

CS4201 Practical 1

160001055

October 2019

1 Introduction

This practical involved the implementation of a lexer and a parser for the programming language *Oreo*, given an informal specification of the language and various examples. Furthermore, as an extension, I have implemented functions and comments handling, extensive error handling and a primitive type checker.

2 Design

The design for a parser is always a challenging one, as factors such verifiability, correct error handling and similar often conflict. As such, I decided to work on a multi-layer design stage. The lowest level, the scanner, simply strips whitespace, and isolate strings and comments. On top of it the lexer recognizes the tokens and passes them to the parser, whose job to actually properly parse the grammar is simplified. Once the parse tree is constructed, the validator runs on it, ensuring that each rule is well formed. Each stage is independently tested, to ensure correctness, and to make sure that changes at each level will not effect the others.

2.1 Iteration

The entire application is built as a on demand stream of characters. Starting from the original character streams, each layers asks the underlying one for the next element, and waits until the computation is completed for the result. This has the added flexibility that, when the input is of considerable size, it is not needed to allocate for the scanning and lexing phase, which could considerably impact memory usage. It also guarantees that the grammar to be parsed is **LL(1)**, as backtracking is impossible. In exchange, computationally this can be slightly less efficient, as operation are interleaved hurting branch prediction. However, since in a full compiler the lexing and the parsing are almost never the bottleneck, I felt the benefits outweighed the cons.

2.2 Scanning

The scanning operation essentially isolates strings and comments in the code, Other than this it strips white spaces and chunks together characters that are not separated by white-space. Further than this, it also collects the line numbers of the scanned item, to be passed to the next stages for error reporting. The scanning is designed to be strictly $O(n)$ in the size of the input, never actually ever needing to backtrack. Also it is nice to note that comments are actually propagated to next stages, and not simply removed. In theory this could later allow for comment-based syntax extensions¹.

2.3 Lexing

Lexing is substantially more complex than scanning, as it cannot rely on white-space for token separation, and as such needs to have a more complex solution. The lexer we have queries the lexer, and operates on each resulting item in order to generate a corresponding stream of tokens. All of these streams of tokens are then flattened in a contiguous stream that then the parser operates on. As in scanning, the lexer is made so to only ever need to read each character in the input once, which, as the implementation shows, is harder than it looks.

2.4 Grammar

In order to allow for the parser to never need to backtrack, I needed to convert the grammar to **LL(1)**. This involved eliminating left recursion, specifying order of operations, removing ambiguity and computing start and follow sets. This was not particularly challenging for most of the grammar, except for the expression part. In particular, preserving the typing in case of bool expression proved complicated. In order to solve this, I decided to move the type checking to a later stage (as it is done in most compilers) and to work with a simpler fully unambiguous grammar. In particular the grammar we resulted with follows. Note that cursive terms are non terminal, bold ones are tokens and typed ones are tokens that might take multiple values (typically literals). Also, not that comments are not included here, as the parser filters them out by default.

$$\begin{aligned} P &\rightarrow \mathbf{program} \textit{id compound} \\ \textit{compound} &\rightarrow \mathbf{begin} \textit{statement statements end} \\ \textit{statements} &\rightarrow \textit{statement statements} | \epsilon \\ \textit{statement} &\rightarrow \textit{decl} | \textit{print} | \textit{assign_or_fun} | \textit{if} | \textit{while} | \textit{fun} \\ \textit{decl} &\rightarrow \mathbf{var} \textit{id opt_expr}; \\ \textit{opt_expr} &\rightarrow \textit{:= expr} | \epsilon \\ \textit{print} &\rightarrow \mathbf{print} \textit{expr}; | \mathbf{println} \textit{expr}; | \mathbf{get} \textit{id}; \end{aligned}$$

