

Bootstrapping the Lateralus Compiler from Python

bad-antics | July 2024 | Compilers

Abstract

This paper describes the process of bootstrapping the Lateralus compiler using Python as the initial implementation language. The bootstrap process covers lexer, parser, type checker, ownership checker, and RISC-V code generator implementation, followed by self-hosting through a three-phase bootstrapping process.

1 Introduction

Bootstrapping a compiler involves writing the compiler for a language in the language itself. Before the language exists, a bootstrap compiler must be written in an existing language. This paper describes bootstrapping the Lateralus compiler using Python as the bootstrap implementation language.

The bootstrap process has three phases: Phase 1 writes a minimal Lateralus compiler in Python. Phase 2 rewrites the compiler in Lateralus and compiles it using the Phase 1 compiler. Phase 3 uses the self-hosted compiler to compile itself, verifying correctness through comparison.

Python is chosen as the bootstrap language for its rapid prototyping capabilities, rich standard library, and excellent string manipulation. The Phase 1 compiler prioritizes correctness and completeness over performance.

2 Language Subset Selection

The bootstrap subset includes the minimum language features needed to write a compiler: functions, structures, enumerations, pattern matching, ownership (move semantics), borrowing, pipelines, basic I/O, and string manipulation.

Excluded from the bootstrap subset: generics (monomorphized manually), closures (replaced with function pointers), `async/await`, operator overloading, and the full standard library. These features are added after bootstrapping.

The subset is designed to be self-sufficient: the compiler for the full language can be written using only the subset features. This property is verified by implementing the full compiler in the subset.

3 Lexer Implementation

The lexer converts source text into a stream of tokens. Each token has a type (keyword, identifier, literal, operator, punctuation), a lexeme (the source text), and a source location (file, line, column).

Keyword recognition uses a hash table mapping keyword strings to token types. Keywords include:

pipe, fn, let, mut, if, else, match, struct, enum, return, loop, while, for, in, break, continue, and use.

Numeric literal lexing handles: integers (42), hexadecimal (0xFF), binary (0b1010), octal (0o17), and floating-point (3.14, 1.0e-5). Underscores in numeric literals (1_000_000) are stripped during lexing.

String literal lexing handles: regular strings (double-quoted), raw strings (r-prefix, no escape processing), and multi-line strings (triple-quoted). Escape sequences include: newline, tab, backslash, and Unicode.

Comment handling skips: line comments (`//` to end of line) and block comments (`/*` to `*/`). Block comments nest, allowing commented-out code that contains block comments.

4 Parser Implementation

The parser transforms the token stream into an abstract syntax tree (AST). The parser uses recursive descent with operator precedence climbing for expressions.

Module parsing: a module consists of a sequence of items: function definitions, structure definitions, enumeration definitions, use declarations, and constant definitions.

Expression parsing uses precedence climbing. Operators are grouped by precedence level: assignment (`=`), logical or (`||`), logical and (`&&`), comparison (`==`, `!=`, `<`, `>`, `<=`, `>=`), addition (`+`, `-`), multiplication (`*`, `/`, `%`), unary (`-`, `!`), and pipeline (`|>`).

Pipeline expression parsing: the `|>` operator has its own precedence level. `'a |> f |> g'` parses as `'g(f(a))'`. Pipeline expressions can chain arbitrarily.

Pattern matching parsing: the match expression contains a scrutinee expression and a list of arms. Each arm has a pattern and a body expression. Patterns include: literal, variable binding, tuple destructuring, enum variant, and wildcard (`_`).

5 Type Checking

Type inference assigns types to all expressions and variables. The type checker uses bidirectional type checking: known types propagate downward (checking mode) and inferred types propagate upward (synthesis mode).

Structure types are nominal: two structures with identical fields but different names are distinct types. Structure field access is checked at compile time. Missing or extra fields are type errors.

Enumeration types define a set of named variants. Each variant can carry data (like a tuple). Pattern matching on enumerations must be exhaustive: all variants must be covered.

Function type checking verifies that: argument types match parameter types, the return type matches the declared return type, and all code paths return a value of the correct type.

Pipeline type checking verifies that the output type of each stage matches the input type of the next

stage. Type mismatches in pipeline connections are reported with the specific stage types.

6 Ownership and Borrow Checking

Ownership checking ensures that every value has exactly one owner at any time. Moving a value transfers ownership. Using a moved value is a compile-time error.

Borrow checking ensures that references to values are valid for their entire lifetime. Mutable references are exclusive (no other references exist). Immutable references can coexist.

Lifetime analysis in the bootstrap compiler uses a simplified region-based approach. Each scope defines a region. References created in a scope cannot outlive that scope.

Drop insertion adds destructor calls at the end of each scope. Values that own heap-allocated data are dropped when they go out of scope. The drop order is the reverse of the declaration order.

7 Code Generation

The bootstrap compiler generates RISC-V assembly. The code generator walks the AST and emits assembly instructions for each node. The output is a .s file that is assembled and linked by the system assembler and linker.

Register allocation uses a simple linear scan algorithm. The allocator assigns variables to registers (a0-a7, s0-s11, t0-t6) and spills excess variables to the stack.

Function calling convention follows RISC-V standard: arguments in a0-a7, return value in a0, callee-saved registers s0-s11, caller-saved registers t0-t6 and a0-a7.

Stack frame layout: saved return address (ra), saved frame pointer (s0), callee-saved registers, local variables, and temporary spill slots. The frame pointer points to the saved ra.

Pipeline code generation transforms 'a |> f |> g' into sequential function calls: compute a, call f with the result, call g with the result. Each stage result is stored in a temporary register.

8 Self-Hosting Process

Phase 2 rewrites the compiler in Lateralus. The Lateralus source code is compiled by the Phase 1 (Python) compiler to produce a native binary. This native binary is the first self-hosted compiler.

Compiler comparison verifies Phase 2 correctness: both the Phase 1 and Phase 2 compilers compile the same test suite. The outputs are compared byte-by-byte. Differences indicate bugs in Phase 2.

Phase 3 compiles the Phase 2 source code using the Phase 2 compiler itself. The resulting binary is compared with the Phase 2 binary. Identical output proves that the compiler correctly compiles itself.

Test suite coverage includes: lexer tests (all token types), parser tests (all grammar productions),

type checker tests (type errors and successes), and code generation tests (all instruction patterns).

9 Optimization Passes

Constant folding evaluates constant expressions at compile time. Arithmetic on literals, string concatenation of literals, and boolean expressions with constant operands are folded.

Dead code elimination removes unreachable code: code after return statements, unused functions, and branches with constant conditions.

Common subexpression elimination identifies repeated computations and replaces them with a single computation stored in a temporary variable.

10 Conclusion

Bootstrapping the Lateralus compiler from Python demonstrates that the language is expressive enough for large-scale software development. The self-hosting process validates the compiler's correctness and the language's practicality.

3.1 Lexer Architecture

The lexer is implemented as a pipeline stage: it takes a character stream and produces a token stream. The pipeline interface enables lazy tokenization: tokens are produced on demand.

Character classification uses lookup tables for performance. ASCII characters are classified into categories: alphabetic, digit, whitespace, operator, and punctuation. Unicode characters use the Unicode category database.

Error recovery in the lexer skips to the next whitespace or delimiter when an unrecognized character sequence is encountered. The error is recorded with its location for later reporting.

Token lookahead provides peek (look at the next token without consuming it) and peek_ahead (look multiple tokens ahead). Lookahead is implemented using a small buffer of pre-lexed tokens.

Source location tracking maintains the current file, line, and column. Newlines increment the line counter and reset the column. Tabs advance the column to the next multiple of 4.

Interpolated string lexing handles string literals with embedded expressions: 'Hello, {name}!'. The lexer produces alternating string literal and expression tokens.

