

# MLIR: A Panacea for ML Compiler Challenges?

Vishakha Agrawal vishakha.research.id@gmail.com

**Abstract**—MLIR (Multi-Level Intermediate Representation) has emerged as a promising solution to address many of the challenges faced in AI/ML compiler development. By providing a flexible and extensible intermediate representation, MLIR enables the decoupling of frontend and backend compiler components, facilitating greater modularity and reusability. However, while MLIR offers significant advantages, it is not a panacea for all compiler-related issues in the AI/ML domain. Several challenges persist, including the need for efficient memory management, optimal kernel fusion etc. MLIR's ability to integrate with existing compiler frameworks and support diverse AI/ML workloads also remains an open question. This paper critically examines MLIR's strengths and limitations in addressing AI/ML compiler challenges, providing a comprehensive analysis of its potential.

**Keywords** - MLIR, LLVM, memory management, kernel fusion, GPU, accelerators, code generation, XLA, ODS

## I. INTRODUCTION

The ML compiler landscape is characterized by a multitude of frameworks (e.g., TensorFlow, PyTorch, ONNX), compilers (e.g., TensorFlow XLA, Nvidia TensorRT), and hardware platforms (e.g., CPUs, GPUs, TPUs). This diversity leads to several challenges. Compilation complexity is a significant issue, as each framework and hardware platform requires customized compilation, resulting in duplicated effort and increased maintenance costs. Additionally, performance optimization challenges arise when optimizing AI/ML models for diverse hardware platforms, often requiring manual tuning and expertise. Furthermore, limited portability is a concern, as AI/ML models are often tightly coupled with specific frameworks and hardware platforms, making it difficult to deploy them on different environments. Inefficient resource utilization is also a problem, as the lack of standardized compilation and optimization techniques can lead to suboptimal resource utilization, resulting in wasted computational resources and energy. Finally, integrating new hardware platforms or architectures can be cumbersome, requiring significant modifications to existing compilation and optimization pipelines. These challenges underscore the need for a unified and flexible compilation framework that can efficiently target diverse AI/ML workloads and hardware platforms, ultimately enabling more widespread adoption and innovation in the field. While MLIR offers a promising solution to these challenges, its adoption is not without complexities. A deeper exploration of MLIR reveals additional considerations that must be carefully navigated to fully harness its potential.

## II. MLIR: A UNIFIED INTERMEDIATE REPRESENTATION

MLIR (Multi-Level Intermediate Representation) offers a flexible and extensible framework for representing computations at multiple levels of abstraction[7]. This unified approach addresses the fragmentation prevalent in AI/ML compiler stacks, where proprietary solutions hinder interoperability. MLIR shares similarities with LLVM[5] as a compiler infrastructure, but whereas LLVM excels in scalar optimizations and homogeneous compilation, MLIR expands its focus to incorporate a broader range of data structures and algorithms. Specifically, MLIR treats tensor algebra, graph representations, and heterogeneous compilation as first-class citizens, enabling a more comprehensive and flexible approach to compiler design. Additionally, MLIR's modular architecture is designed to seamlessly integrate with existing frameworks, such as the polyhedral framework ([4], [1]) which provides a powerful basis for analyzing and optimizing complex loop structures. While MLIR aims to provide a unified infrastructure for ML compilers, other approaches like Tensor Comprehensions focus on domain-specific languages and polyhedral compilation to automatically generate optimized GPU kernels for tensor operations[8].

### A. Key Features: MLIR's key features as described in [6]

- 1) **Multi-level Design** : At its core, MLIR's most distinctive feature is its multi-level design, which supports multiple abstraction levels within the same IR. This unique approach enables seamless representation from high-level operations down to low-level instructions, allowing for gradual and selective lowering of operations as needed and solves the problem highlighted in [2] of transparent multi-level view.
- 2) **Strong Type System** : This flexibility is complemented by MLIR's strong type system, which provides robust and extensible support for both primitive and complex data types. The type system enables compile-time checking and verification, while also accommodating dialect-specific custom types.
- 3) **SSA-based Representation** : MLIR employs a Static Single Assignment (SSA)[3] based representation, where each value is defined exactly once. This design choice simplifies data flow analysis and optimizations while maintaining explicit def-use relationships throughout the code.
- 4) **Dialect Framework** : The dialect framework is cornerstone of MLIR's architecture, providing mechanisms to define custom operations and types. This framework enables domain-specific optimizations and allows different levels of abstraction to coexist harmoniously, with

examples including the Standard dialect, LLVM dialect, and GPU dialect.