¹I am not advocating this here, it is almost certainly a terrible idea, but you *could* do it

$$\begin{aligned}
& assign_or_fun \rightarrow id\ tuple_or_assign; \\
& tuple_or_assign \rightarrow :=\ expr\ |\ (c_args) \\
& if \rightarrow \mathbf{if}\ (expr)\ \mathbf{then}\ compound\ \mathbf{else}; \\
& \quad else \rightarrow \mathbf{else}\ compound\ |\ \epsilon \\
& while \rightarrow \mathbf{while}\ (expr)\ compound; \\
& fun \rightarrow \mathbf{procedure}\ id\ (d_args) compound; \\
& \quad d_args \rightarrow \mathbf{var}\ id\ d_args'\ |\ \epsilon \\
& \quad d_args' \rightarrow ,\ \mathbf{var}\ id\ d_args'\ |\ \epsilon \\
& \quad expr \rightarrow term\ expr' \\
& \quad expr' \rightarrow \wedge\ term\ expr'\ |\ \epsilon \\
& \quad term \rightarrow factor\ term' \\
& \quad term' \rightarrow \leq_R\ factor\ term'\ |\ \epsilon \\
& \quad factor \rightarrow product\ factor' \\
& \quad factor' \rightarrow \oplus\ product\ factor'\ |\ \epsilon \\
& \quad product \rightarrow atom\ product' \\
& \quad product' \rightarrow \otimes\ atom\ product'\ |\ \epsilon \\
& \quad atom \rightarrow unit\ |\ \mathbf{not}\ atom \\
& unit \rightarrow (expr)\ |\ id\ func_args\ \mathbf{str}\ |\ \mathbf{int}\ |\ \mathbf{bool} \\
& \quad func_args \rightarrow (c_args)\ |\ \epsilon \\
& \quad c_args \rightarrow expr\ c_args'\ |\ \epsilon \\
& \quad c_args' \rightarrow ,\ expr\ c_args'\ |\ \epsilon \\
& \quad \wedge \rightarrow \mathbf{and}\ |\ \mathbf{or} \\
& \quad \leq_R \rightarrow < \ |\ \leq \ |\ > \ |\ \geq \ |\ == \\
& \quad \oplus \rightarrow + \ |\ - \\
& \quad \otimes \rightarrow * \ |\ /
\end{aligned}$$

This grammar, despite being really verbose, does do operator precedence correctly, grouping operators of the same level of precedence on the right. The operator grouping is as follows: $\wedge, \leq_R, \oplus, \otimes, \mathbf{not}$ so the following expression: " $\mathbf{not} 1 * 2 + 3 < 4 \mathbf{and} 5$ " should parse as $((\mathbf{not}(1 * 2) + 3) < 4) \mathbf{and} 5$ and the following " $a - b - c$ " will parse as $(a - (b - c))$.

2.5 Parsing

The parser is a simple recursive descent one, which progressively builds a typed abstract syntax tree. It is mostly a direct translation of the grammar above, operating on the stream of tokens produced by the lexer. Each nonterminal f corresponds to a function of the form $f : \text{TokenStream} \rightarrow \text{Ok } F | \text{Err } \text{Error}$ where F is a typed node that represent the terminal. Composing this functions quite naturally allows for both good error reporting and flexible parsing.

2.6 Validation

In the original grammar, there were different production rules for boolean and expressions, such that every boolean was an expression and not vice-versa. Despite much trying, I wasn't able to construct an equivalent grammar that was **LL(1)** and still was able to uphold this condition. In fact I speculate that such a grammar does not exist, but this is just speculation. In order to still produce errors when invalid expressions are used as booleans, I added an extra validation step, which in a way mimics the way a type-checker in a complete compiler works. The validator walks the abstract syntax trees and checks the following conditions holds:

$$\begin{aligned} \text{if}(e) \dots &\implies \mathbf{B}(e) \\ \text{while } (e) \dots &\implies \mathbf{B}(e) \\ \text{expr} &\implies \mathbf{V}(\text{expr}) \end{aligned}$$

Where we define:

$$\text{recV}(f, \text{tail}) = \begin{cases} \text{tail} = \epsilon & \mathbf{true} \\ \text{otherwise} & f(\text{tail}) \end{cases}$$