Operator lexing handles multi-character operators: ==, !=, <=, >=, |>, ->, =>, &&, ||, <<, >>, and compound assignment operators (+=, -=, *=, /=).

Indentation-sensitive mode (optional) uses indentation to delimit blocks instead of braces. The lexer inserts synthetic INDENT and DEDENT tokens based on indentation changes.

Token stream filtering removes whitespace and comment tokens for the parser. The raw token

stream (including whitespace) is available for formatting tools.

Lexer performance: the Python lexer processes approximately 1 million characters per second. The self-hosted lexer processes 50 million characters per second.

4.1 AST Node Types

Expression nodes: literal (integer, float, string, boolean), identifier, binary operation, unary operation, function call, pipeline, field access, index, tuple, array, match, if-else, block, and lambda.

Statement nodes: let binding (with optional type annotation), assignment, expression statement, return, break, continue, while loop, for loop, and loop (infinite).

Item nodes: function definition (with parameters, return type, and body), structure definition (with fields), enumeration definition (with variants), use declaration, and constant definition.

Pattern nodes: literal pattern, identifier pattern (binds the matched value), tuple pattern, enum variant pattern (with nested patterns), wildcard pattern, and or-pattern (p1 | p2).

Type annotation nodes: named type, tuple type, function type, reference type (& and &mut), array type, and generic type application (e.g., Vec[T]).

AST pretty printing formats the tree as indented text for debugging. Each node type has a custom printing format. Pretty printing is used to verify parser correctness.

Source span attachment: every AST node carries the source span (start and end location) of the original source text. Source spans enable precise error messages.

AST visitor pattern provides a base class with visit methods for each node type. Compiler passes extend the visitor to implement specific transformations.

AST mutation supports in-place modification for optimization passes. Mutable visitors can replace nodes, insert nodes, and delete nodes.

AST serialization writes the tree to JSON for external tool consumption. The JSON format includes all node types, source spans, and type annotations.

5.1 Type Inference Algorithm

The type inference algorithm uses Hindley-Milner type inference extended with ownership types. Type variables are introduced for unknown types and unified with concrete types during inference.

Unification resolves type constraints by finding a substitution that makes two types equal. Unification handles: concrete types (exact match), type variables (substitution), and function types (recursive unification).

Constraint generation walks the AST and generates type constraints for each expression. A function call generates constraints: argument types equal parameter types, and the result type equals the return type.

Constraint solving processes constraints in order, applying substitutions eagerly. Each solved constraint may simplify subsequent constraints. Unsolvable constraints are reported as type errors.

Type error messages include: the expected type, the actual type, and the source location. For pipeline type mismatches, the message shows both stage types and the pipeline connection point.

Recursive type detection identifies infinite types (type variable appears in its own solution) and reports them as errors. Recursive types are supported explicitly through struct definitions.

Let-polymorphism allows let-bound functions to be used at multiple types. The function's type is generalized (type variables are made polymorphic) at the let binding and instantiated at each use.

Type checking for match expressions verifies exhaustiveness: all possible values of the scrutinee type are covered by the patterns. Missing patterns are reported as warnings.

Type checking for ownership ensures that moved values are not used after the move. The checker tracks the ownership state of each variable: owned, moved, borrowed, or mutably borrowed.

Type environment management: the type checker maintains a stack of scopes. Each scope maps variable names to types. Entering a block pushes a scope. Leaving a block pops a scope.

6.1 Ownership Analysis Details

Move semantics analysis tracks value moves through: assignment (let $y = x$ moves x), function call ($f(x)$ moves x), return (return x moves x), and pipeline ($x \mid > f$ moves x).

Copy types are exempt from move semantics. Integer, float, boolean, and character types are implicitly copied. User-defined types can opt into copy semantics if all fields are copy types.

Partial move analysis tracks moves of individual struct fields. After moving one field, the remaining fields are still accessible but the whole struct is partially moved.

Move in loops analysis ensures that values are not moved in loop iterations after the first. A value moved in a loop body must be re-assigned before the next iteration.

Borrow scope analysis computes the region during which a borrow is active. The borrow starts at the borrow expression and ends at the last use of the reference. Nested borrows create nested regions.

Mutable borrow exclusivity: when a mutable borrow is active, no other borrows (mutable or immutable) of the same value can be active. The checker tracks active borrows for each value.

Temporary value lifetimes: temporary values created in expressions live until the end of the enclosing statement. References to temporaries cannot be stored in variables.

Reborrowing: a mutable reference can be reborrowed as a shorter-lived mutable reference. The original reference is suspended during the reborrow and becomes active again after.

Ownership error messages explain: which value was moved, where it was moved, and where the invalid use occurs. Suggestions for fixing the error are included when possible.

Ownership analysis complexity is linear in the number of statements. Each statement is analyzed once. The analysis visits each variable use to check its ownership state.

7.1 Assembly Generation Details

Integer arithmetic generates RISC-V integer instructions: add, sub, mul, div, rem for signed operations. Unsigned variants use: divu, remu. Shift operations: sll, srl, sra.

Comparison operations generate: beq (branch if equal), bne (branch not equal), blt (branch if less than), bge (branch if greater or equal). Unsigned comparisons use bltu, bgeu.

Boolean operations generate: and, or, xor for bitwise operations. Short-circuit evaluation for && and || uses conditional branches to skip the right operand.

Function call code generation: evaluate arguments into a0-a7, save caller-saved registers that are live across the call, emit the call instruction, restore saved registers, move the result from a0.

Function prologue: save ra and s0 to the stack, set up s0 as frame pointer, allocate stack space for local variables and spill slots. Function epilogue: restore ra and s0, deallocate stack, ret.

If-else code generation: evaluate the condition, branch to the else label if false, generate the then body, jump to the end label, generate the else body at the else label.

While loop code generation: loop start label, evaluate the condition, branch to the end label if false, generate the body, jump to the loop start label.

Match expression code generation: evaluate the scrutinee, test each pattern in order (branch to the next pattern if the test fails), execute the body of the matching arm.

String literal code generation stores string data in the .rodata section. Each string is labeled and its address is loaded using the la pseudo-instruction.

Global variable code generation stores global variables in the .data section (initialized) or .bss section (uninitialized). Access uses the la instruction followed by load/store.

8.1 Phase 1 Implementation Details

The Phase 1 compiler is 8,500 lines of Python. The lexer is 600 lines, the parser is 2,200 lines, the type checker is 2,800 lines, and the code generator is 2,900 lines.

Python AST classes use dataclasses for concise definitions. Each AST node is a dataclass with typed fields. Pattern matching on AST nodes uses isinstance checks.

Error reporting in Phase 1 collects all errors and reports them at the end of each compilation phase. Errors include: the error message, source location, and error category.

Phase 1 compilation speed: 200 lines per second. Compiling the Phase 2 compiler (12,000 lines) takes 60 seconds. The speed is acceptable for bootstrapping but not for production use.

Phase 1 testing uses pytest with 450 test cases. Each test case is a small Lateralus program with expected output. The test runner compiles the program, runs it, and compares the output.

Debugging the Phase 1 compiler uses Python's pdb debugger and extensive print-based tracing. Each compiler phase can be dumped to a file for inspection.

Memory usage of the Phase 1 compiler peaks at 500 MB for the largest compilation unit. Python's garbage collector manages memory. The AST is the largest data structure.

Cross-platform support: the Phase 1 compiler runs on any platform with Python 3.10+. The generated assembly is RISC-V specific and requires a RISC-V toolchain for assembly and linking.

Build system for Phase 1 uses a Makefile that: runs the Python compiler on the Lateralus source, assembles the output with gas (GNU Assembler), and links with ld.

