

# Comparison of Rust Implementation to Original C4 Compiler

Meera Alawadhi

Hamda

May 1, 2025

## Introduction

The C4 compiler, originally implemented in C, is a lightweight, self-contained compiler for a subset of C. The `c4_rust_team_alfaqa` project reimplements C4 in Rust, focusing on lexing and parsing to produce an Abstract Syntax Tree (AST) for C-like code. This report compares the Rust implementation to the original C version, evaluating the impact of Rust's safety features, performance differences, and challenges in maintaining compatibility with C4's behavior.

## Rust's Safety Features and Design Impact

Rust's safety guarantees, including memory safety and strict lifetime management, significantly shaped the design of the Rust implementation while ensuring compatibility with C4's functionality.

### Memory Safety

In the original C version, manual memory management using `malloc` and `free` posed risks of memory leaks, dangling pointers, and buffer overflows, particularly in string handling and token buffering. Rust eliminates these risks through its ownership model and automatic memory management. For example, the lexer in `c4.rs` uses `String` and `Vec<char>` for input processing, ensuring memory is deallocated when variables go out of scope. This removed the need for manual memory cleanup, as seen in C's `free` calls, while maintaining identical tokenization of keywords, identifiers, and strings.

The parser's AST, defined with `enum Expr` and `enum Stmt`, leverages Rust's ownership to manage nested structures (e.g., `Box<Expr>` for recursive expressions) without manual pointer management. This contrasts with C's use of raw pointers, reducing the risk of null pointer dereferences. Compatibility with C4's AST output was preserved by mirroring its structure, ensuring identical parsing of statements like `if (x) { return 5; } else { return 0; }.`

### Lifetimes

Rust's lifetime system enforced strict borrowing rules, impacting the lexer and parser design. In `Lexer::new`, the input string is converted to a `Vec<char>` owned by the `Lexer` struct, avoiding lifetime annotations for borrowed slices. This simplified the design compared to C's `char*`

pointers, which required careful tracking to avoid invalid memory access. The parser's `current_token: Result<Token, String>` ensured tokens were safely borrowed during parsing, preventing data races or invalid references. These features maintained C4's single-pass parsing behavior while enhancing robustness.

## Performance Differences

Quantitative performance measurements were not conducted due to the project's focus on functionality up to parsing. However, qualitative differences can be inferred from Rust's and C's characteristics when compiling the same C code (e.g., `test.c` with `{ x = 5; printf("%s", p); return 0; }`).

Rust's zero-cost abstractions, such as `String` and `Vec`, introduce minimal overhead compared to C's raw arrays and pointers. The Rust lexer's use of `Vec<char>` for input iteration may incur slight allocation costs versus C's direct `char*` traversal, but Rust's iterator optimizations mitigate this. Parsing performance is comparable, as both implementations use single-pass, recursive-descent parsing. Rust's bounds checking on `Vec` accesses adds minor runtime overhead absent in C, but this ensures safety without significant slowdown for typical C4 inputs.

The Rust compiler's optimizations (e.g., `cargo build -release`) produce efficient binaries, likely within 10-20% of C's performance for parsing tasks, based on typical Rust benchmarks. For small inputs like `c4.c`, the difference is negligible. Future work, such as implementing the virtual machine (VM), would require benchmarking to quantify execution speed, but the current parser's performance is sufficient for C4's use case.

## Challenges and Solutions

Replicating C4's behavior in Rust presented several challenges, primarily due to differences in language paradigms and safety requirements.

### Challenge: Replicating C's Flexible Parsing

C4's C implementation uses lenient parsing, tolerating minor syntax errors (e.g., missing semicolons) through pointer-based token skipping. Rust's strict type system and error handling required explicit error management. For example, `parse_stmt` in `c4.rs` enforces semicolon termination for statements, returning `Err` for invalid syntax. To maintain compatibility, the parser was designed to produce identical ASTs for valid inputs, using `Result<Token, String>` to handle errors gracefully, mirroring C4's error reporting without compromising safety.

### Challenge: Handling C's Unsafe String Operations

C4's string handling (e.g., for string literals like `"hello n"`) relies on unsafe `char*` manipulation, risking buffer overflows. Rust's `String` type ensured safe string construction in `Lexer::next`, with escape sequence handling ( `n`,

") implemented using safe push operations. This preserved C4's string tokenization while eliminating overflow risks.

### **Challenge: Test Failures and Debugging**

Initial test failures in `test_parse_block`, `test_parse_error`, `test_parse_expr_stmt`, and `test_parse_if_else_stmt` stemmed from incorrect semicolon handling and `else` branch parsing. These were resolved by simplifying `parse_stmt` to parse expressions first, then check for `=` or `;`, and ensuring `parse_block` correctly handles braces. Rust's E0507 move error in `parse_block` was fixed by borrowing `&self.current_token`, aligning with Rust's ownership rules while maintaining C4's parsing logic.

### **Conclusion**

The Rust implementation of C4 successfully replicates the original's lexing and parsing behavior while leveraging Rust's safety features. Memory safety and lifetimes eliminated common C vulnerabilities, enhancing robustness without sacrificing compatibility. Performance remains comparable for parsing tasks, with minor overhead from Rust's safety checks. Challenges in replicating C4's lenient parsing and unsafe operations were addressed through careful design and error handling. The project, completed up to integration testing, provides a solid foundation for future extensions, such as a VM, while demonstrating Rust's suitability for systems programming tasks like compilers.