways within each library. Parser combinator libraries are no exception to this, as there are usually many different parsers that can conform to the same grammar, each with different performance characteristics.

Although great care was taken to be fair to each parser combinator library in each comparative benchmark, I am not an expert in the performance characteristics of these libraries. However, having written Parslers myself, I am an expert in its performance characteristics. For this reason, the Branflakes comparative benchmarks must not be confused with the fastest possible parsers in each library.

However, each of the libraries benchmarked in the Comparative JSON Benchmarks section have JSON parsers included with a benchmarking suite. These premade parsers were used and can be assumed to be strong comparative benchmarking platforms.

All benchmarks were run on the same machine with an Intel i5 4060k 3.5 GHz CPU and 22 GB of DDR3 1600 MHz RAM.

5.1.2 Comparative JSON Benchmarks

Throughput of Parsers Built with Different Parsing Libraries (JSON)

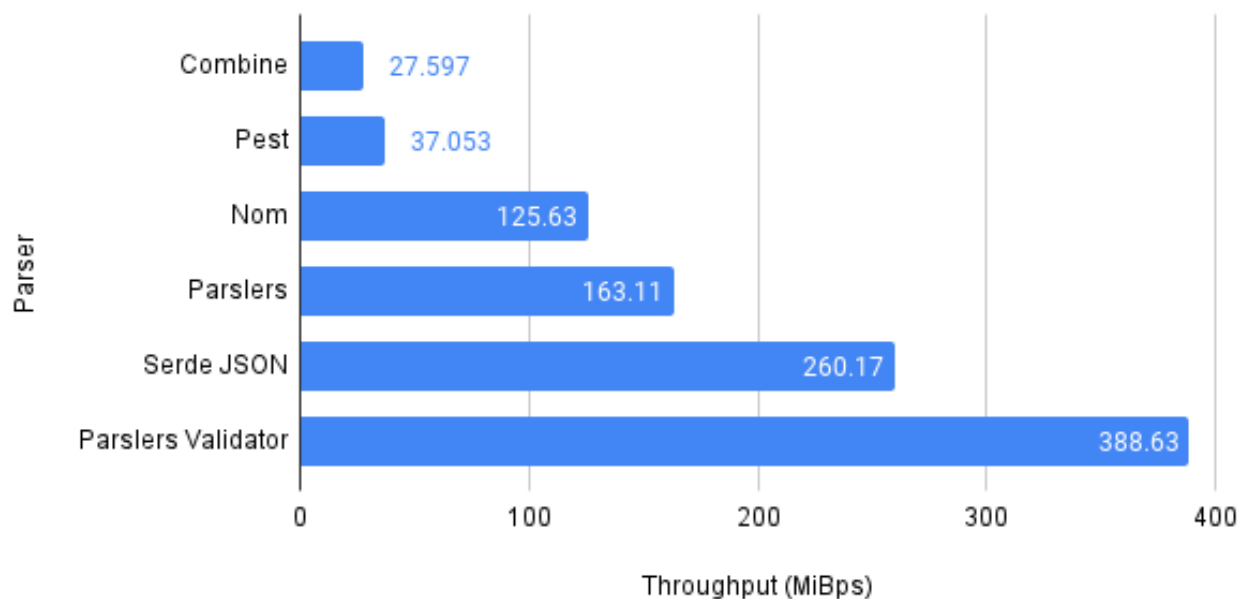


Figure 1. A plot of the throughput of parsers built with Parslers, Nom, Combine, Pest, and Serde JSON each provided by the benchmark suite of the respective library, benchmarked on the Canada JSON dataset.

JSON (JavaScript Object Notation) is one of the most-used formats for machine-to-machine com-

munication. It is commonly used as a part of client-server communication in Web APIs. This is because the grammar is quite simple, but expressive (Pezoa, Reutter, Suarez, Ugarte & Vrgoč 2016) and already has large-scale adoption and language support. This makes it a great candidate for benchmarking, and many general-purpose

In this benchmark, Parslers is put against Rust’s most commonly used general-purpose parsing libraries (Nom, Pest, and Combine), and the most commonly used hand-written parser for JSON (Serde JSON). Nom is generally considered the fastest general-purpose parsing library for the Rust programming language and is also a parser combinator library. As can be seen in 1, The Parslers library outperforms the general-purpose parsing libraries but does not outperform Serde JSON. This is reasonable, as Serde JSON is a hand-written library dedicated to parsing a specific grammar, and is heavily optimised to parse as fast as possible. It is worth noting, however, that the Parslers Validate benchmark does outperform Serde JSON, which does not currently have a way to validate JSON files rather than just parse them into structured JSON objects.

