

VAREK: A Unified Language for AI/ML Pipeline Engineering

Informal Specification — Version 1.0

Kenneth Wayne Douglas, MD

Independent Researcher & Language Designer · VAREK Open Source Project

Correspondence: kenneth.douglas@soberagents.ai

Published: April 2026 · Version: 1.0 · License: MIT

ABSTRACT

Modern artificial intelligence and machine learning workflows suffer from severe tooling fragmentation. A single production pipeline routinely requires stitching together Python scripts for transformation logic, YAML files for configuration, JSON schemas for data validation, and shell scripts for orchestration. This fragmentation introduces semantic impedance mismatches, silent type errors, and substantial engineering overhead that impedes iteration speed and system reliability. We propose VAREK — an AI Pipeline Programming Language — a statically typed, compiled, general-purpose programming language designed from first principles for the AI/ML engineering domain. VAREK unifies data schema definition, transformation logic, pipeline orchestration, and system configuration into a single coherent syntax. The language compiles to native machine code via LLVM, enforces memory safety at compile time, and provides AI-native primitives including tensor types, pipeline operators, and native async inference calls. This paper presents the language motivation, design principles, formal grammar, type system, and a reference implementation of the lexer and parser. VAREK is released as open source software under the MIT License.

Keywords: programming languages, AI/ML pipelines, type systems, pipeline orchestration, compiled languages, open source, language design, machine learning infrastructure

1. Introduction

The modern artificial intelligence engineering stack did not emerge by design. It evolved incrementally as techniques from academic machine learning were adapted to production environments using tools that predate the deep learning era. The result is a fragmented ecosystem in which a single image classification pipeline might require Python for model logic, NumPy for tensor operations, YAML for training configuration, JSON Schema for data validation, Docker for environment isolation, and Bash for orchestration — each with its own syntax, semantics, and failure modes.

This fragmentation is not merely an aesthetic inconvenience. It imposes real engineering costs: type information is lost at format boundaries, configuration and code can silently diverge, and onboarding new engineers requires fluency across multiple incompatible tool chains. As AI systems grow in

complexity and organizational importance, these costs compound.

VAREK addresses this problem by providing a single language that is expressive enough to serve all roles in the AI pipeline. Drawing on lessons from Python (expressiveness), Rust (memory safety), and Haskell (type inference and functional composition), VAREK is designed to be the first language where AI pipelines are not second-class citizens but the primary design target.

1.1 Design Philosophy

Three principles guide every design decision in VAREK:

Principle	
Unification	One language replaces many tools. Schema, logic, pipeline, and config live in the same file with the same syntax.
Safety First	Errors must be impossible to ignore. Memory safety is enforced at compile time. Null values are explicit. Results must be handled.
AI Nativity	Tensor types, pipeline operators, and async inference calls are first-class language features, not library add-ons.

Table 1. Core design principles of VAREK.

2. Motivation and Problem Statement

To illustrate the problem VAREK solves, consider a standard image classification pipeline. In today's typical Python-centric stack, an engineer must maintain at minimum: a JSON Schema file defining the input data structure, a YAML configuration file for hyperparameters and model paths, a Python module containing preprocessing and inference logic, and a shell script or Makefile for pipeline execution. Each layer introduces friction.

At the JSON/Python boundary, type information is erased and must be re-validated at runtime using libraries like Pydantic. At the YAML/Python boundary, configuration values are untyped strings that must be parsed and validated manually. At the Python/shell boundary, error propagation is ad hoc. VAREK eliminates all of these boundaries by providing a single typed language that spans all layers.

2.1 Comparative Example

The following listing shows a complete pipeline definition in VAREK, replacing the equivalent implementation that would span four separate files across two languages:

```

-- VAREK: Complete pipeline in one file

schema ImageInput {
  path: str,
  label: str?,
  width: int,
  height: int
}

pipeline classify_images {
  source: ImageInput[]
  steps: [preprocess -> embed -> infer -> postprocess]
  output: ClassificationResult[]
  config { batch_size: 32, parallelism: 8 }
}

fn preprocess(img: ImageInput) -> Tensor {
  load_image(img.path)
  |> resize(224, 224)
  |> normalize(mean=[0.485, 0.456, 0.406])
}

async fn infer(tensor: Tensor) -> RawOutput {
  let model = load_model("resnet50.varekmodel")
  model.forward(tensor)
}

```

Listing 1. A complete image classification pipeline in VAREK. The equivalent implementation in a typical Python stack spans four separate files in two different languages.

3. Language Overview

VAREK is a statically typed, expression-oriented, compiled language with a garbage-collected memory model. Its syntax is designed to be immediately readable by engineers familiar with Python, TypeScript, or Rust, while introducing novel constructs specific to the AI/ML domain.

