# Learn Compiler Design Through xDSL:
# A Practical Course for AI/ML Engineers

Djordje Todorovic

August 2025

# Contents

# 1   Preface: Who This Course Is For and Why Compilers Matter

## 1.1   Prerequisites and Intended Audience

This course is designed for engineers and developers who want to understand the fundamental concepts of compiler design through hands-on experience with xDSL. To get the most out of this material, you should have:

- **Python Proficiency**: Strong knowledge of Python is essential, as xDSL is a Python-native framework. You should be comfortable with classes, decorators, type hints, and the Python standard library.

- **Basic Programming Knowledge**: Understanding of fundamental programming concepts like functions, variables, control flow, data structures, and algorithms is required.

- **Machine Learning and AI Basics**: Familiarity with ML/AI concepts is crucial - you should understand tensors, neural networks, computational graphs, and optimization techniques. Experience with frameworks like PyTorch or TensorFlow will help you appreciate the compiler optimizations we'll explore.

- **Mathematical Foundations**: Basic understanding of linear algebra and discrete mathematics will be helpful, especially when working with optimization algorithms and graph transformations.

If you're an ML engineer curious about what happens "under the hood" when your models are compiled and optimized, or a systems programmer interested in building efficient code transformation tools, this course will bridge that gap using familiar Python syntax.

## 1.2   Compilers: The Invisible Revolution That Changed Computing

Compilers are perhaps the most transformative technology in the history of computing, yet they work so seamlessly that we rarely think about them. To understand their revolutionary impact, consider this: the Linux kernel, which powers billions of devices worldwide, was initially written in assembly language - a tedious, error-prone process where programmers had to think in terms of individual CPU instructions and memory addresses. Each line of assembly code corresponded directly to a machine instruction, making even simple tasks require hundreds of lines of code.

Then came the evolution of the C compiler. As compilers matured, they could translate high-level C code into assembly with equal or sometimes even better performance than hand-written assembly. This transformation was revolutionary - Linux could be rewritten in C, making it portable across different architectures, easier to maintain, and accessible to a broader community of developers. What once required intimate knowledge of specific CPU architectures could now be expressed in readable, maintainable code. The compiler handled the complex translation, optimization, and architecture-specific details automatically.

This same revolution continues today in machine learning. Just as C compilers freed developers from assembly, ML compilers like XLA, TVM, and MLIR free ML engineers from writing architecture-specific kernels. Your PyTorch model can run efficiently on CPUs, GPUs, TPUs, or custom accelerators, all thanks to sophisticated compiler technology working behind the scenes.

## 1.3   The Core Mission: Represent and Transform

At its heart, a compiler's job is elegantly simple yet profoundly powerful: to **represent** programs in structured forms and to **transform** these representations to achieve specific goals. Think of a compiler as a sophisticated translator that not only converts between languages but also understands the meaning of what it's translating deeply enough to improve it along the way.

This representation and transformation paradigm involves:

- **Representation**: Converting source code into structured data (Abstract Syntax Trees, Intermediate Representations, Control Flow Graphs) that capture the program's semantics precisely

- **Analysis**: Understanding data dependencies, control flow, memory access patterns, and optimization opportunities within these representations

- **Transformation**: Systematically modifying these representations to optimize performance, reduce resource usage, or target specific hardware architectures

- **Preservation**: Ensuring that transformations maintain the original program's correctness and semantics

## 1.4  The Power of Serialize, Deserialize, and Verify

A crucial but often overlooked aspect of compiler infrastructure is the ability to serialize, deserialize, and verify intermediate representations. This capability fundamentally changes how we can work with compilers and is a cornerstone of modern compiler design.

**Why These Capabilities Matter:**

- **Serialization** allows us to save the compiler's intermediate state to disk, enabling separate compilation, caching of optimization results, and distribution of compiled modules. You can pause compilation, save the IR, and resume later - or on a different machine entirely.

- **Deserialization** lets us load previously compiled modules, compose them together, and build large systems incrementally. This is essential for modular compilation and linking.