- 5) Operation Definition Specification (ODS) : The Operation Definition Specification (ODS) system in MLIR offers declarative specification of operations with automated code generation. This reduces boilerplate code and potential errors while providing built-in verification of operation properties.
- 6) Region and Block Structure : The region and block structure of MLIR supports nested regions within operations, enabling natural representation of control flow, both structured and unstructured, and facilitating complex nested computations.
- 7) Generic Optimization Framework : MLIR's generic optimization framework includes a pattern-based rewriting infrastructure, dialect conversion framework, and common optimization passes, all managed by an extensible pass management system.
- 8) Integration capabilities : MLIR is designed to work seamlessly with existing compilers, bridging different compilation systems and supporting various frontend and backend technologies, particularly through its integration with LLVM.
- 9) Open Source and Community-Driven : The open-source nature of MLIR, being part of the LLVM project, has fostered an active development community and a growing ecosystem of tools and dialects. This is supported by comprehensive documentation and continuous community contributions.
- 10) Analysis and Transformation Tools : The analysis and transformation tools include built-in debugging capabilities, visualization tools for IR inspection, and an extensive testing framework.

#### B. Strength: MLIR's strengths includes

- 1) Unification: Providing a unified representation for AI/ML models, simplifying compilation and optimization.
- 2) Modularity: Enabling the development of specialized dialects for specific AI/ML domains.
- 3) Community: Boasting a growing community of developers and users.
- 4) Heterogeneity: While various attempts have been made for to address hardware heterogeneity, MLIR allows high-level languages to harness heterogeneity through extensible operations and types, while providing a shared infrastructure for transforming these constructs into optimized, target-specific code.

### III. LIMITATIONS AND CHALLENGES: A DEEP TECHNICAL ANALYSIS

The Multi-Level Intermediate Representation (MLIR) framework, while powerful and innovative, faces several significant challenges and limitations that impact its adoption and effectiveness. This analysis explores these challenges in depth,

examining their technical implications and potential impact on the broader compiler ecosystem.

#### 1) Fundamental Architectural Challenges:

- One of the most pressing challenges in MLIR centers around dialect interoperability. The framework's ability to support multiple dialects, while powerful, has led to what many practitioners call "dialect explosion." This phenomenon occurs when projects create numerous custom dialects, leading to a complex web of translations and interactions. The maintenance overhead becomes substantial as teams must manage compatibility between these dialects, often requiring intricate conversion patterns and careful consideration of semantic preservation.
- The translation complexity between dialects presents another significant challenge. When converting between different abstraction levels, developers must carefully manage information preservation while maintaining performance. This becomes particularly challenging when dealing with dialects that operate at significantly different abstraction levels or have fundamentally different semantic models. The complexity of these translations can lead to performance bottlenecks and potential loss of important program information during the conversion process.

#### 2) Type System and Operation Semantics:

- MLIR's type system, while flexible, faces limitations in expressing complex type relationships across dialects. Developers often struggle to represent sophisticated type constraints that span multiple dialects or involve complex interactions between different levels of abstraction. This becomes particularly evident when working with domain-specific types that don't map cleanly to standard MLIR types.
- Operation semantics present another significant challenge. The framework struggles to fully capture and preserve complex operation semantics during transformations, particularly when dealing with side effects or intricate control flow patterns. This limitation can make it difficult to implement certain optimizations or ensure behavioral consistency across transformations.

#### 3) Implementation and Performance Challenges:

- The compilation pipeline in MLIR faces several significant technical challenges. Pass management becomes increasingly complex as projects grow, with developers struggling to maintain optimal pass ordering while ensuring all necessary transformations are applied correctly. The inter-dependencies between passes can create subtle bugs and per-

formance issues that are difficult to diagnose and resolve.

- Memory usage represents another critical challenge. MLIR's intermediate representation can consume significant memory, particularly when dealing with large programs or complex transformation sequences. This becomes especially problematic during dialect conversion and optimization passes, where multiple versions of the IR might need to be maintained simultaneously. The memory overhead can become a bottleneck in large-scale applications.

#### 4) Development and Tooling Ecosystem:

- The development experience with MLIR presents its own set of challenges. The learning curve is steep, with developers needing to understand not only the core concepts but also the intricacies of dialect design and transformation implementation. The API for defining new operations and transformations often requires significant boilerplate code, which can be error-prone and time-consuming to maintain.
- Debugging support in MLIR, while improving, remains limited compared to more mature compiler frameworks. Developers often struggle to track the effects of transformations and understand error messages, particularly when dealing with complex transformation sequences or dialect conversions. The lack of sophisticated debugging tools can significantly slow down development and optimization efforts.