And then we recursively (using pattern matching) define \mathbf{V}, \mathbf{B} as follows:

$$\mathbf{V}(\text{expr}) \text{ where } \text{expr} \rightarrow \text{term expr}' = \begin{cases} \text{expr}' = \epsilon & \mathbf{V}(\text{term}) \\ \text{otherwise} & \mathbf{B}(\text{term}) \text{ and } \mathbf{B}(\text{expr}') \end{cases}$$

$$\mathbf{B}(\text{expr}') \text{ where } \text{expr}' \rightarrow \wedge \text{term expr}'_1 = \mathbf{B}(\text{term}) \text{ and } \text{rec}(\mathbf{B}, \text{expr}'_1)$$

$$\mathbf{V}(\text{term}) \text{ where } \text{term} \rightarrow \text{factor term}' = \mathbf{V}(\text{factor}) \text{ and } \text{rec}(\mathbf{V}, \text{term}')$$

$$\mathbf{B}(\text{term}) \text{ where } \text{term} \rightarrow \text{factor term}' = \begin{cases} \text{term}' = \epsilon & \mathbf{B}(\text{factor}) \\ \text{otherwise} & \mathbf{V}(\text{factor}) \text{ and } \mathbf{V}(\text{term}') \end{cases}$$

$\mathbf{V}(term')$ where $term' \rightarrow_{\leq R} factor$ $term'_1 = \mathbf{V}(factor)$ and $rec(\mathbf{V}, term'_1)$

$\mathbf{V}(factor)$ where $factor \rightarrow product$ $factor' = \mathbf{V}(product)$ and $rec(\mathbf{V}, factor')$

$\mathbf{B}(factor)$ where $factor \rightarrow product$ $factor' = \mathbf{B}(product)$ and $factor' = \epsilon$

$\mathbf{V}(factor')$ where $factor' \rightarrow \oplus product$ $factor'_1 = \mathbf{V}(product)$ and $rec(\mathbf{V}, factor'_1)$

$\mathbf{V}(product)$ where $product \rightarrow atom$ $product' = \mathbf{V}(atom)$ and $rec(\mathbf{V}, product')$

$\mathbf{B}(product)$ where $product \rightarrow atom$ $product' = \mathbf{B}(atom)$ and $product' = \epsilon$

$\mathbf{V}(product')$ where $product' \rightarrow \otimes atom$ $product'_1 = \mathbf{V}(atom)$ and $rec(\mathbf{V}, product'_1)$

$\mathbf{V}(atom)$ where $atom \rightarrow unit = \mathbf{V}(unit)$

$\mathbf{V}(atom)$ where $atom \rightarrow \mathbf{not}$ $atom_1 = \mathbf{B}(atom_1)$

$\mathbf{B}(atom)$ where $atom \rightarrow unit = \mathbf{B}(unit)$

$\mathbf{B}(atom)$ where $atom \rightarrow \mathbf{not}$ $atom_1 = \mathbf{B}(atom_1)$

$\mathbf{V}(unit)$ where $unit \rightarrow \text{id } func_args \mid \text{str} \mid \text{int} \mid \text{bool} = \mathbf{true}$

$\mathbf{V}(unit)$ where $unit \rightarrow (expr) = \mathbf{V}(expr)$

$\mathbf{B}(unit)$ where $unit \rightarrow (expr) = \mathbf{B}(expr)$

$\mathbf{B}(unit)$ where $unit \rightarrow \text{id } func_args \mid \text{bool} = \mathbf{true}$

$\mathbf{B}(unit)$ where $unit \rightarrow \text{str} \mid \text{int} = \mathbf{false}$

Using this simple (yet verbose) rules, we can ensure that every expression is valid, the operand of boolean operators are indeed booleans and that the conditions in control flows are of the right type.

2.7 Error Handling

Error handling is done independently in each stage of the parser. The scanner is able to report unclosed strings and comments. The lexer is able to report when some characters cannot be recognized as tokens (i.e. `=` when not in `==`, `<=`, `=>`, `:=`). Similarly the parser reports most of the error in the grammar and generally it is quite insightful. Up to this point, the error information contains the line number and a succinct description of what the compiler was expecting and what it was found. The further validation, which was added later after the grammar changes, is able to validate typing very effectively, but since the abstract syntax tree does not track line information it is not able to point out where the error is. In order to accommodate that the best solution would be a two layer parse tree, one strongly typed as the one we have now and one typeless one which would track line information. That would require a lot of changes in the parsing, which I did not have time to implement. While the scanner and the lexer can lazily build the stream and report every error at the end, the parser works eagerly reporting errors as soon as they come. While it would be possible to parse as much as possible before reporting errors, it would have involved a lot of change as:

- Introducing an **Error** node in the typed tree
- Allowing the **Error** to be a children of any other node

Which would have made any further operation on the syntax tree more cumbersome. Furthermore, this way integrated neatly with Rust's syntax and try operator, which significantly cut down the lines of code.