5.1.3 Comparative Branflakes Benchmarks

Throughput of Parsers Built with Different Parsing Libraries
(Branflakes)

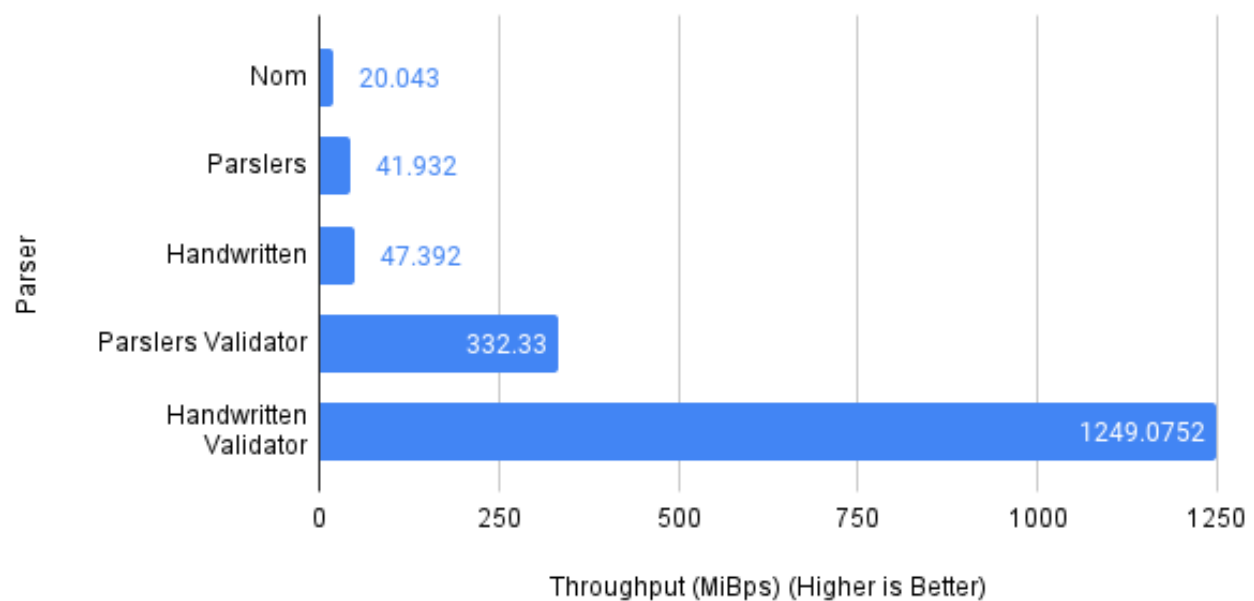


Figure 2. A plot of the throughput of parsers built with Parslers, Nom, a hand-written parser, and a hand-written validator each provided by the benchmark suite of the respective library, benchmarked on the Lost Kingdom Branflakes dataset.

Branflakes is a simple Turing-complete programming language. It is composed of an unsigned, 8-bit-valued tape memory, and a pointer to the current tape location. Values on this tape can be manipulated with eight total symbols:

- + Increment the value at the address of the current tape pointer by 1.
- Decrement the value at the address of the current tape pointer by 1.
- < Decrement the tape pointer by 1.
- > Increment the tape pointer by 1.

- . Print out the value at the address of the current tape pointer.
- , Read a character and set the value at the address of the current tape pointer to it.
- [If the value at the address of the current tape pointer is 0, then jump to the next matching].
-] If the value at the address of the current tape pointer is 0, then jump to the preceding matching [.

The simplicity of the grammar makes it a good benchmarking tool for comparing different parsers. Furthermore, it is a great grammar for writing an optimised hand-written parser and validator. For this benchmark, a hand-written parser was developed, lacking any recursion and is stack-based to be performant. This serves as a good comparison between hand-written parsers and Parslers in the context of a simple grammar.

Looking at the results in 2, the results show that Parslers is faster than Nom in this context as well. The handwritten parser is faster than Parslers, but surprisingly only by about 13%. However, with the optimisations available by knowing the grammar, the handwritten validator is faster than Parslers' validator by about 376%. It is worth noting that the handwritten validator has to be written separately from the actual handwritten parser, which means there is extra overhead in writing a separate structured parser and validator. This overhead does not exist for Parslers.

5.1.4 Staging and Optimisation Benchmarks

Parslers is a staged parser combinator library. However, each of the primitive parser combinators packaged with Parslers has the ability to parse before staging occurs. This makes the question of how much of a performance increase can be found in staging alone without any further optimisations applied. A second element of interest is how much different optimisations improve the performance of resulting parsers, and whether or not parsers for different grammars will benefit more from certain optimisations or others.

In order to test this, benchmarks were set up to compare different optimisation configurations, to be tested on both Branflakes and JSON parsers. These include:

- Unstaged
- Staged, but Unoptimised
- Reduction Analysis Only
- Usage Analysis Only
- Fully Optimised

From both the JSON and Branflakes benchmarks in 4 and 3, it is clear that staging is hugely beneficial to the performance of parsers, regardless of applied optimisations. This makes sense when considering that compilation is split up into compilation units, which are compiled and optimised separately. However, after staging large sections of the parser exist within the same function, and will be optimised as a single compilation unit. Furthermore, there is no dynamic dispatch in the staged parsers, whereas the unstaged parsers can heavily use dynamic dispatch in order to fix problems discussed in Lazy Parsers and Named Parsers.