Version control strategy: the Phase 1 source is maintained in a separate repository. After bootstrapping succeeds, the Phase 1 compiler is archived and development continues on the self-hosted compiler.

9.1 Optimization Details

Constant propagation replaces variable uses with their constant values when the variable is assigned a constant and never reassigned. Propagation enables further constant folding.

Strength reduction replaces expensive operations with cheaper ones: multiplication by powers of two becomes left shift. Division by powers of two becomes right shift (for unsigned) or arithmetic right shift (for signed).

Function inlining replaces a function call with the function body. Inlining is applied to small functions (fewer than 20 instructions) called from a single site. Inlining eliminates call overhead.

Tail call optimization converts tail-recursive calls to loops. A tail call is a function call that is the last operation in the function. The optimization reuses the current stack frame.

Copy propagation replaces uses of a variable that was assigned from another variable with the original variable. Copy propagation reduces register pressure by eliminating unnecessary copies.

Register coalescing merges two virtual registers that are connected by a move instruction. If the merged register does not conflict with other registers, the move is eliminated.

Instruction selection uses pattern matching to select RISC-V instructions. Patterns match AST sub-trees to instruction sequences. Complex patterns generate fewer instructions than node-by-node selection.

Peephole optimization scans the generated assembly for: redundant load-store pairs, jumps to the next instruction, and unnecessary register moves. Each pattern has a replacement.

Loop-invariant code motion moves computations that produce the same result on every iteration to before the loop. The optimization reduces the number of instructions executed per iteration.

Optimization pipeline: constant propagation -> dead code elimination -> copy propagation -> register coalescing -> instruction selection -> peephole optimization. Each pass is run once.

2.1 Bootstrapping Theory

A T-diagram represents a compiler: the source language at the top-left, the target language at the top-right, and the implementation language at the bottom. Bootstrapping chains T-diagrams.

The bootstrapping chain: Lateralus->RISC-V (in Python) compiles Lateralus->RISC-V (in Lateralus) to produce Lateralus->RISC-V (native). The native compiler compiles itself to verify.

Trusting trust: Ken Thompson's 1984 paper showed that a compiler can contain a Trojan horse that replicates through bootstrapping. Diverse double-compilation detects such attacks.

Diverse double-compilation compiles the compiler with two independently developed compilers. The outputs are compared. A Trojan horse in one compiler would produce different output.

Tombstone diagrams extend T-diagrams with interpreters and programs. The bootstrapping process can use interpreters for intermediate stages.

Incremental bootstrapping adds language features one at a time. Each increment uses only features available in the previous stage. This avoids circular dependencies.

Feature ordering for incremental bootstrapping: basic types and functions first, then structures and enums, then ownership checking, then pipelines, then pattern matching, then the remaining features.

Bootstrap compiler size minimization: the Phase 1 compiler implements only the features needed to compile Phase 2. Extra features increase Phase 1 complexity without benefit.

Historical bootstrapping examples: GCC was bootstrapped from a Pastel compiler, Go was bootstrapped from C, and Rust was bootstrapped from OCaml.

Bootstrapping verification: the Phase 3 output must be identical to the Phase 2 output. This fixed-point property proves that the compiler produces consistent output.

4.2 Error Recovery in Parsing

Synchronization tokens are used for error recovery. After a parse error, the parser skips tokens until it finds a synchronization token (semicolon, closing brace, or keyword). Parsing resumes from there.

Error productions in the grammar handle common mistakes. For example, 'if (condition)' with parentheses (C-style) produces a helpful error: 'Parentheses not needed around if condition'.

Error node insertion: when the parser expects an expression but finds something else, it inserts an error node in the AST. The error node carries the error message and allows parsing to continue.

Maximum error count limits the number of errors before the parser gives up. The default limit is 50 errors. Cascading errors (caused by a single mistake) inflate the error count.

Error message quality: messages include the expected token(s), the found token, and the source location. For ambiguous situations, multiple suggestions are provided.

Parenthesis matching: the parser tracks unmatched opening parentheses, brackets, and braces. An unmatched opening delimiter at the end of the file reports the location of the opening delimiter.

Operator precedence errors: using `=` where `==` was intended is detected by the type checker (assignment in expression context). The error message suggests the correct operator.

Missing semicolon recovery: if a statement parse fails, the parser tries inserting a semicolon before the current token. If the parse succeeds with the inserted semicolon, a helpful error is reported.

Indentation-based recovery: when braces are missing, the parser uses indentation to guess block boundaries. This recovery works well for common formatting styles.

Error recovery testing uses a corpus of programs with known errors. Each program should produce exactly the expected error messages. New error recovery features are tested against this corpus.

5.2 Ownership Type System

The ownership type system adds ownership qualifiers to the base type system. Each type has a qualifier: owned (default), borrowed (`&T`), or mutably borrowed (`&mut T`).

Function signatures declare ownership for each parameter: owned parameters take ownership, borrowed parameters borrow, and mutably borrowed parameters borrow mutably.

Return type ownership: functions that return owned values transfer ownership to the caller. Functions that return references must declare the reference lifetime.

Ownership inference for local variables: the type checker infers ownership transfers from the usage pattern. Explicit annotations are required only for function signatures.

Ownership and generics: generic type parameters carry ownership information. A generic container `Vec[T]` owns its elements. Borrowing a `Vec` borrows all elements.

Ownership and enums: each enum variant owns its contained data. Pattern matching on an owned enum moves the contained data out. Pattern matching on a borrowed enum borrows the contained data.

Ownership and closures: closures capture variables by reference or by move. The capture mode is inferred from usage. Explicit move capture is available with the `move` keyword.

Self parameter ownership: methods declare `self` as owned (`self`), borrowed (`&self`), or mutably borrowed (`&mut self`). The caller's value is moved, borrowed, or mutably borrowed accordingly.

Ownership checking performance: the ownership analysis runs in $O(n)$ time where n is the number of statements. The analysis uses a single forward pass over the control flow graph.

Comparison with Rust's borrow checker: Lateralus uses a simplified region-based approach without

explicit lifetime annotations. The simplification handles common cases but rejects some programs Rust would accept.

7.2 Memory Layout

Integer types: i8 (1 byte), i16 (2 bytes), i32 (4 bytes), i64 (8 bytes). Unsigned variants: u8, u16, u32, u64. All integers are stored in little-endian byte order on RISC-V.

Floating-point types: f32 (4 bytes, IEEE 754 single precision), f64 (8 bytes, IEEE 754 double precision). Floating-point values are stored in the FPU registers (f0-f31).

Boolean type: stored as a single byte. True is 1, false is 0. Booleans are zero-extended when loaded into registers.

Tuple type: fields are stored contiguously in memory with natural alignment. A tuple (i32, i64) is stored as: 4 bytes (i32), 4 bytes padding, 8 bytes (i64) = 16 bytes total.

Struct type: fields are stored in declaration order with natural alignment. Padding is inserted between fields as needed. The struct size is rounded up to the alignment of its largest field.

Enum type: a tag byte followed by the largest variant's data. The tag identifies the active variant. The total size is 1 byte (tag) + padding + max variant size, rounded to alignment.

Reference type: stored as a pointer (8 bytes on RV64). The pointer contains the virtual address of the referenced value. Null references are not allowed.

Array type: stored as a pointer (8 bytes) and a length (8 bytes). The pointer points to heap-allocated memory containing the array elements. Elements are stored contiguously.

String type: stored as a pointer (8 bytes), length (8 bytes), and capacity (8 bytes). Strings are heap-allocated and UTF-8 encoded.

Stack allocation: local variables are allocated on the stack in the function's stack frame. Aggregate types (structs, tuples) are stored on the stack unless they escape the function.