2.8 Testing

Since abstract syntax trees can grow to considerable sizes, especially when a grammar is so granularly defined as ours, and since parsing can present many pitfalls I wanted a testing solution that allowed for easily adding and verifying new cases, without having to manually validate each one. To this avail, I took a page out of the rust analyzer² project, and used a snapshot based solution. Each test generates a snapshot of the output of that test, and on each change the user gets a chance to review said snapshot to verify that the behaviour is as expected. In particular, adding each of the examples as a snapshot test allowed me to be sure I was compliant to the specification. Also, since the testing independently verified each layer of the parse, I could verify that changes in one section would not negatively effect the others. Also, without running the program, the snapshots are version controlled and serve as documentation for what the parsing will result in.

²<https://github.com/rust-analyzer/rust-analyzer>

2.9 Outputs

For the output requirements of the program, I decided to aim for a machine readable format. Since my artistic skills are not in a great place, this seemed like the best solution, and also this felt like a decent way to allow for other programs to use the output in a proper way. As for machine readable format, I allow the user to choose between YML, JSON and RON (Rusty Object Notation). While neither of them is particularly pretty, I personally quite like the YML output. Also, since the serialization is powered by the `serde`³ framework, it would be trivial to add any more of the supported data types.

3 Implementation

3.1 Rust

The choice of Rust for this project was dictated by various factors:

- I already had experience with Rust in parsers, thanks to some contributions I did on Rust analyzer
- The snapshot testing library works very well for ASTs
- A strong type system allows for the ASTs to be completely typechecked
- Performance
- Static memory safety
- Pattern Matching
- Try error handling

Some of these, like performance and testing, are simply nice to have, and don't effect the program too much. Others, like the strong type checking and the pattern matching, are simply invaluable, as they prevent you from shooting yourself in the foot too often. In particular, without the Try operator, the size of the parser might have doubled or tripled, and with it the complexity overhead.

3.2 Zero Copy

One of the most interesting features of this program is that I actively aimed to reduce allocations and copying to a minimum. In fact, the only allocations in the entire core are those that are needed to avoid infinitely recursive types, and for the statement list. Even the identifier token never allocates, but instead keeps a reference to the data in the original input string. As for copying, I can guarantee that for complex types it never happens unless it is in an error path.

³<https://serde.rs/#data-formats>

3.3 Iteration

The parser operates on a `TokenStream`, which is simply a trait wrapper over the `Peekable<Iterator<Item = Token>>` trait. What does this imply, together with the zero allocation condition, is that the grammar must be **LL(1)** and that the parser never backtrack. In fact, this condition can be statically checked by the compiler, which I think is awesome. The above follows from the fact that we can at most only peek one token ahead in the stream, and from the fact that a general iterator cannot be rewinded. Furthermore, the no allocation property implies we cannot store a dynamic amount of tokens to further use as rewind stack, and as such we must be satisfying the desired properties.

3.4 Benchmarking

While it was not required for the practical, I found it interesting to write a couple of benchmarks for the parser. I have three benches for parsing, one which test a multitude of linear assignments, one which tests some very deeply nested ifs and one which tests a single extremely complex expression. Further, I have some benchmarks that test performance for lexing and scanning. My findings are as follows, first of all I can comfortably handle 32 level of nesting before a stack overflow, which is nice⁴. Secondly, scanning and lexing account for remarkably big slice of the time the program spends on parsing. Third, throughput seems to decrease with the increase of input size, which seems to be caused by the decreased data locality (at least so I speculate). The result are as follows (throughput in MiB/s).

Statements	Scanner	Lexer	Parsing
128	4.3	28.7	13.2
256	2.2	29.0	13.0
512	1.1	27.8	11.8
1024	0.587	28.6	7.3
2048	0.288	28.3	6.4

Table 1: Throughput on Linear Program

Nesting	If	Expr
2	37.8	27.9
4	33.8	21.2
8	31.1	19.0
16	33.3	18.0
32	30.8	17.8

Table 2: Throughput on Non Linear Program

⁴Using an explicit stack would allow for an almost limitless number of levels, but anyone using more than 32 nested statement is a madman already

A couple of notes, first of all the result of the benchmark are a bit unfair, as the scanner is the only one which effectively iterates trough each of the input bytes over which the throughput is calculated. The lexer already (for example) skips string, comments and whitespace. Finally the parser operates only on prepackaged tokens, which speeds it up quite a bit. Also, I suggest to take the benchmarks in the nested cases with a grain of salt, as since the input size is smaller they are more prone to performance changes. In order to isolate the results of each stage, we collect the result of the previous stage, so not to have the on demand computation we used to have before.