Throughput of Parslers with Different Optimisations (JSON)

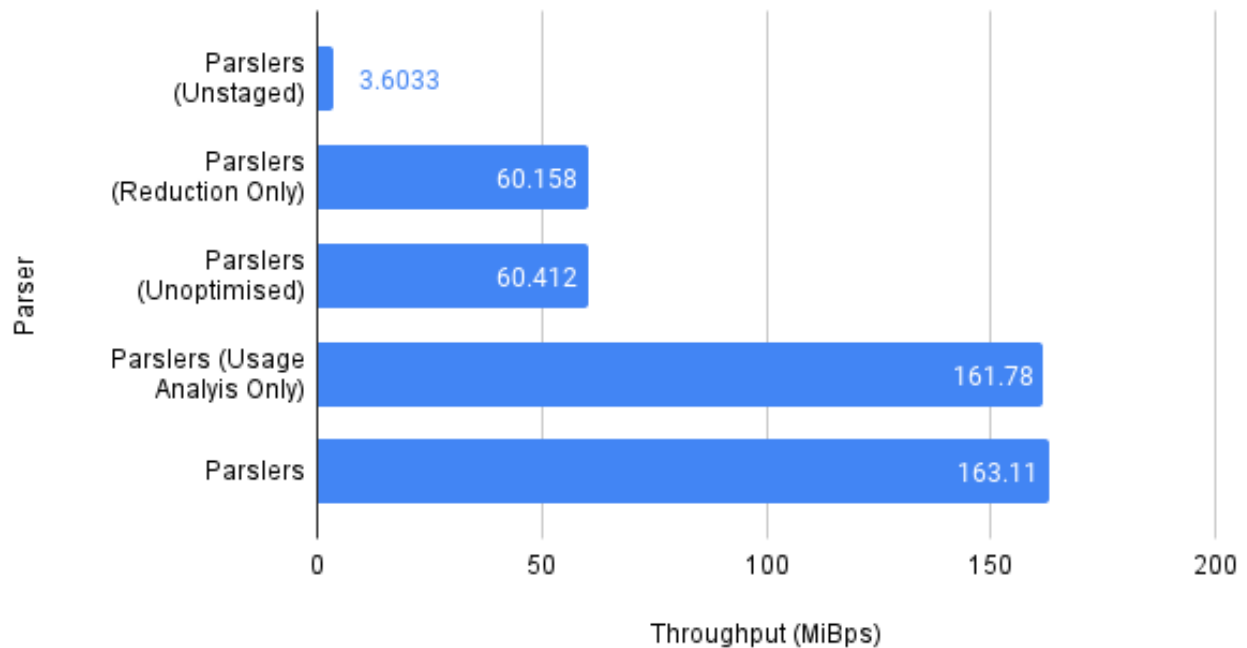


Figure 3. A plot of the throughput of parsers built with Parslers, each with different optimisation configurations, benchmarked on the Canada JSON dataset.

There is also a difference in which optimisation provides the largest performance benefit between the Branflakes parser and the JSON parser. In the Branflakes parser (see 4), reduction analysis provides the greatest benefit in performance improvement, but usage analysis makes no significant change to performance. On the other hand, the JSON parser (see 3) benefits greatly from usage analysis, but almost not at all from reduction analysis. This can be explained by the parsers written for the benchmarks having very different properties. For example, the Branflakes parser does not do any wasteful computation that can be optimised by the usage analysis optimiser but does have a more verbose structure that can be optimised using reduction analysis. However, the JSON parser does not change much during reduction analysis but does make heavy use of the `Recognise` parser, which can be heavily optimised using usage analysis.

5.1.5 Conclusion

Parslers clearly provides performance gains in comparison to the fastest general-purpose parsing library, `nom`. This shows that staging, and the opportunity for optimisations that staging provides, can help considerably close the gap between hand-written parsers and high-level general-purpose parsing libraries. Although Parslers is still not as fast as hand-written parsers (and may never be in the general case), there are still several opportunities for further optimisations available. These optimisations will be expanded upon in the Future Work section.

5.2 Expressive Power

Moving on to the discussion around the expressiveness of Parslers as a parser combinator library. As stated in 4.1, expressiveness in the context of Parslers means the range of grammars that can be expressed. As a part of the requirements, Parslers must be able to express any grammar expressed as a PEG grammar.

In order to demonstrate that Parslers is at least as expressive as PEG grammars, each of the avail-

Throughput of Parslers with Different Optimisations (Branflakes)