8.2 Phase 2 Rewrite Process

The Phase 2 rewrite translates each Python module to a Lateralus module. The translation is mostly mechanical: Python classes become Lateralus structs, Python methods become functions, and Python lists become arrays.

Data structure translation: Python dictionaries become hash maps. Python sets become hash sets. Python tuples become Lateralus tuples. Python None becomes Lateralus Option type.

Error handling translation: Python try/except becomes Lateralus Result types with pattern matching. Each function that can fail returns a Result. Callers handle errors explicitly.

String processing translation: Python's string methods (split, join, replace, find) are reimplemented as Lateralus functions. Regular expressions are replaced with hand-written parsers where possible.

I/O translation: Python's file I/O becomes Lateralus I/O with explicit error handling. Buffered reading and writing use Lateralus's BufReader and BufWriter types.

Testing translation: Python pytest tests become Lateralus test functions. Each test function is marked with a #[test] attribute. The test runner discovers and executes all test functions.

Performance comparison: the Phase 2 compiler is 50x faster than Phase 1 (200 vs 10,000 lines per second). The speedup comes from native code execution and efficient data structures.

Lines of code comparison: Phase 1 is 8,500 lines of Python. Phase 2 is 12,000 lines of Lateralus. The increase is due to explicit type annotations, error handling, and ownership annotations.

Bug discovery during Phase 2: the translation process discovered three bugs in Phase 1: an off-by-one error in the lexer, a missing case in the type checker, and an incorrect register allocation decision.

Phase 2 development time: the rewrite took 6 weeks of full-time work. Most time was spent on: ownership annotation (2 weeks), testing (2 weeks), and debugging code generation differences (2 weeks).

3.2 Token Design

Token categories: keywords (30 tokens), operators (25 tokens), punctuation (12 tokens), literals (5 token types), identifiers (1 token type), and special tokens (EOF, error, newline).

Operator tokens include: arithmetic (+, -, *, /, %), comparison (==, !=, <, >, <=, >=), logical (&&, ||, !), bitwise (&, |, ^, <<, >>), assignment (=, +=, -=, *=, /=, %=), and pipeline (|>).

Punctuation tokens include: parentheses ((,)), brackets ([,]), braces ({, }), comma (,), semicolon (;), colon (:), dot (.), double colon (::), arrow (->), and fat arrow (=>).

Literal tokens carry their value: integer literals (i64 value), float literals (f64 value), string literals (String value), character literals (char value), and boolean literals (bool value).

Token equality: two tokens are equal if they have the same type and the same value. Token equality is used by the parser for matching expected tokens.

Token display: each token type has a human-readable display string used in error messages. For example, 'expected ;, found }' uses the display strings of semicolon and right brace tokens.

Token interning: identifier and keyword strings are interned in a global string table. Interning ensures that identical strings share the same memory and can be compared by pointer.

Token stream operations: peek (look ahead), advance (consume current token), expect (consume if matching, error if not), and match_any (consume if any of several types).

Token buffer: the lexer maintains a small buffer (4 tokens) for lookahead. Most grammar productions need at most 2 tokens of lookahead. The buffer is refilled lazily.

Performance measurement: lexing 100,000 tokens takes 50 milliseconds in Phase 1 (Python) and 1 millisecond in Phase 2 (native). The 50x speedup validates the bootstrapping approach.

6.2 Borrow Checker Implementation

The borrow checker runs after type checking. It analyzes the ownership and borrowing state of each variable at each point in the program. The analysis uses dataflow analysis on the control flow graph.

Control flow graph (CFG) construction transforms the AST into a graph of basic blocks. Each block is a straight-line sequence of statements. Branches and loops create edges between blocks.

Dataflow state: at each program point, the state maps each variable to: owned (value present), moved (value moved away), borrowed (immutable reference exists), or mutably-borrowed (mutable reference exists).

Transfer functions define how each statement updates the dataflow state. Assignment introduces an owned value. Move transfers ownership. Borrow creates a borrowed state. Drop removes the value.

Join function for control flow merges: if a variable is owned on one path and moved on another, the variable is moved after the join (conservative approximation).

Loop analysis iterates the dataflow analysis until the state at the loop header reaches a fixed point. Most loops converge in two iterations.

Error reporting: when the borrow checker detects a violation, it reports: the invalid operation, the conflicting operation (the borrow or move that prevents it), and both source locations.

Borrow checker complexity: the analysis visits each basic block once per dataflow iteration. With at most 2 iterations per loop, the total complexity is $O(n * d)$ where n is the number of statements and d is the loop nesting depth.

Escape analysis determines whether a value can be allocated on the stack (does not escape the function) or must be allocated on the heap (escapes through return or storage in a data structure).

The borrow checker handles: struct field borrows (borrowing a field does not borrow the whole struct), array index borrows (borrowing an element requires a borrow of the array), and reborrowing.

9.2 Backend Optimization

Basic block ordering arranges blocks to minimize the number of unconditional jumps. Fall-through paths are preferred. The ordering algorithm uses a depth-first traversal of the CFG.

Branch optimization inverts branch conditions to enable fall-through. If the taken branch is more common than the not-taken branch, the condition is inverted and the blocks are swapped.

Instruction scheduling reorders independent instructions to avoid pipeline stalls. On RISC-V, load instructions have a 1-cycle latency. Scheduling places independent instructions between a load and its use.

Register spill code optimization minimizes the number of spill and reload instructions. Spill slots are reused for variables with non-overlapping live ranges. Spill stores are hoisted out of loops.

Calling convention optimization: tail calls are detected and converted to jumps. Functions that do not call other functions (leaf functions) skip the frame pointer setup.

ABI-compatible struct passing: small structs (up to 2 registers) are passed in registers. Larger structs are passed by pointer. The calling convention follows the RISC-V psABI.

Immediate optimization: constants that fit in 12 bits use immediate instructions (addi instead of add). Constants that fit in 20 bits use LUI. Larger constants use LUI + ADDI.

Zero register optimization: RISC-V's x0 register is always zero. The code generator uses x0 for: zero initialization, comparison with zero, and conditional moves.

Multiplication optimization: multiplication by small constants (2, 3, 4, 5, 7, 8, 9) uses shift and add sequences instead of the mul instruction, which has higher latency.

Division optimization: unsigned division by constants uses multiplication by the reciprocal. The reciprocal is computed at compile time. The optimization replaces a multi-cycle division with a single multiplication.

10.1 Standard Library Bootstrapping

The standard library is bootstrapped incrementally. The Phase 2 compiler includes a minimal runtime: memory allocation (malloc/free), I/O (read/write), and process control (exit).

String library implementation provides: creation, concatenation, slicing, searching, replacement, formatting, and conversion (to/from integers). Strings are UTF-8 encoded.

Collection library implements: Vec (growable array), HashMap (hash table), HashSet (hash set), and BTreeMap (balanced tree). Each collection is generic over its element type.

I/O library provides: File (file operations), BufReader (buffered reading), BufWriter (buffered writing), stdin/stdout/stderr (standard streams), and path manipulation.

Error handling library provides: Result (success or error), Option (present or absent), and the ? operator (early return on error). Error types implement a common Error trait.

Memory allocation library provides: Box (heap allocation), Rc (reference counting), and Arc (atomic reference counting). The allocator uses the system malloc with a Lateralus wrapper.

Formatting library provides: format strings with positional and named arguments, Display trait for user-facing output, and Debug trait for developer-facing output.

Test framework provides: test discovery (functions marked with #[test]), assertion macros (assert, assert_eq, assert_ne), and test result reporting.

Build system library provides: module dependency resolution, incremental compilation (recompile

only changed modules), and linking.

Documentation generation extracts documentation comments (`///`) from source code and generates HTML documentation. Cross-references between types and functions are linked automatically.