- **Verification** ensures that the IR is well-formed and obeys the type system and semantic rules. This catches errors early, enables safe transformations, and provides guarantees about program behavior. Without verification, compiler bugs could silently corrupt programs.

The xDSL infrastructure we're about to explore has these capabilities built into its core. Every operation, every transformation, and every intermediate state can be serialized to a human-readable textual format, loaded back, and verified for correctness. This isn't just a convenience feature - it's fundamental to building robust, composable compiler pipelines. You can inspect the IR at any stage, debug transformations by examining the before-and-after states, and even hand-write IR for testing.

This approach contrasts sharply with traditional compilers where the intermediate states are often opaque binary structures locked inside the compiler's memory. With xDSL (inspired by MLIR), the entire compilation pipeline becomes transparent, debuggable, and extensible. You can serialize the IR after each optimization pass, analyze what changed, and even replay specific transformations. This visibility and control make compiler development more like software engineering and less like black magic.

As you progress through this course, you'll come to appreciate how these fundamental capabilities - serialize, deserialize, and verify - enable you to build reliable, maintainable, and powerful compilation tools. They transform the compiler from a monolithic black box into a modular, inspectable, and trustworthy system.

## 2 Introduction: Compilers as Representation and Transformation Engines

### 2.1 Welcome to the World of Compilers

If you're coming from an AI/ML background, you already understand the power of transforming data through layers of abstraction. Compilers are remarkably similar: they take programs written in one representation and systematically transform them into another, more optimized or more executable form.

Think of a compiler as a sophisticated pipeline that:

1. **Represents** programs as structured data (Abstract Syntax Trees, Intermediate Representations)

2. **Transforms** these representations through a series of optimization passes

3. **Analyzes** code to understand dependencies, patterns, and optimization opportunities

4. **Generates** efficient target code for specific hardware

### 2.2 Why xDSL?

xDSL is a Python-native compiler framework that makes compiler concepts accessible to Python programmers. Unlike traditional compiler frameworks written in C++ (like LLVM), xDSL allows you to:

- Build compilers using familiar Python syntax

- Prototype and experiment quickly

- Understand compiler internals without low-level complexity

- Leverage the entire Python ecosystem for analysis and visualization

### 2.3 Course Philosophy: Learn by Building

This course takes a hands-on approach. Instead of starting with theory, we'll build small compilers and optimization passes from day one. Each concept will be introduced through practical examples that you can run, modify, and experiment with.

### 2.4 What You'll Learn

By the end of this course, you'll understand:

- How compilers represent programs internally (IR - Intermediate Representations)

- What operations and operators mean in a compiler context

- How to write analysis passes that extract information from code

- How to implement transformations that optimize programs

- The connection between compiler optimizations and ML model optimization

# 3 Chapter 1: Understanding Intermediate Representations (IR)

## 3.1 What is an IR?

An Intermediate Representation (IR) is the compiler's internal language for representing programs. Just as neural networks use tensors to represent data, compilers use IRs to represent code structure and semantics.

Key properties of IRs:

- **Abstract**: Higher-level than machine code, but more structured than source code

- **Explicit**: Makes implicit operations visible (e.g., memory allocations, type conversions)

- **Analyzable**: Designed for easy pattern matching and transformation

- **Hierarchical**: Can represent nested structures (functions, loops, blocks)

## 3.2 The Compiler Pipeline: From Source to Machine Code

Before diving into xDSL specifics, let's understand how compilers work. A compiler is traditionally divided into three main phases, each with specific responsibilities:

### 3.2.1 Frontend: Understanding Your Code

The frontend translates source code into an intermediate representation. It consists of:

**1. Lexical Analysis (Lexer/Tokenizer)**

- Breaks source code into tokens (keywords, identifiers, operators)

- Like splitting a sentence into words

- Example: "x = 42 + y" becomes tokens: [ID:x, ASSIGN, NUM:42, PLUS, ID:y]

**2. Syntax Analysis (Parser)**

- Builds an Abstract Syntax Tree (AST) from tokens

- Checks grammatical structure

- Like diagramming a sentence to understand its structure

```
# Source: x = 42 + y
# AST representation:
Assignment(
    left=Variable("x"),
    right=BinaryOp(
        op="+",
        left=Number(42),
        right=Variable("y")
    )
)
```
Listing 1: Simple AST Example

**3. Semantic Analysis**

- Type checking: Is this operation valid?