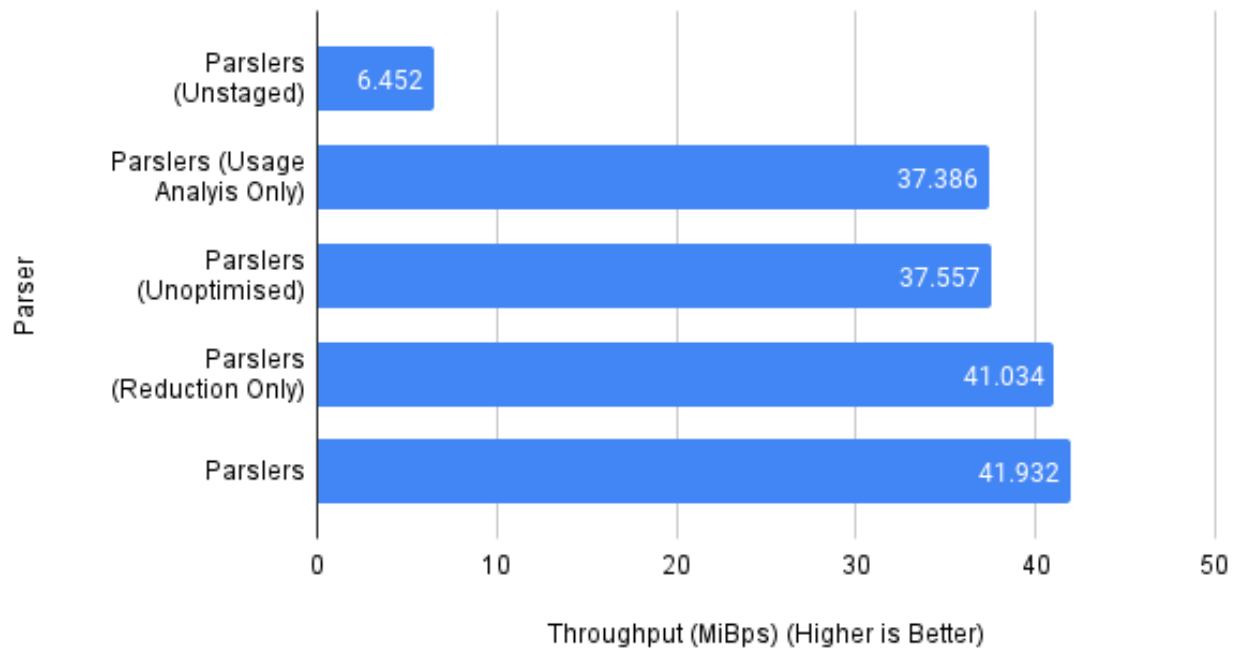


Figure 4. A plot of the throughput of parsers built with Parslers, each with different optimisation configurations, benchmarked on the Lost Kingdom Branflakes dataset.

able operators in PEG will be discussed, and a corresponding combinator will be given in Parslers, demonstrating the equivalence of the two.

It is worth noting that in PEG, failure always results in backtracking the whole consumed input. In Parslers, this is not implicit in order to improve performance. Backtracking upon failure must be explicitly added via the `Attempt` combinator.

Recall that PEG grammars follow the following operators are available:

5.2.1 Sequence

In PEG, the sequence operator takes two parsing expressions e_1 and e_2 , and performs them in sequence in the form e_1e_2 . That is, e_2 is parsed if e_1 succeeds, and either e_1 or e_2 failing also results in the whole sequence failing.

In Parslers, sequencing can be replicated using either the `Before` or `Then` combinators. In this form, `e1.attempt().then(e2.attempt())` or `e1.attempt().before(e2.attempt())` are semantically identical to e_1e_2 .

5.2.2 Ordered Choice

In PEG, ordered choice is in the form e_1/e_2 . This operator attempts to parse e_1 , but if e_1 fails then the consumed input is backtracked, and then e_2 is parsed.

In Parslers, backtracking is explicit except in the case of `Satisfy`, which backtracks upon failure. Although `Or` does fit, there is the caveat that it does not automatically backtrack upon the failure of the first parser. However, Parslers does have the `Attempt` combinator. Therefore, `e1.attempt().or(e2.attempt())` is equivalent to e_1/e_2 .

5.2.3 Zero/One-or-More

In PEG, zero-or-more and one-or-more are written in the form e^* and e^+ respectively. They attempt to parse several matches of the parser e until failure is reached, never backtracking.

In Parslers, the section 4.7 covers these semantics and extensions made to Parslers to optimise these sorts of parsers. `Loop(pure(()), e.attempt())` and `Loop(e.attempt(), e.attempt())` are equivalent to zero-or-more and one-or-more operators respectively.

5.2.4 Optional

In PEG, the optional operator is written in the form $e^?$, and will attempt to parse either zero or one occurrence of e in the input stream.

In Parslers, this is achievable by writing a combinator in the form `e.attempt().then(pure(()))`. However, Parsers comes built-in with the `opt(e)` auxiliary combinator that is less verbose.

5.2.5 And/Not Predicates

The and predicate is an operator written $\&e$. The semantics are to parse e , but always to backtrack upon success or failure. This is look-ahead semantics. Similarly, the not predicate is written $!e$, and fails if parser e succeeds, never consuming input.

In Parslers, both of these are achievable via the primitive combinators `Look` and `NegLook`, and can be written as `e.look()` and `e.neg_look()` to represent and not operators respectively. These maintain the backtracking semantics found in both of these operators.