#### 5) Tool Integration and Production Readiness:

- Integration with existing compiler frameworks and tools presents significant challenges. While MLIR is designed to be flexible, incorporating it into established compilation pipelines can require substantial engineering effort. Performance overhead in framework bridges and version compatibility issues can make integration particularly challenging for production environments.
- Production readiness remains a concern for many potential adopters. API stability issues can make it difficult to maintain long-term compatibility, while performance monitoring and profiling tools are still maturing. The lack of standardized benchmarks makes it challenging to assess the impact of changes and optimizations reliably.

#### 6) Future Challenges and Scaling Concerns:

- As machine learning models continue to grow in size and complexity, MLIR faces increasing challenges in handling very large IR representations efficiently. The framework must evolve to support distributed compilation scenarios and manage memory more efficiently for big models. These scaling challenges become particularly acute when dealing with modern transformer architectures and their billions of parameters.

- Support for emerging hardware architectures presents another ongoing challenge. As new accelerators and specialized processors emerge, MLIR must adapt to effectively target these platforms while maintaining performance portability. This requires continuous evolution of the framework's hardware abstraction capabilities and optimization strategies.

By acknowledging both the benefits and limitations of MLIR, researchers and developers can better understand its potential as a unified compiler infrastructure for AI/ML and work towards addressing its challenges.

## IV. POTENTIAL SOLUTIONS

While MLIR faces significant challenges in various areas, many potential solutions exist or are being developed. This analysis presents practical and theoretical solutions to the major challenges identified in MLIR implementation and usage, along with potential road maps for future improvements.

- The challenge of dialect explosion and interoperability can be addressed through several systematic approaches. Implementing a standardized dialect interface framework would allow different dialects to communicate through well-defined protocols. This could be achieved by developing a common abstraction layer that serves as a bridge between dialects.
- The development experience can be improved through better tooling, implementing stability frameworks and performance monitoring system.
- To address scaling challenges, distributed compilation systems can be implemented.
- There is also room for AI-Assisted Optimization within the MLIR framework to enable more advanced and automated optimizations, essentially creating an AI-assisted optimization capability within the compiler infrastructure. This will allow for better analysis of complex operations and generation of highly optimized code for diverse hardware platforms.

## V. CONCLUSION

MLIR has the potential to address many of the AI/ML compiler issues by providing a unified, modular, and extensible intermediate representation. However, its adoption and maturity are still ongoing processes. While MLIR is not a silver bullet that solves all AI/ML compiler issues, it is a significant step towards unifying the compilation process and improving performance optimization.

## REFERENCES

- [1] Riyadh Baghdadi, Jessica Ray, Malek Ben Romdhane, Emanuele Del Sozzo, Abdurrahman Akkas, Yunming Zhang, Patricia Suriana, Shoaib Kamil, and Saman Amarasinghe. Tiramisu: A polyhedral compiler for expressing fast and portable code. In *2019 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 193–205. IEEE, 2019.
- [2] Gogul Balakrishnan and Thomas Reps. Wysinwyx: What you see is not what you execute. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 32(6):1–84, 2010.

- [3] Ron Cytron, Jeanne Ferrante, Barry K Rosen, Mark N Wegman, and F Kenneth Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 13(4):451–490, 1991.
- [4] Tobias Grosser, Armin Groesslinger, and Christian Lengauer. Polly—performing polyhedral optimizations on a low-level intermediate representation. *Parallel Processing Letters*, 22(04):1250010, 2012.
- [5] Chris Lattner and Vikram Adve. Llvm: A compilation framework for lifelong program analysis & transformation. In *International symposium on code generation and optimization, 2004. CGO 2004.*, pages 75–86. IEEE, 2004.
- [6] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. Mlir: A compiler infrastructure for the end of moore’s law. *arXiv preprint arXiv:2002.11054*, 2020.
- [7] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Arnaud Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. Mlir: Scaling compiler infrastructure for domain specific computation. In *CGO 2021*, 2021.
- [8] Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal, Zachary Devito, William S Moses, Sven Verdoolaege, Andrew Adams, and Albert Cohen. The next 700 accelerated layers: From mathematical expressions of network computation graphs to accelerated gpu kernels, automatically. *ACM Transactions on Architecture and Code Optimization (TACO)*, 16(4):1–26, 2019.