- Symbol resolution: What does this variable refer to?

- Example: Can't add a string to a number, variables must be declared

### 3.2.2 Middle-end: Optimization Central

The middle-end works on the IR to optimize code:

- **Platform-independent**: Optimizations work for any target

- **Examples**: Dead code elimination, constant folding, loop optimization

- **Multiple passes**: Each optimization is a separate pass over the IR

### 3.2.3 Backend: Generating Machine Code

The backend translates optimized IR to target-specific code:

- **Instruction selection**: Choose machine instructions

- **Register allocation**: Assign variables to CPU registers

- **Instruction scheduling**: Order operations for performance

## 3.3 Learning from Giants: LLVM and MLIR

### 3.3.1 LLVM: The Industry Standard

LLVM (Low Level Virtual Machine) is the most widely used compiler infrastructure:

- Powers Clang (C/C++), Swift, Rust, and many other languages

- Provides a well-defined IR (LLVM IR) that many frontends target

- Excellent tutorial: Kaleidoscope Tutorial

- Shows how to build a complete compiler for a simple language

```
1 ; LLVM IR for: int add(int a, int b) { return a + b; }
2 define i32 @add(i32 %a, i32 %b) {
3 entry:
4   %sum = add i32 %a, %b
5   ret i32 %sum
6 }
```
Listing 2: LLVM IR Example

### 3.3.2 MLIR: Multi-Level IR

MLIR (Multi-Level Intermediate Representation) extends LLVM's concepts:

- Developed by Google, now part of LLVM project

- Supports multiple IRs (dialects) in the same infrastructure

- Used by TensorFlow, PyTorch, and other ML frameworks

- Tutorial: Toy Language Tutorial

- xDSL is heavily inspired by MLIR but implemented in Python

```
1  // High-level ML operation
2  %result = "tf.MatMul"(%a, %b) : (tensor<2x3xf32>, tensor<3x4xf32>)
3                                  -> tensor<2x4xf32>
4
5  // After lowering to linalg dialect
6  %result = linalg.matmul ins(%a, %b : memref<2x3xf32>, memref<3x4xf32>)
7                          outs(%c : memref<2x4xf32>)
8
9  // After lowering to loops
10 scf.for %i = 0 to 2 {
11   scf.for %j = 0 to 4 {
12     scf.for %k = 0 to 3 {
13       // actual computation
14     }
15   }
16 }
```

Listing 3: MLIR Example - Multiple Abstraction Levels

## 3.4   Why xDSL Uses IR Instead of AST

While traditional compilers start with AST, xDSL (like MLIR) works directly with IR because:

- **Uniformity**: All transformations work on the same structure

- **Composability**: Easy to combine different languages/dialects

- **Reusability**: Optimizations work across different source languages

- **Analysis-friendly**: SSA form makes many analyses trivial

Think of it this way:

- **AST**: Like a sentence diagram - shows structure

- **IR**: Like assembly with types - shows computation

## 3.5   SSA Form: Single Static Assignment

xDSL uses SSA (Single Static Assignment) form, where each variable is assigned exactly once. This is similar to functional programming and makes many analyses simpler.

```
1  # Traditional form
2  x = 5
3  x = x + 1
4  y = x * 2
5
6  # SSA form (conceptual)
7  x_0 = 5
8  x_1 = x_0 + 1
9  y_0 = x_1 * 2
```

Listing 4: Traditional vs SSA Form

Why SSA matters:

- **No ambiguity**: Each value has exactly one definition

- **Easy analysis**: Use-def chains are explicit

- **Better optimization**: Many optimizations become simpler

- **Parallel-friendly**: Dependencies are clear
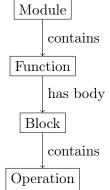
9