5.2.6 Conclusion

The above arguments show that an isomorphism exists between PEG and Parslers, as any PEG grammar can be converted into a Parslers parser by the steps shown above. This, therefore, implies that Parslers must be at least as expressive as PEG. However, the context sensitivity introduced by the selective combinator `Branch` extends the expressive power of Parslers beyond PEG. This requirement is satisfied, and the Future Work section will discuss how the expressive power of Parslers can be improved upon.

5.3 Correctness

The correctness of Parslers is also defined in 4.1, and requires that Parslers maintains the type safety of user-written parser before and after staging. It also requires that Parslers apply restrictions on the types of parsers to be malformed before staging - in this way, disallowing the user to write malformed parsers.

This section does not concern itself with the correctness of the compiler, as this would require formal verification that Rust does not currently support. However, given that the optimisations performed are lawful analysis or usage analysis, these do not change the type of the resulting parser or its semantics.

There is a caveat to the above statement. Part of usage analysis is removing unnecessary computation. To demonstrate how this may be a problem, take a look at the following parser:

```
1 use std::sync::atomic::*;
2 static atomic: AtomicUsize = AtomicUsize::new(0);
3
4 #[reflect]
5 fn to_u32(c: char) -> u32 {
```

```

6     atomic.fetch_add(1, Ordering::SeqCst);
7     c as u32
8 }
9
10 fn parser() -> impl Parsler<Output = u32> + Clone {
11     match_char('a').map(to_u32)
12 }

```

As can be seen, the function `to_u32` mutates a global variable during execution. Under normal circumstances, this function will continue to be executed as the user expects. However, usage analysis on a similar parser, `parser().then(pure(()))` will remove the application of `to_u32` to the parsed character. Therefore, this function will never run.

Because of this, Parslers makes it a requirement that functions used in the creation of the Parser are pure and side-effect free. However, in Rust it is not yet possible to restrict functions to be pure apart from using the `const` keyword. Although using `const` is a viable solution, Rust considers most forms of allocation as non-const, and requires that functions that the user may call within a `const` function are also labelled `const`. This would become a major disadvantage to the usability of Parslers, and so was omitted. Therefore, it is possible that if functions produce side-effects, then the optimisations within Parslers may alter the behaviour of the program.

Apart from the caveat mentioned above, Parslers still requires that users write well-defined parsers at compile-time due to the type restrictions introduced by the `Parsler` trait (and its implementations over the primitive combinators). Therefore, Parslers can be considered correct in this perspective.

5.4 Usability

Parslers is a parser combinator library designed to be performant and easy to use. The core idea is to remove the burden of optimising parsers from the users while maintaining the type-safe Rust frontend.

There are a number of points that can be discussed here, which will be split up into subsections and discussed individually.

5.4.1 Staging

Staging is the heart of how Parslers can perform optimisations and native code generation at compile time. Without staging, Parslers could not be as performant as it is. However, staging does not come without issues of usability.

Firstly, staging requires some form of runtime reflection. This is because Parslers generates native Rust code, so pure values and referenced functions cannot make the jump from the `build.rs` file to the generated module without some way to serialize them and recreate them using Rust code. Although the `Reflect` trait, along with its macros make this possible, reflection does not yet exist for closures, as they cannot be inspected at compile time or at runtime. This means that Parslers can only accept a subset of the Rust programming language.

Furthermore, for larger parsers, the compile times can really bloat. One of the major advantages to the staging of parsers is that it restricts the number of code generation units the compiler can use,

giving more opportunity to potential optimisations being applied. However, this is at the cost of compile times - an issue the Rust compiler already faces. With sufficiently large parsers, this can cause a lot of problems for language servers used in code editors' incremental compile checks.

Finally, structures used within parsers must be declared in a separate crate. This is because the `build.rs` file is unable to access constructs defined in the crate's binary. Likewise, it does not make sense to define these structures directly in the `build.rs` file and use reflection to include them in the compiled Rust code. This might be an issue for users who do not want to use workspaces or split up their programs into separate crates. However, recent shifts within the Rust community suggest that for larger projects, it is idiomatic to split programs up into separate crates to make use of faster compile times and separation of concerns. Therefore, this requirement may not be an issue for users as the Rust community moves closer towards a modularised program architecture.

5.4.2 The Reflect Trait

The development of the `Reflect` trait has largely allowed this project to be possible. It allows functions and values to jump the staging into the generated Rust code. However, its introduction has come with disadvantages as well as advantages to the usability of the library. One of the biggest issues, which to my knowledge, cannot be fixed while using the `Reflect` trait. The 6.1 section discusses how this may be fixed.

Take, for example, the following function:

```
1  #[reflect]
2  fn zip<A, B>(a: A) -> impl FnOnce(B) -> (A, B) {
3      move |b| (a, b)
4  }
```

This is a function which takes a value and will return a function that converts the second value into a tuple of the two values. This is allowed without the `##[reflect]` macro, but not with it.