4.3 Grammar Specification

The grammar is specified in EBNF (Extended Backus-Naur Form). The specification is 200 lines and covers all syntactic constructs. The grammar is unambiguous and LL(2).

Expression grammar with precedence: each precedence level is a separate production. Lower precedence levels reference higher ones. The highest level is primary expressions (literals, identifiers, parenthesized).

Statement grammar: a statement is either a let binding, an expression statement (expression followed by semicolon), or a control flow statement (if, while, for, loop, return, break, continue).

Item grammar: an item is a function definition, struct definition, enum definition, use declaration, const definition, or impl block. Items can appear at the top level or inside modules.

Disambiguation rules: the parser resolves ambiguities using: operator precedence (defined by the expression grammar), the longest match rule (identifiers consume as many characters as possible), and context (`<` is comparison in expressions, generic in types).

Grammar evolution: new syntax is added through a proposal process. Each proposal includes: motivation, grammar changes, examples, and migration plan. The grammar version is tracked.

Grammar testing: a corpus of valid and invalid programs tests the grammar. Valid programs must parse successfully. Invalid programs must produce specific error messages.

Left recursion elimination: the grammar avoids left recursion for recursive descent compatibility. Operator precedence climbing handles left-associative binary operators.

Dangling else resolution: the grammar associates else with the nearest unmatched if. This is enforced by the parser's greedy matching of the else clause.

Semicolon inference: the parser optionally infers semicolons at end-of-line when the next line starts a new statement. Inference reduces syntactic noise while maintaining unambiguous parsing.

7.3 Linker and Object Format

The compiler generates ELF object files. Each compilation unit produces a `.o` file containing: text section (code), data section (initialized globals), bss section (uninitialized globals), and symbol table.

Symbol visibility: public symbols are accessible from other modules. Private symbols are local to the compilation unit. The compiler marks symbols as global or local in the ELF symbol table.

Relocation entries describe how the linker should fix up addresses. Common relocation types: `R_RISCV_CALL` (function call), `R_RISCV_HI20/LO12` (address loading), and `R_RISCV_BRANCH`

(conditional branch).

Static linking combines all object files into a single executable. The linker resolves symbol references, applies relocations, and generates the final ELF executable.

Linker script specifies memory layout: text at 0x80000000 (for bare-metal) or default (for Linux). The script defines: entry point, section placement, and memory regions.

Debug information generation uses DWARF format. Debug info includes: line number tables (source line to instruction mapping), variable locations, and type descriptions.

Incremental linking: the build system tracks module dependencies and relinks only when object files change. Incremental linking reduces build time for large projects.

Link-time optimization (LTO) passes the intermediate representation (instead of machine code) to the linker. The linker runs optimization passes across module boundaries.

Cross-compilation: the compiler generates code for a target architecture that differs from the host. Cross-compilation is used for: embedded systems, different operating systems, and testing.

Binary size optimization: the linker removes unused functions and data (garbage collection). Section-level granularity ensures that only referenced sections are included in the executable.

8.3 Testing Infrastructure

Unit tests for each compiler component: 150 lexer tests, 200 parser tests, 250 type checker tests, 200 ownership checker tests, and 300 code generation tests. Total: 1,100 unit tests.

Integration tests compile complete programs and compare output. The test suite includes 200 programs ranging from simple expressions to complex data structures and algorithms.

Snapshot testing captures compiler output (AST, type annotations, generated assembly) and compares against saved snapshots. Snapshot tests detect unintended output changes.

Fuzz testing generates random Lateralus programs and checks that the compiler does not crash. The fuzzer uses grammar-based generation to produce syntactically plausible programs.

Cross-compilation testing compiles the test suite on x86_64 and runs on RISC-V (via QEMU). Cross-compilation tests verify that the compiler and runtime work correctly on the target platform.

Performance regression testing measures compilation speed and executable size. Performance is tracked over time. Regressions exceeding 10% are investigated.

Error message testing verifies that incorrect programs produce specific error messages. Each test case contains a program with a deliberate error and the expected error text.

Bootstrap testing: the full bootstrapping process (Phase 1 -> Phase 2 -> Phase 3) is run as a CI job. Phase 3 output must be identical to Phase 2 output.

Code coverage measurement tracks which lines of the compiler are exercised by the test suite.

Coverage target is 90% for all components. Uncovered code is reviewed for dead code or missing tests.

Mutation testing introduces small changes (mutations) to the compiler and verifies that tests detect the changes. Mutations that survive (undetected) indicate weak test coverage.

2.2 Python Implementation Strategy

Python class hierarchy mirrors the Lateralus type system. Abstract base classes define the AST node interface. Concrete classes implement specific node types.

Visitor pattern in Python uses double dispatch: each AST node has an `accept` method that calls the appropriate visit method on the visitor. Each compiler pass is a visitor.

Python type hints (from `typing` module) provide static type checking for the compiler code. `Mypy` verifies type consistency. Type hints serve as documentation for the Phase 2 rewrite.

Python performance optimization: memoization for frequently computed results (string hashing, type comparison), generator-based iteration (lazy evaluation), and `__slots__` for memory-efficient AST nodes.

Python testing with hypothesis generates random AST nodes for property-based testing. Properties include: parsing followed by printing produces the original source, and type checking is idempotent.

Error handling in Python uses exception hierarchies: `CompileError` (base), `LexError`, `ParseError`, `TypeError`, and `OwnershipError`. Each error carries source location and descriptive message.

Python packaging: the Phase 1 compiler is packaged as a pip-installable package. The package includes: compiler module, driver script, standard library source, and test suite.

Command-line interface uses `argparse` for: source file specification, output file specification, optimization level, debug output, and verbose mode.

Configuration management: compiler options are stored in a `Config` dataclass. The config is passed to each compiler phase. Options include: optimization level, target architecture, and debug flags.

Logging in Phase 1 uses Python's logging module. Each compiler phase has its own logger. Verbosity levels: `ERROR` (compile errors), `WARNING` (non-fatal issues), `INFO` (phase progress), and `DEBUG` (detailed internals).

7.4 Pipeline Code Generation

Pipeline expressions are the distinguishing feature of Lateralus code generation. The `|>` operator threads data through a sequence of transformations, each compiled as a function call.

Simple pipeline: `'x |> f |> g'` compiles to: load `x` into `a0`, call `f` (result in `a0`), call `g` (result in `a0`). No temporary storage is needed when stages are simple function calls.

Pipeline with closures: `'x |> map(fn(a) a + 1) |> filter(fn(a) a > 0)'` compiles to: create closure objects

on the stack, pass the closure and input to map, pass the closure and result to filter.

Pipeline fusion optimization: consecutive map operations 'x |> map(f) |> map(g)' are fused into 'x |> map(fn(a) g(f(a)))'. Fusion eliminates intermediate allocations.

Pipeline error propagation: 'x |> f? |> g?' compiles to: call f, check for error (branch to error handler if so), call g, check for error. The ? operator generates branch-on-error code.

Pipeline parallelization: the compiler can detect independent pipeline stages and generate parallel code using RISC-V hardware threads. This is an experimental optimization.

Pipeline debugging: debug builds insert trace calls between pipeline stages. Each trace records the stage name, input value, and output value. Traces are written to a log file.

Pipeline type specialization: generic pipeline functions are monomorphized (specialized for each concrete type). Specialization eliminates virtual dispatch overhead.

Multi-value pipelines: stages that produce multiple outputs use tuple types. 'x |> split |> (map(f), map(g)) |> merge' generates code that unpacks the tuple, processes each element, and repacks.

Pipeline stage inlining: small pipeline stages (fewer than 10 instructions) are inlined at the pipeline site. Inlining eliminates function call overhead for simple transformations.