## 3.6 Your First xDSL Program

Let's create a simple IR program using xDSL:

```python
from xdsl.context import Context
from xdsl.dialects import builtin, arith, func
from xdsl.ir import Operation, Block, Region
from xdsl.builder import Builder

# Create a context (like a namespace for dialects)
ctx = Context()
ctx.load_dialect(builtin.Builtin)
ctx.load_dialect(arith.Arith)
ctx.load_dialect(func.Func)

# Build a simple function that adds two numbers
@Builder.implicit_region
def create_add_function():
    # Create function signature: (i32, i32) -> i32
    func_type = func.FunctionType.from_lists([builtin.i32, builtin.i32
    ], [builtin.i32])

    with func.FuncOp("add", func_type, visibility="public").body:
        # Get function arguments
        arg0 = Block.current().args[0]
        arg1 = Block.current().args[1]

        # Create add operation
        result = arith.Addi(arg0, arg1).result

        # Return the result
        func.Return(result)

module = builtin.ModuleOp([create_add_function()])
print(module)
```

Listing 5: Creating IR in xDSL

## 3.7 Understanding IR Structure

The IR in xDSL follows a hierarchical structure:

```
Module
  | contains
Function
  | has body
Block
  | contains
Operation
```

## 3.8 Dialects: Modular IR Design

xDSL organizes operations into *dialects* - collections of related operations and types. This is similar to how PyTorch has different modules (nn, optim, etc.).

Common dialects in xDSL:

- `arith`: Arithmetic operations (add, multiply, divide)

- `func`: Function definitions and calls

- `memref`: Memory operations (like NumPy arrays)

- `scf`: Structured control flow (loops, conditionals)

- `linalg`: Linear algebra operations (matrix multiply, convolution)

# 4 Chapter 2: Operations and Operators

## 4.1 What is an Operation?

An operation in xDSL is the fundamental unit of computation. Each operation:

- Has a name (e.g., `arith.addi` for integer addition)

- Takes input values (operands)

- Produces output values (results)

- May have attributes (compile-time constants)

- Can contain regions (nested blocks of operations)

## 4.2 Creating Custom Operations

Let's define a custom operation for matrix operations:

```python
from xdsl.irdl import (
    IRDLOperation,
    OpResult,
    Operand,
    OperandDef,
    ResultDef,
    op_def,
    result_def
)
from xdsl.dialects.builtin import TensorType

@op_def
class MatMulOp(IRDLOperation):
    """Matrix multiplication operation"""
    name = "my_dialect.matmul"

    # Define operands (inputs)
    lhs = operand_def(TensorType)
    rhs = operand_def(TensorType)

    # Define results (outputs)
    result = result_def(TensorType)

    def __init__(self, lhs: Value, rhs: Value):
        super().__init__(operands=[lhs, rhs],
                         result_types=[infer_matmul_type(lhs, rhs)])
```

```
28    def verify(self):
29        """Verify that the operation is well-formed"""
30        # Check that dimensions are compatible
31        lhs_shape = self.lhs.type.shape
32        rhs_shape = self.rhs.type.shape
33
34        if lhs_shape[-1] != rhs_shape[0]:
35            raise ValueError("Incompatible matrix dimensions")
```

Listing 6: Defining Custom Operations

## 4.3 Operation Semantics

Every operation has well-defined semantics - what it means to execute that operation:

```
1  # Arithmetic operations
2  add_op = arith.Addi(a, b)  # result = a + b
3  mul_op = arith.Muli(a, b)  # result = a * b
4
5  # Memory operations
6  alloc_op = memref.Alloc.get(shape=[10, 20], element_type=f32)  #
       Allocate 10x20 array
7  load_op = memref.Load.get(memref_value, indices)  # Load from memory
8  store_op = memref.Store.get(value, memref_value, indices)  # Store to
       memory
9
10 # Control flow operations
11 for_op = scf.For(lower_bound, upper_bound, step) # for loop
12 if_op = scf.If(condition, has_else=True)  # if-then-else
```