The compiler will throw the following error:

```
1  error[E0207]: the type parameter 'B' is not constrained by the impl trait,
2                self type, or predicates
3      --> parsers-test/build.rs:282:15
4      |
5  282 | fn zip<A, B>(a: A) -> impl FnOnce(B) -> (A, B) {
6      |                ^ unconstrained type parameter
```

In order to figure out why the compiler thinks this, the code generated by `##[reflect]` must be inspected:

```
1  #[allow(non_camel_case_types)]
2  #[allow(incorrect_ident_case)]
3  #[derive(Clone, Copy, Debug, Hash)]
4  pub struct zip;
5
```

```

6  impl<A, B> FnOnce<(A,)> for zip {
7      type Output = impl FnOnce(B) -> (A, B);
8      extern "rust-call" fn call_once(self, (a,): (A,)) -> Self::Output {
9          move |b| (a, b)
10     }
11 }

```

The issue is in the line `impl<A, B> FnOnce<(A,)> for zip // ...`, where we declare that the trait is generic over `A` and `B`, but do not constrain what `B` is bound by. This is a limitation of the `FnOnce` trait, which has not been kept up to date with the addition of generic associated types into the language. A workaround could be implemented by adding `B` as a generic to the `zip` struct, and rewriting the generated code to:

```

1  #[allow(non_camel_case_types)]
2  #[allow(incorrect_ident_case)]
3  #[derive(Clone, Copy, Debug, Hash, Default)]
4  pub struct zip<B>(std::marker::PhantomData<B>);
5
6  impl<A, B> FnOnce<(A,)> for zip<B> {
7      type Output = impl FnOnce(B) -> (A, B);
8      extern "rust-call" fn call_once(self, (a,): (A,)) -> Self::Output {
9          move |b| (a, b)
10     }
11 }

```

Which would work in theory. However, in practice, this would ruin the usage of `zip` as being identical to a function. A user would instead have to write something akin to `zip::<_>::default()` in order to use `zip` in their parser, which is not ergonomic. The 6.1 section will describe how this can be overcome with later additions of Parslers.

Finally, `Reflect` poses another issue. All structs used within Parslers parsers must implement the `Reflect` trait. This is an issue only because Rust does not allow foreign traits to be implemented on foreign types (those are, traits/types not declared in the crate). This is to disallow the same trait from being implemented on the same type in two separate crates used as dependencies in a main binary. There are workarounds to this, but this is less an issue with Parslers and more of an issue with the general Rust ecosystem. It is still worth noting that this can be considered a usability issue.

5.4.3 Conclusion

Although the review of Parslers' usability flaws has been critical, these issues that have been pointed out do not make Parslers a difficult library to work with. The code that must be written is still very close to pure Rust code, and the separation between the `build.rs` file and the generated Rust code should not deter many from writing code using Parslers. These are simply flaws that may be improved upon in the future while still remaining small enough that they should not be a deal breaker to those wishing to use Parslers within their own projects.

6 CONCLUSIONS

During this project, I have shown that it is possible to use the Rust programming language to implement a staged selective parser combinator library as a domain-specific language. Furthermore, I have shown that the use of staging in this way can produce parsers that are faster than the state-of-the-art libraries currently provided in the Rust ecosystem in the benchmarking process.

Although imperfect, I believe Parslers helps paves the way towards using staging to transform domain-specific languages written in Rust into native Rust code, while performing domain-specific optimisations on user-generated structures all at compile time.

During development, Parslers started off as a Rust macro DSL that used a custom syntax to describe parsers. However, I wanted to allow users to describe their parsers strictly as Rust code, using Rust's type system to prevent malformed parsers from being written accidentally. So Parslers started using Rust's staged compilation model to have parsers written as typed Rust structures, and compile them into native Rust code to be included in the final binary. This, however, needed a way to include functions and values declared in the `build.rs` file to be transferred over to the compiled Rust code, leading the way for something new in Rust: runtime reflection. The `Reflect` trait and macros became an entirely separate library which Parslers was built on top of.

With further work, Parslers could become a part of the Rust ecosystem. Furthermore, the `Reflect` package could help power more staged domain-specific languages in the Rust ecosystem, assisting in compile-time optimisations of user-written code.

6.1 Future Work

During subsequent versions of Parslers, there are several things I would improve upon.

6.1.1 The `Reflect` Trait

Although the `Reflect` trait was the enabling feat that allowed Parslers to make the jump between the different stages of compilation, it is not perfect. This has been discussed in the evaluation section previously.

One of the key issues with `Reflect` is that it still lacks semantic information. This could, in theory, be fixed by running `Rust doc`, a tool that automatically generates documentation from Rust code. This documentation can be outputted in JSON format and subsequently analysed in order to retrieve the declarations of functions and values used in Parslers at compile-time. The reason why I elected not to use this tool for Parslers during the project is that `Cargo` (the Rust build tool) creates a file lock when running commands that affect the target directory (such as building, checking, or documenting). This is to prevent race conditions in the target directory. However, towards the end of this project, another project by the name of `cargo px` was launched, which aliases `cargo` for the purposes of executing `cargo` commands during compile-time. If `cargo px` was launched prior to the start of this project rather than towards the end, this may have been a better route to go down rather than implementing runtime reflection for the purposes of Parslers.