5.3 Type Checking Pipelines

Pipeline type checking verifies that each stage's output type matches the next stage's input type. The checker processes stages left-to-right, propagating the type through the pipeline.

Stage function type extraction: for a stage 'f', the checker looks up f's type. If f has type 'A -> B', the stage takes A as input and produces B. Partial application is handled by currying.

Pipeline type error messages show the pipeline chain with types annotated at each connection point. Mismatched types are highlighted with an arrow indicating the connection.

Generic pipeline functions: 'identity |> map[T](f)' requires instantiating T with the concrete type flowing through the pipeline. Type inference resolves T from the pipeline context.

Pipeline with method syntax: 'x |> .len()' is sugar for 'x |> fn(a) a.len()'. The type checker resolves the method on the type flowing through the pipeline.

Recursive pipeline type checking: pipelines can contain nested pipelines. The checker processes nested pipelines first, then uses their result types in the outer pipeline.

Pipeline ownership checking: each stage receives ownership of its input (unless the input is borrowed). The borrow checker verifies that borrowed pipeline inputs are not used after the pipeline.

Pipeline type inference with holes: the programmer can write 'x |> f |> _ |> g' where _ is inferred. The checker searches for a function that connects f's output type to g's input type.

Bidirectional pipeline checking: when the expected output type of a pipeline is known (from context),

the checker propagates this type backward to resolve ambiguities.

Pipeline type checking performance: type checking a pipeline of n stages takes $O(n)$ time. Each stage type check is constant time (function type lookup and unification).

9.3 Optimization Metrics

Compilation speed benchmarks: Phase 1 compiles 200 lines/second. Phase 2 compiles 10,000 lines/second. Phase 2 with optimizations compiles 5,000 lines/second (optimizations add overhead).

Code size benchmarks: unoptimized output averages 50 bytes per source line. Optimized output averages 35 bytes per source line. Dead code elimination provides the largest size reduction.

Runtime performance benchmarks: optimized Lateralus code runs within 2x of equivalent C code. Unoptimized code runs within 5x of C. Pipeline fusion is the most impactful optimization.

Compilation memory usage: Phase 1 uses 500 MB for the largest compilation unit (the compiler itself). Phase 2 uses 50 MB for the same input. The reduction is due to efficient data structures.

Optimization pass timing: constant propagation takes 5% of optimization time, dead code elimination takes 10%, copy propagation takes 5%, register allocation takes 40%, and instruction selection takes 40%.

Register allocation quality: the linear scan allocator produces code with 15% more spills than optimal (computed by integer linear programming). The quality is acceptable for compilation speed.

Instruction count metrics: a typical function of 50 source lines generates 200 instructions unoptimized and 140 instructions optimized. Peephole optimization removes 15% of instructions.

Branch prediction metrics: branch optimization (condition inversion and block reordering) improves branch prediction accuracy from 75% to 88% on typical workloads.

Cache performance: basic block ordering improves instruction cache hit rate from 92% to 97%. The improvement is significant for large programs with many function calls.

Link-time optimization impact: LTO provides an additional 10% code size reduction and 8% runtime speedup. The benefit comes from cross-module inlining and dead code elimination.

10.2 Future Compiler Improvements

Incremental compilation: recompile only functions whose source or dependencies changed. The compiler maintains a dependency graph and invalidates affected functions.

Parallel compilation: compile independent modules in parallel using worker threads. Module dependency analysis determines the compilation order and parallelism opportunities.

LLVM backend: generate LLVM IR instead of RISC-V assembly. The LLVM backend enables: optimized code for multiple architectures, battle-tested optimization passes, and debugging support.

x86-64 backend: direct x86-64 code generation for development machines. This eliminates the need for RISC-V emulation during development and testing.

Wasm backend: generate WebAssembly for running Lateralus programs in web browsers. The Wasm backend enables: web-based IDEs, playground websites, and web application development.

Incremental type checking: type check only changed functions and their dependents. The type checker caches type information for unchanged functions.

Better error recovery: implement error correction that suggests fixes for common mistakes. The compiler would suggest: missing semicolons, typos in identifiers, and incorrect operator usage.

Profile-guided optimization: use runtime profiling data to guide optimization decisions. Hot functions are optimized more aggressively. Cold functions are optimized for size.

Auto-vectorization: detect loops that can be vectorized and generate RISC-V vector instructions (V extension). Auto-vectorization improves performance for numerical computation.

Compile-time evaluation: evaluate constant expressions and pure functions at compile time. Compile-time evaluation enables: computed constants, string formatting at compile time, and configuration validation.

6.3 Ownership in Practice

Common ownership patterns: builder pattern (returns self by move), iterator pattern (borrows the collection), factory pattern (returns owned values), and observer pattern (borrows the subject).

Ownership migration guide for Python developers: Python uses garbage collection and shared references. Lateralus requires explicit ownership. Key differences: no aliasing of mutable data, explicit clone for copying, and scope-based deallocation.

Ownership debugging techniques: the compiler emits ownership annotations in debug builds. Each variable is annotated with its ownership state at each program point. A visualization tool displays the annotations.

Ownership and unsafe code: the unsafe block disables ownership checking. Unsafe code is used for: raw pointer manipulation, FFI calls, and performance-critical sections. Unsafe blocks are reviewed carefully.

Ownership cost analysis: ownership checking adds 5% compilation time overhead. Runtime overhead is zero (all checks are at compile time). Memory overhead is reduced (no garbage collector, no reference counting).

Common ownership errors and solutions: use-after-move (clone the value or restructure the code), double borrow (reduce borrow scope or use Cell), and lifetime errors (restructure to remove the problematic reference).

Ownership and concurrency: ownership prevents data races. Shared mutable state requires

synchronization (Mutex, RwLock). The type system enforces that shared state is properly synchronized.

Interior mutability: Cell (for Copy types) and RefCell (for any type) allow mutation through shared references. RefCell checks borrow rules at runtime and panics on violation.

Ownership and FFI: foreign functions receive raw pointers. The caller is responsible for managing ownership of data passed to foreign functions. Wrapper functions provide safe interfaces.

Automatic reference counting: the Rc type provides shared ownership with automatic deallocation when the last reference is dropped. Rc is used for tree structures and shared configuration.

3.3 Unicode Support

The lexer supports Unicode identifiers. Identifiers can contain letters, digits, and underscores from any Unicode script. The first character must be a letter or underscore.

Unicode normalization: identifiers are normalized to NFC (Canonical Decomposition followed by Canonical Composition) before comparison. This ensures that visually identical identifiers are treated as equal.

Confusable character detection: the compiler warns when two identifiers differ only by confusable characters (e.g., Latin 'a' and Cyrillic 'a'). This prevents security issues from visual ambiguity.

String encoding: source files are UTF-8. String literals contain UTF-8 encoded text. Character literals contain Unicode code points (up to 4 bytes in UTF-8).

Unicode escape sequences: string literals support `\u{XXXX}` for Unicode code points. The code point is validated (must be a valid Unicode scalar value). Invalid code points are compile errors.

Source file encoding detection: the compiler expects UTF-8. A BOM (byte order mark) at the start of the file is ignored. Non-UTF-8 bytes produce lexer errors with the byte offset.

Unicode categories for identifier classification: Lu (uppercase letter), Ll (lowercase letter), Lt (titlecase letter), Lm (modifier letter), Lo (other letter), Nl (letter number), Mn, Mc, Nd, Pc.

Right-to-left text handling: the compiler processes source code left-to-right regardless of Unicode directionality. Bidirectional control characters in source code are rejected with warnings.

Unicode in comments: comments can contain any Unicode text. Documentation comments (`///`) are extracted as Markdown and may contain Unicode formatting characters.

Unicode performance: character classification uses precomputed lookup tables for ASCII (fast path) and binary search on Unicode range tables for non-ASCII (slow path). ASCII-only source files pay no Unicode overhead.