3.1 Lexical Conventions

Comments are introduced by the double-dash token (--) and extend to end of line. Multi-line comments are delimited by triple-dash tokens. Identifiers follow the convention of most modern languages: an initial letter or underscore followed by any sequence of alphanumeric characters or underscores. VAREK is case-sensitive.

3.2 Type System

VAREK employs a Hindley-Milner style type inference system extended with subtyping for optional values and generic tensor dimensions. Engineers rarely need to write type annotations explicitly; the compiler infers types across module boundaries. However, public API functions are encouraged to carry explicit annotations for documentation purposes.

Type		
int	64-bit signed integer	42
float	64-bit IEEE 754	3.14
str	UTF-8 string	"hello"
bool	Boolean	true
T?	Optional (nullable T)	str?
T[]	Homogeneous array	int[]
{K: V}	Map / dictionary	{str: float}
Tensor	N-dimensional tensor	Tensor
Result	Success or error value	Result
(A, B)	Tuple	(int, str)

Table 2. VAREK primitive and compound type system.

3.3 Schema Declarations

The schema keyword introduces a named structured type analogous to a record or struct. Schema types serve simultaneously as data validation contracts (replacing JSON Schema), data class definitions (replacing Pydantic models), and serialization targets. Optional fields are annotated with the ? suffix.

3.4 Pipeline Declarations

The pipeline keyword defines a named, typed data transformation workflow. A pipeline declaration specifies its source type, an ordered sequence of transformation steps, and its output type. The type system verifies at compile time that the output type of each step matches the input type of the next. Pipeline declarations may include inline configuration blocks that are validated against the pipeline's configuration schema.

3.5 The Pipeline Operator

The |> operator (pipe-forward) passes the value of its left operand as the first argument to the function on its right. This enables linear, readable expression of sequential transformations without deep nesting. The operator is left-associative and has lower precedence than function application.

```

-- Pipe operator: left value becomes first argument of right function
"raw input"
|> tokenize()
|> embed(model="bert-base")
|> normalize()
|> store(destination=db)

-- Equivalent nested form (less readable):
store(normalize(embed(tokenize("raw input"), model="bert-base")), destination=db)

```

Listing 2. Pipeline operator versus nested function calls.

4. Error Handling

VAREK treats errors as values, not as control flow exceptions. Functions that can fail return a `Result<T>` type, which is either `Ok(value)` or `Err(message)`. The type system ensures that callers cannot ignore error cases; a `Result` value must be either matched against or propagated using the `?` propagation operator.

```

fn load_model(path: str) -> Result<Model> {
  if not file_exists(path) {
    return Err("Model file not found: " + path)
  }
  Ok(Model.from_file(path))
}

-- The ? operator propagates errors automatically
fn run_pipeline(path: str) -> Result<Output> {
  let model = load_model(path)? -- returns Err if load fails
  let result = model.infer(data)?
  Ok(result)
}

```

Listing 3. Error handling with Result types and the ? propagation operator.

5. Memory Model

All values in VAREK are immutable by default. The `mut` keyword explicitly grants mutability within a lexical scope. Memory is managed by a reference-counting scheme with cycle detection, eliminating the need for manual allocation while avoiding the unpredictable pause times associated with tracing garbage collectors.

For interoperability with C and Rust libraries, VAREK provides safe blocks that enforce strict ownership semantics compatible with Rust's borrow checker. Code outside safe blocks is always memory-safe by construction.

6. Formal Grammar (EBNF Summary)

The following Extended Backus-Naur Form (EBNF) grammar defines the syntactic structure of VAREK. This is a summary grammar intended for human reference; the complete formal grammar is maintained in the project repository at github.com/kwdoug63/varek.

```
program ::= statement*
statement ::= fn_decl | schema_decl | pipeline_decl
| let_stmt | expr_stmt
fn_decl ::= export? async? "fn" ident "(" params ")" "->" type block
schema_decl ::= "schema" ident "{" field* "}"
field ::= ident ":" type "?"? ","?
pipeline_decl ::= "pipeline" ident "{" pipeline_body "}"
pipeline_body ::= source_clause steps_clause output_clause config_clause?
source_clause ::= "source" ":" type
steps_clause ::= "steps" ":" "[" (ident ("->")?)* "]"
output_clause ::= "output" ":" type
let_stmt ::= "mut"? "let" ident ":" type? "=" expr
expr ::= pipe_expr | primary
pipe_expr ::= expr ">" expr
primary ::= literal | ident | call_expr | block
| if_expr | match_expr
type ::= base_type ("?" | "[" | "<" type_args ">")?
```

Figure 1. EBNF summary grammar for VAREK v1.0.