6.1.2 Performance Optimisations

Parslers currently stands as the fastest general-purpose parsing framework for the Rust programming language. However, it is still far away from approaching the performance of hand-written parsers. Through implementing further optimisation passes, I believe the speed of hand-written parsers can be approached in some cases.

One of the large potential optimisations that could be explored as a part of the next iteration would be consumption analysis. Consumption analysis works by factoring out length checks in the parser

where possible. Take the following parser for example:

```
1 match_char('a').then(match_char('b')).then(match_char('c')).then(match_char('d'))
```

The generated code for this parser looks like the following:

```
1 pub fn parser(input: &mut std::str::Chars) -> Result<char, &'static str> {
2     {
3         {
4             let old_input = input.clone();
5             input
6                 .next() // Check for next character here
7                 .ok_or_else(|| "Found EOF when character was expected")
8                 .and_then(|c| {
9                     if (f0)(c) {
10                        Ok(())
11                    } else {
12                        *input = old_input;
13                        Err("expected a specific character")
14                    }
15                })
16        }
17        .and_then(|_| {
18            let old_input = input.clone();
19            input
20                .next() // Check for next character here
21                .ok_or_else(|| "Found EOF when character was expected")
22                .and_then(|c| {
23                    if (f1)(c) {
24                        Ok(())
25                    } else {
26                        *input = old_input;
27                        Err("expected a specific character")
28                    }
29                })
30        })
31    }
32    .and_then(|_| {
33        let old_input = input.clone();
34        input
35            .next() // Check for next character here
36            .ok_or("Found EOF when character was expected")
37            .and_then(|c| {
38                if (f2)(c) {
39                    Ok(c)
40                } else {
41                    *input = old_input;
42                    Err("Expected a specific character")
43                }
44            })
45    })
46 }
```

```
44         })  
45     })  
46 }
```

The exact semantics of the generated code is not of much importance for this example. However, notice that `input.next()` is called three times (corresponding to the three characters matched during the parser with `match_char`). Each occurrence of `input.next()` makes a check to see if the end of the input stream has been reached. However, it is worth noting that the parser will always fail if the input stream has less than three characters. Therefore, it stands to reason that these length checks may be factored out of the parser at compile time via an optimisation pass on the parser. This exact analysis is detailed further in (Willis, Wu & Schrijvers 2022).

Parsers for languages such as Rust, C, C++, Java, and many others will have a large number of checks for keywords whose lengths are known at compile-time to the parser. This optimisation could provide performance benefits to those parsers.

Furthermore, parsers such as `match_char('a').or(match_char('b'))`, where a character is read on one branch of an `Or` combinator and the other, may also benefit from having only one length check at the beginning of the `Or` sub-parser execution.

6.1.3 Error Messages

For many programming language compiler frontends, parsers must produce detailed error messages in the case a syntax error is encountered. In the future, I hope to add this detailed error messaging into Parslers automatically. Using staging, reasonable error messages could be inferred and embedded into the parser at compile-time. This would provide great utility to users of Parslers, as this could potentially remove the need to format error messages themselves, as is required in several other parsing libraries. The notable exception to this is Chumsky (Barretto 2021), which produces very readable error messages upon failure. However, Chumsky's performance characteristics are rather poor, which does not make a good use case for high-performance parsing.

6.1.4 Asynchronous Parsing

Another avenue that would be exciting to explore is the idea of asynchronous parsing. That is, to parse an input stream as it is being received from an asynchronous executor. This is possible thanks to Rust's strong `async/await` functionality being implemented using `Future`s, which allows functions to be paused and continued as asynchronous data is required/given.

An interesting use case is for low-latency applications such as market making, where the faster you are able to parse, analyse and make a decision on data, the better you can perform in the markets. Being able to write fast custom parsers that can process data as it is being received opens up the avenue for reducing latency in data analysis pipelines.

Another example of how this may be useful is parsing data from a disk using buffered asynchronous readers, whereby a double buffer is used to parse data and fetch data at the same time. This would remove the need to store the whole data in memory, but at no cost of complexity to the user.

For grammars with finite look-ahead, the asynchronous parsers would require a ring buffer of size equal to the look-ahead in order to backtrack successfully upon failure. Using consumption analysis, the parser could be analysed to figure out the buffer size required at compile time, without any input from the user.

6.2 Example of the Branflakes Parser

This section will include an annotated example of a Branflakes parser built with Parslers. Branflakes is chosen because the grammar is small and simple, but complex enough that it can demonstrate how a user would interact with Parslers in their own projects.