8.4 Bootstrapping Verification

Triple compilation test: compile the compiler source with Phase 1 (P1), Phase 2 (P2), and Phase 2

output (P3). Verify P2 output equals P3 output. This proves the compiler is a fixed point.

Differential testing: compile the test suite with P1 and P2. Run both outputs. Compare execution results. Any difference indicates a bug in either P1 or P2.

Assembly comparison: compare the assembly output of P1 and P2 for each function. Differences are expected (different register allocation, different instruction selection) but semantics must match.

Performance comparison: measure P2 and P3 compilation speed and output quality. P2 and P3 should have identical performance (they are the same binary).

Memory safety verification: run P2 under Valgrind (via RISC-V emulation) to detect memory safety violations. Common issues: buffer overflows, use-after-free, and memory leaks.

Stress testing: compile large generated programs (100,000+ lines) with P2 to test stability. Generated programs exercise: deeply nested expressions, large match statements, and long pipelines.

Correctness proof sketch: if $P2 = P3$ and P2 passes all tests, then the compiler correctly compiles itself. Combined with comprehensive testing, this provides strong confidence in correctness.

Regression testing: each code change triggers the full bootstrap process. Regressions are detected immediately. The bootstrap CI job takes 15 minutes (1 minute P1, 10 minutes P2, 4 minutes testing).

Binary reproducibility: the same source code, compiler version, and flags produce the same binary. Reproducibility is verified by building on different machines and comparing outputs.

Bootstrapping documentation: the bootstrapping process is documented in a step-by-step guide. The guide includes: prerequisites, build commands, verification steps, and troubleshooting.

10.3 Lessons Learned

Start with a complete language specification before implementing. Changes to the specification during implementation cause rework in both the Python and Lateralus compilers.

Python is an excellent bootstrapping language. Rapid prototyping, easy debugging, and rich libraries accelerate Phase 1 development. The 50x performance gap is acceptable for bootstrapping.

Ownership checking is the hardest compiler component to implement correctly. The borrow checker required the most debugging time and the most test cases.

Testing infrastructure investment pays off early. Comprehensive tests caught bugs during Phase 2 development that would have been difficult to debug in the bootstrapped compiler.

The Phase 2 rewrite is easier than writing the compiler from scratch. The Phase 1 implementation serves as a detailed specification. Most bugs are mechanical translation errors.

Pipeline semantics are natural for compiler passes. The compiler itself is structured as a pipeline: source -> tokens -> AST -> typed AST -> IR -> assembly. Each pass is a pipeline stage.

Error message quality matters. Users spend most of their time reading error messages. Investing in

clear, precise, and helpful error messages reduces user frustration.

Register allocation dominates code quality. Moving from simple linear scan to a better allocator would improve performance significantly. This is the highest-priority future optimization.

Cross-platform development requires RISC-V emulation. QEMU provides adequate emulation speed for testing but is too slow for development. A native x86-64 backend would eliminate this bottleneck.

Community feedback drives language evolution. Early adopters identified usability issues with ownership annotations and pipeline syntax that led to significant language improvements.

Documentation is part of the compiler. Doc comments are compiled into the binary metadata. The compiler and documentation share the same source of truth.

Performance measurement is essential from the start. Without baseline measurements, it is impossible to evaluate optimization effectiveness. The compiler tracks 20 performance metrics.

4.4 Pipeline Parsing Details

Pipeline parsing uses left-to-right associativity. 'a |> f |> g |> h' parses as '((a |> f) |> g) |> h'. Each intermediate result is the input to the next stage.

Pipeline with arguments: 'a |> f(b, c)' means 'f(a, b, c)'. The pipeline input is inserted as the first argument. Explicit placeholder syntax `_` allows other positions: 'a |> f(b, _, c)'.

Pipeline blocks: 'a |> { let x = f(_); g(x) }' allows multi-statement pipeline stages. The block receives the pipeline input as `_` and its result is the block's value.

Pipeline method chains: 'x |> .method1() |> .method2()' is sugar for method calls. The dot-prefix indicates a method call on the pipeline value.

Pipeline error handling: 'x |> f? |> g?' uses the `?` operator for early return on error. Each stage that returns `Result` is unwrapped, and errors propagate up the pipeline.

Nested pipelines: 'x |> f(y |> g |> h)' contains a nested pipeline as a function argument. The inner pipeline is parsed first and its result type flows into the outer pipeline.

Pipeline formatting: the formatter places each pipeline stage on its own line, indented. Long pipelines are readable with one transformation per line.

Pipeline precedence: `|>` has lower precedence than arithmetic and comparison but higher than assignment. This allows: 'let result = x |> f |> g' without parentheses.

Pipeline debugging annotations: '#[trace]' before a pipeline inserts runtime tracing for each stage. Trace output shows: stage name, input value, output value, and execution time.

Pipeline desugaring: before type checking, pipeline expressions are desugared to nested function calls. 'a |> f |> g' becomes 'g(f(a))'. The desugared form is type-checked normally.

7.5 Runtime System

The runtime system provides: memory allocation (malloc/free wrappers), panic handling (stack unwinding), and program startup (argument parsing, environment setup, main function call).

Memory allocator: the runtime wraps the system allocator (glibc malloc) with: allocation tracking (for leak detection in debug mode), alignment support, and out-of-memory handling.

Panic handling: when a panic occurs, the runtime: prints the panic message and source location, unwinds the stack (calling destructors for each frame), and exits with code 1.

Stack unwinding uses frame pointer chains. Starting from the current frame, the unwinder follows the saved frame pointers to reconstruct the call stack. Each frame's destructor is called.

Program entry point: the runtime's `_start` function: initializes the heap, parses command-line arguments, calls `main()`, and exits with `main`'s return value.

Debug runtime: in debug mode, the runtime inserts: bounds checks on array access, null reference checks, integer overflow checks, and stack overflow detection.

Runtime size: the minimal runtime is 2 KB of code. The debug runtime is 8 KB. The full runtime with standard library is 50 KB. Size is measured after dead code elimination.

Signal handling: the runtime installs handlers for: SIGSEGV (memory access violation), SIGFPE (arithmetic error), SIGABRT (assertion failure), and SIGBUS (alignment error).

Environment access: the runtime provides functions to: read environment variables, get the program path, get the current directory, and get the system time.

Thread support: the runtime provides: thread creation (`pthread_create` wrapper), mutex (`pthread_mutex`), condition variable (`pthread_cond`), and atomic operations (RISC-V AMO instructions).

5.4 Type System Extensions

Trait system: traits define shared behavior. A trait declares a set of method signatures. Types implement traits by providing method bodies. Trait bounds constrain generic types.

Associated types: traits can declare associated types. Implementors provide concrete types. Associated types reduce the number of generic parameters: `Iterator` has one associated type `Item`.

Where clauses: complex trait bounds use where clauses. `'fn f[T](x: T) where T: Display + Clone'` constrains `T` to implement both `Display` and `Clone`.

Trait objects: dynamic dispatch through trait pointers. A trait object is a fat pointer: data pointer + vtable pointer. Trait objects enable heterogeneous collections.

Operator overloading through traits: `Add`, `Sub`, `Mul`, `Div`, `Rem`, `Neg`, `Not`, `Index`, `Deref`. Implementing `Add` for a type enables the `+` operator.

Coherence rules: each trait implementation must be in the same module as either the trait or the type. This prevents conflicting implementations from different modules.

Blanket implementations: `'impl[T: Display] ToString for T'` provides `ToString` for all types that implement `Display`. Blanket implementations reduce boilerplate.

Sealed traits: traits marked as sealed cannot be implemented outside their defining module. Sealed traits enable exhaustive matching on trait implementors.