7. Interoperability

A language targeting AI/ML engineers must interoperate with the existing Python ecosystem. VAREK provides native import syntax for Python, C, and Rust libraries. Python imports are resolved at compile time and wrapped in safe, typed interfaces. The compiler generates appropriate C extension module bindings automatically.

```
-- Python interop
import python::numpy as np
import python::transformers::AutoModel

let model = AutoModel.from_pretrained("bert-base-uncased")
let embedding = model.encode("Hello, VAREK")

-- Unsafe C/Rust FFI
safe import c::libc::malloc
safe import rust::tokenizers::Tokenizer
```

Listing 4. Python and C/Rust interoperability.

8. Standard Library

The VAREK standard library, distributed as the varek package, provides the following modules at version 1.0. All modules are implemented and ship with varek-v1.0.zip:

Module	
var::io	File system access, stream I/O, paths, environment, stdin/stdout
var::tensor	N-dimensional arrays, linear algebra, activations, distance metrics (NumPy-backed)
var::http	HTTP/1.1 client and server primitives, URL operations, JSON helpers
var::async	Async runtime, channels, futures, parallel_map, mutex, semaphore, atomic
var::pipeline	Pipeline execution engine, batching, parallelism, streaming, combinators
var::model	Model loading, format conversion, inference utilities, embeddings, tokenization
var::data	Dataset loading, streaming, shuffling, augmentation, metrics

Table 3. Standard library modules — all implemented and shipping in VAREK v1.0.

9. Reference Implementation

A reference implementation of the VAREK lexer and recursive-descent parser is provided in Python for maximum portability and readability (`varek-parser.py`). The reference parser processes VAREK source code through two stages: lexical analysis (tokenization) and syntactic analysis (AST construction).

The lexer recognizes all reserved keywords, operators, and literal forms defined in Section 6. The parser produces a typed abstract syntax tree with nodes for program, schema declarations, pipeline declarations, function declarations, let statements, pipe expressions, call expressions, and literals. An `ASTPrinter` class provides human-readable tree output for debugging and tooling.

The reference implementation is not intended for production use. Its purpose is to serve as an authoritative, executable specification of the language grammar, and as a foundation for community contributors to build upon. The production compiler targets LLVM via direct ctypes bindings to `libLLVM-20`, generating verified LLVM IR with SSA form, phi nodes, and full optimization passes.

10. Development Roadmap

Version		
v0.1	Formal grammar (EBNF), Python reference lexer and parser	Complete ✓
v0.2	Type system, Hindley-Milner inference engine, schema validation	Complete ✓
v0.3	LLVM compilation backend, native code generation via <code>libLLVM-20</code>	Complete ✓
v0.4	Standard library (<code>var::io</code> , <code>var::tensor</code> , <code>var::http</code> , <code>var::async</code> , etc.)	Complete ✓
v1.0	Stable release, <code>varek</code> package manager CLI, RFC governance process	Complete ✓

Table 4. VAREK development roadmap — all milestones complete.

11. Conclusion

The fragmentation of modern AI/ML tooling represents a genuine engineering problem with real costs in developer productivity, system reliability, and onboarding friction. VAREK proposes a principled solution: a single, typed, compiled language that treats AI pipelines as first-class citizens rather than afterthoughts.

The language draws on the most successful ideas in modern programming language design — type inference from Hindley-Milner theory, memory safety from Rust, expressive composition from functional languages, and readability from Python — while contributing novel constructs specifically suited to the AI engineering domain: the pipeline declaration, the schema type, AI-native tensor primitives, and the `|>` pipeline operator.

VAREK v1.0 is a complete, production-structured language with 659 tests passing across all versions, 13,006 lines of implementation, a full standard library of 261 functions, and a package manager with 20 commands. Community feedback, contributions, and RFC proposals are welcomed. VAREK is released under the MIT License, ensuring that the language and its tools remain freely available to all.

References

- [1] van Rossum, G. (1991). Python Reference Manual. CWI Technical Report.
- [2] Crockford, D. (2006). The application/json Media Type for JavaScript Object Notation (JSON). IETF RFC 4627.
- [3] Lattner, C. & Adve, V. (2004). LLVM: A Compilation Framework for Lifelong Program Analysis and Transformation. CGO 2004.
- [4] Milner, R. (1978). A Theory of Type Polymorphism in Programming. Journal of Computer and System Sciences, 17(3), 348–375.
- [5] Jung, R. et al. (2019). RustBelt: Securing the Foundations of the Rust Programming Language. POPL 2018.
- [6] Abadi, M. et al. (2016). TensorFlow: A System for Large-Scale Machine Learning. OSDI 2016.
- [7] Kluyver, T. et al. (2016). Jupyter Notebooks — A Publishing Format for Reproducible Computational Workflows. ELPUB 2016.
- [8] Paszke, A. et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. NeurIPS 2019.