Firstly, in a separate crate, define the structures used in the results of the parsers. Use the `Reflected` macro to automatically derive the `Reflect` trait on these structures.

```
1 use parsers_lib::reflect::Reflect;
2 use parsers_macro::*;
3
4 #[derive(Clone, Debug, Reflected)]
5 pub struct BranflakesProgram(pub Vec<Branflakes>);
6
7 #[derive(Clone, Debug, Reflected)]
8 pub enum Branflakes {
9     Add,
10    Sub,
11    Left,
12    Right,
13    Read,
14    Print,
15    Loop(BranflakesProgram),
16 }
```

Then, construct the parser, and use the `Builder` struct provided to compile the parsers together, apply optimisations, and compile them to a `branflakes.rs` file in the target build directory.

```
1 use parsers_branflakes::Branflakes;
2 use parsers_lib::{
3     builder::Builder,
4     parser::{auxiliary::*, *},
5     reflect::*,
6 };
7 use parsers_macro::reflect;
8
9 #[reflect]
10 fn branflakes_val(
11     p: Vec<parsers_branflakes::Branflakes>,
12 ) -> parsers_branflakes::BranflakesProgram {
13     parsers_branflakes::BranflakesProgram(p)
14 }
15
16 #[reflect]
17 fn branflakes_loop(p: parsers_branflakes::BranflakesProgram)
18     -> parsers_branflakes::Branflakes {
19     parsers_branflakes::Branflakes::Loop(p)
20 }
```



```

21 fn branflakes_program()
22   -> impl Parsler<Output = parsers_branflakes::BranflakesProgram> + Clone {
23     let left = match_char('<').then(pure(Branflakes::Left));
24     let right = match_char('>').then(pure(Branflakes::Right));
25     let add = match_char('+').then(pure(Branflakes::Add));
26     let sub = match_char('-').then(pure(Branflakes::Sub));
27     let print = match_char('.').then(pure(Branflakes::Print));
28     let read = match_char(',').then(pure(Branflakes::Read));
29     let loop_ = || {
30       match_char('[')
31         .then(branflakes_program())
32         .map(branflakes_loop)
33         .before(match_char(']'))
34     };
35     name("branflakes", move || {
36       many(
37         left.or(right)
38           .or(add)
39           .or(sub)
40           .or(print)
41           .or(read)
42           .or(loop_()),
43       )
44       .map(branflakes_val)
45     })
46   }
47
48 fn main() {
49   let out_dir = std::env::var("OUT_DIR").unwrap();
50   Builder::new("branflakes", branflakes_program())
51     .add_parser("branflakes_validate", branflakes_program().then(pure(())))
52     .reduce()
53     .usage_analysis()
54     .build(&format!("{out_dir}/branflakes.rs"));
55 }

```

Finally, use the `include!()` macro to add the generated file to the binary, and call the named parser directly.

```

1 mod branflakes {
2   #![allow(warnings, unused)]
3   include!(concat!(env!("OUT_DIR"), "/branflakes.rs"));
4 }
5
6 fn main() {
7   let data = std::fs::read_to_string("LostKng.b").unwrap();
8
9   let chars = &mut data.chars();

```

```
10  let parsed = branflakes::branflakes(chars);
11  println!("Output: {:?}", parsed.is_ok());
12  println!("Remainder: {:?}", chars.as_str());
13
14  let chars = &mut data.chars();
15  let validate = branflakes::branflakes_validate(chars);
16  println!("Validate: {:?}", validate.is_ok());
17  println!("Remainder: {:?}", chars.as_str());
18 }
```

And with about 90 total lines of code, a ready-to-use Branflakes parser and validator are ready to be used.

REFERENCES

- Allott, Nicholas; Lohndal, Terje & Rey, Georges. 2021, Synoptic Introduction, *A Companion to Chomsky*, pp. 1–17.
- Anderson, Carter. 2020, *Bevy: A refreshingly simple data-driven game engine built in Rust*. Available: <https://bevyengine.org/>.
- Barretto, Joshua. 2021, *Chumsky: A parser library for humans with powerful error recovery*. Available: <https://crates.io/crates/chumsky>.
- Griesemer, Russ, Robert. Cox & Fitzpatrick, Brad. 2015, *cmd/compile/internal/gc: recursive-descent parser*. Available: <https://go-review.googlesource.com/c/go/+16665>.
- Hutton, Graham & Meijer, Erik. 1996, *Monadic parser combinators*.
- Jäger, Gerhard & Rogers, James. 2012, Formal language theory: refining the Chomsky hierarchy, *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 367, no. 1598, pp. 1956–1970.
- Muchnick, Steven et al.. 1997, *Advanced compiler design implementation*, Morgan kaufmann.
- Pezoa, Felipe; Reutter, Juan L; Suarez, Fernando; Ugarte, Martín & Vrgoč, Domagoj. 2016, Foundations of JSON schema, In: *Proceedings of the 25th International Conference on World Wide Web*, International World Wide Web Conferences Steering Committee, pp. 263–273.
- Willis, Jamie; Wu, Nicolas & Pickering, Matthew. 2020, Staged selective parser combinators, *Proc. ACM Program. Lang.*, vol. 4, no. ICFP, pp. 120:1–120:30. Available: <https://doi.org/10.1145/3409002>.
- Willis, Jamie; Wu, Nicolas & Schrijvers, Tom. 2022, Oregano: staging regular expressions with Moore Cayley fusion, In: Nadia Polikarpova, ed., *Haskell '22: 15th ACM SIGPLAN International Haskell Symposium, Ljubljana, Slovenia, September 15 - 16, 2022*, ACM, pp. 66–80. Available: <https://doi.org/10.1145/3546189.3549916>.