Also, just to have some pictures, here are some flamegraphs: As you can

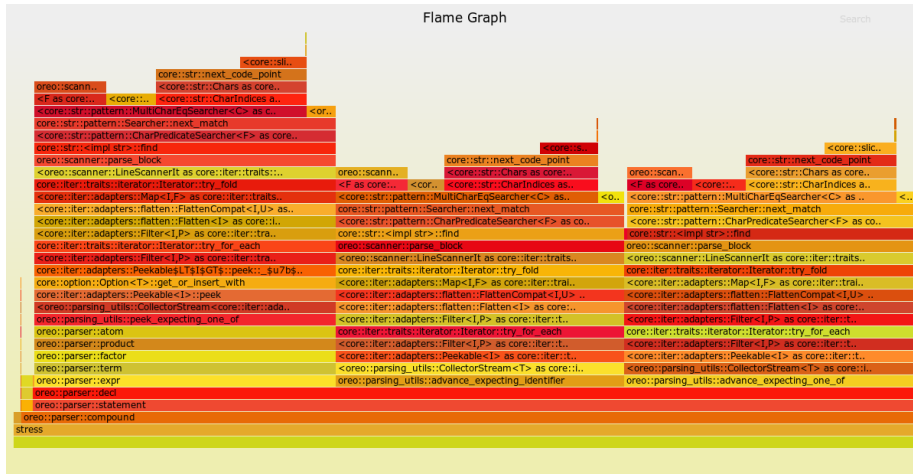


Figure 1: Flamegraph on demand parsing

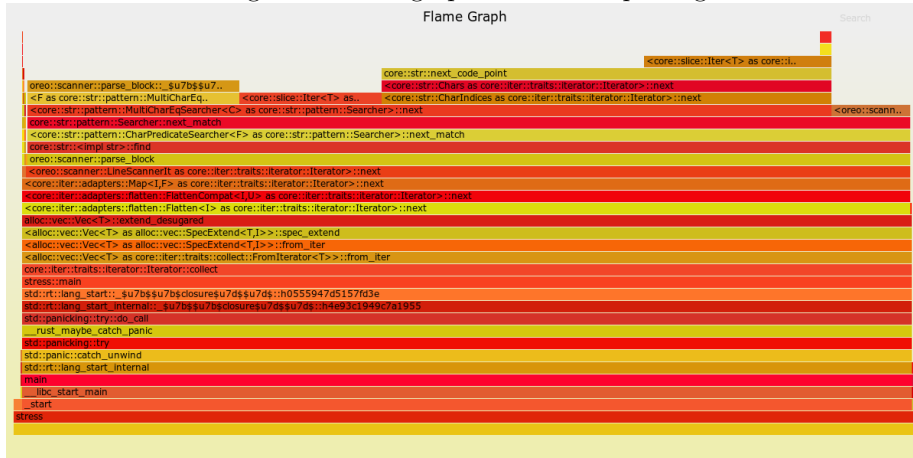


Figure 2: Flamegraph parsing after collecting tokens

see, in the on demand model each section takes approximately the same time,

put the actual impact of parsing is much more evident in the model where we precompute everything ahead of parsing. As you can see, (or actually as you can't see), the parsing in the second graph accounts for approximately 0.5% of the running time, and as such is definitely not the bottleneck. Most of the time is spent looking for characters in strings, which is what happens in both the scanner and the lexer.

4 Conclusion

I thoroughly enjoyed this practical, as it finally allowed me to practically develop a parser for a programming language which, despite not being terribly complex, had many real life problems to be addressed. The development of a grammar from the specification presented many challenges, and finding a workaround for some of the issues that it had was also a stimulating exercise. Having an occasion to use a solid and efficient testing framework also made this extremely nice. Finally, while performance was never a focus for this practical, I still found it worthwhile to analyze, as it showed the discrepancy between expectations and code.