Trait inheritance: a trait can require other traits. `'trait Eq: PartialEq'` means that implementing `Eq` requires implementing `PartialEq` first.

Derive macros: common traits (`Debug`, `Clone`, `PartialEq`, `Eq`, `Hash`) can be automatically derived. The compiler generates the implementation based on the struct's fields.

8.5 Cross-Platform Bootstrapping

Primary development platform: x86-64 Linux. The Phase 1 compiler runs natively. The Phase 2 compiler runs under RISC-V emulation via QEMU user mode.

QEMU user mode emulation: QEMU translates RISC-V instructions to x86-64 instructions. Emulation overhead is approximately 10x. Compiling the compiler under emulation takes 10 minutes.

ARM64 support: the Phase 1 compiler generates RISC-V assembly regardless of the host platform. ARM64 developers use QEMU for testing, same as x86-64 developers.

macOS support: the Phase 1 compiler runs on macOS with Python 3.10+. QEMU user mode is replaced by QEMU system mode with a minimal Linux kernel for testing.

Windows support: the Phase 1 compiler runs on Windows with Python 3.10+. WSL2 provides Linux compatibility for QEMU and the RISC-V toolchain.

Docker-based development: a Docker image contains: Python 3.12, RISC-V GCC toolchain, QEMU, and all build dependencies. The image ensures reproducible builds across platforms.

CI/CD pipeline: GitHub Actions runs the full bootstrap process on each push. The CI job uses the Docker image. Bootstrap failures block merges.

Native RISC-V testing: the bootstrap is periodically tested on RISC-V hardware (SiFive Unmatched board). Native testing verifies that emulation accurately models the hardware.

Cross-compilation targets: the compiler generates code for: `rv64gc` (general purpose with compressed instructions), `rv64im` (integer with multiply), and `rv32im` (32-bit embedded).

Target-specific testing: each target is tested with its own test suite. Target-specific features (floating point, compressed instructions, atomics) have dedicated tests.

Build artifact caching: CI caches the Phase 1 output (Phase 2 binary) to speed up subsequent builds. The cache is invalidated when Phase 1 source or the bootstrap process changes.

Release engineering: release builds use: full optimization, LTO, and binary stripping. The release binary is 500 KB. Debug builds are 2 MB with full debug information.

9.4 Whole-Program Analysis

Module dependency graph: the compiler constructs a directed acyclic graph of module dependencies. Cycles in the dependency graph are compile errors. The graph determines compilation order.

Whole-program dead code elimination: starting from the main function, the compiler traces reachable functions. Unreachable functions and their dependencies are removed from the binary.

Whole-program type specialization: generic functions are specialized for each concrete type used in the program. Unused specializations are eliminated. This reduces binary size.

Call graph construction: the compiler builds a call graph for optimization decisions. Direct calls create edges. Indirect calls (through function pointers) create edges to all possible targets.

Escape analysis across modules: the compiler determines whether values allocated in one module escape to another. Non-escaping values can be stack-allocated.

Whole-program constant propagation: constants defined in one module are propagated to uses in other modules. This enables: configuration-driven optimization and compile-time feature flags.

Module interface extraction: each module exports a summary of its public types and functions. The summary enables separate compilation: clients compile against the summary, not the source.

Link-time code generation: instead of generating machine code during compilation, the compiler generates an intermediate representation. The linker converts IR to machine code with cross-module optimization.

Whole-program analysis limitations: analysis time grows with program size. For programs exceeding 100,000 lines, whole-program analysis is optional. Per-module analysis is the default.

Incremental whole-program analysis: the compiler caches analysis results and invalidates only affected results when source changes. Incremental analysis reduces rebuild time for large programs.

6.4 Lifetime Elision Rules

Lifetime elision reduces annotation burden. The compiler applies three rules: (1) each reference parameter gets its own lifetime, (2) if there is exactly one input lifetime, it is assigned to all output references, (3) if one parameter is `&self` or `&mut self`, the self lifetime is assigned to all output references.

Elision covers 95% of function signatures. Only functions with multiple input lifetimes that return references require explicit lifetime annotations.

Lifetime annotation syntax: `'fn longest[a](x: &'a str, y: &'a str) -> &'a str'` declares lifetime 'a shared

between inputs and output.

Lifetime bounds: 'where 'a: 'b' means lifetime 'a outlives lifetime 'b. Bounds are needed when a reference stores another reference with a different lifetime.

Static lifetime: '&static references live for the entire program. String literals have static lifetime. Global constants have static lifetime.

Lifetime in structures: a structure containing a reference must declare the lifetime: 'struct Parser['a] { input: &'a str }'. The structure cannot outlive the referenced data.

Lifetime variance: references are covariant in their lifetime parameter. A reference with a longer lifetime can be used where a shorter lifetime is expected.

Lifetime splitting: the borrow checker can split a mutable reference into non-overlapping field references. Each field reference has an independent lifetime.

Higher-rank lifetimes: 'fn apply(f: fn(&i32) -> &i32, x: &i32) -> &i32' uses higher-rank lifetimes. The function f works for any lifetime of its argument.

Lifetime debugging: the compiler can print inferred lifetimes for debugging. The output shows: each reference's lifetime, the constraints between lifetimes, and the constraint solution.

Lifetime error visualization: error messages include diagrams showing the overlapping regions of conflicting lifetimes. Visual representation helps developers understand complex lifetime interactions.

Lifetime and closures: closures that capture references have lifetimes constrained by the captured references. A closure cannot outlive any reference it captures.

10.4 Development Tooling

Language server protocol (LSP) implementation: the compiler provides an LSP server for IDE integration. Features include: code completion, go-to-definition, find references, and inline diagnostics.

Formatter: the lateralus-fmt tool formats source code according to the official style guide. Formatting is deterministic: running the formatter twice produces the same output.

Linter: the lateralus-lint tool checks for common mistakes and style violations. Lint rules include: unused variables, unreachable code, shadowed names, and naming conventions.

Package manager: the lateralus-pkg tool manages dependencies. Features include: version resolution, dependency locking, and registry publishing. Packages are distributed as source code.

REPL: the lateralus-repl provides interactive evaluation. The REPL compiles each expression to native code and executes it. The compilation overhead is amortized by caching compiled expressions.

Debugger integration: the compiler generates DWARF debug information. GDB and LLDB can debug

Lateralus programs with: breakpoints, variable inspection, stack traces, and stepping.

Profiler integration: the compiler generates frame pointer metadata for sampling profilers (perf, Instruments). Call graph profiling identifies hot functions.

Documentation generator: lateralus-doc extracts documentation from source comments and generates HTML. Cross-references between types, functions, and modules are linked automatically.

Build system: lateralus-build manages multi-module projects. Features include: dependency resolution, parallel compilation, incremental building, and test execution.

Playground: a web-based Lateralus playground compiles and runs code in the browser via WebAssembly. The playground enables: sharing code snippets, interactive tutorials, and documentation examples.

Editor plugins: official plugins exist for VS Code, Vim/Neovim, and Emacs. Plugins provide: syntax highlighting, LSP integration, snippet expansion, and build integration.

Continuous integration templates: pre-built CI configurations for GitHub Actions, GitLab CI, and Jenkins. Templates include: build, test, lint, format check, and documentation generation.

References

- [1] Appel, A. Modern Compiler Implementation in ML. Cambridge, 1998.
- [2] Nystrom, R. Crafting Interpreters. Genever Benning, 2021.
- [3] Aho, A. et al. Compilers: Principles, Techniques, and Tools. Addison-Wesley, 2006.
- [4] Wirth, N. Compiler Construction. Addison-Wesley, 1996.
- [5] Cooper, K. and Torczon, L. Engineering a Compiler. Morgan Kaufmann, 2011.