Listing 7: Operation Semantics Example

## 4.4 Attributes vs Operands

Understanding the difference between attributes and operands is crucial:

- **Attributes**: Compile-time constants (shapes, types, flags)

- **Operands**: Runtime values (variables, intermediate results)

```
1  # Constant has an attribute (the value)
2  const = arith.Constant.from_int_and_width(42, 32)   # 42 is an attribute
3
4  # Add has operands (runtime values)
5  result = arith.Addi(x, y)  # x and y are operands
6
7  # Alloc has attributes (shape) but no operands
8  mem = memref.Alloc.get(shape=[100], element_type=f32)  # shape is
       attribute
```

Listing 8: Attributes vs Operands

# 5 Chapter 3: Analysis Passes

## 5.1 What is Analysis?

Analysis passes extract information from IR without modifying it. They answer questions like:

- Which variables are used where? (Use-Def Analysis)

- Which operations can run in parallel? (Dependency Analysis)

- What are the loop bounds? (Range Analysis)

- Which operations are dead code? (Liveness Analysis)

## 5.2 Writing a Simple Analysis Pass

Let's write an analysis that counts operations by type:

```python
from xdsl.passes import Pass
from xdsl.ir import Operation, OpResult
from collections import defaultdict

class OperationCounterPass(Pass):
    """Count operations by type in the IR"""

    name = "count-ops"

    def apply(self, module: builtin.ModuleOp) -> None:
        op_counts = defaultdict(int)

        # Walk through all operations in the module
        for op in module.walk():
            op_counts[op.name] += 1

        # Print analysis results
        print("Operation counts:")
        for op_name, count in sorted(op_counts.items()):
            print(f"  {op_name}: {count}")

        # Store results for use by other passes
        self.op_counts = op_counts
```

Listing 9: Operation Counter Analysis

## 5.3 Use-Def Chains

Use-Def chains track where values are defined and used:

```python
class UseDefAnalysis(Pass):
    """Analyze use-def relationships"""

    def apply(self, module: builtin.ModuleOp) -> None:
        # Build use-def chains
        for op in module.walk():
            for result in op.results:
                print(f"Value {result} defined by {op.name}")
                print(f"  Used by: {[use.operation.name for use in
    result.uses]}")

    def find_unused_values(self, module: builtin.ModuleOp) -> list[
    OpResult]:
        """Find values that are never used"""
        unused = []
        for op in module.walk():
            for result in op.results:
```

13

```
16                    if not result.uses:
17                        unused.append(result)
18            return unused
```

Listing 10: Use-Def Analysis

## 5.4 Dominance Analysis

Dominance tells us which operations must execute before others:

```
1  from xdsl.ir import Block
2  from xdsl.irdl.dominance import DominanceInfo
3
4  class DominanceAnalysis(Pass):
5      """Analyze control flow dominance"""
6
7      def apply(self, module: builtin.ModuleOp) -> None:
8          # For each function in the module
9          for func_op in module.body.ops:
10             if isinstance(func_op, func.FuncOp):
11                 dom_info = DominanceInfo(func_op.body)
12
13                 # Check dominance relationships
14                 for block in func_op.body.blocks:
15                     dominators = dom_info.get_dominators(block)
16                     print(f"Block {block} dominated by: {dominators}")
```

Listing 11: Dominance Analysis

## 5.5 Dataflow Analysis

Dataflow analysis propagates information through the program:

```
1  class ConstantPropagationAnalysis(Pass):
2      """Analyze which values are compile-time constants"""
3
4      def apply(self, module: builtin.ModuleOp) -> None:
5          known_constants = {}
6
7          # First pass: identify constants
8          for op in module.walk():
9              if isinstance(op, arith.Constant):
10                 known_constants[op.result] = op.value.value
11
12         # Iteratively propagate constants
13         changed = True
14         while changed:
15             changed = False
16             for op in module.walk():
17                 if isinstance(op, arith.Addi):
18                     # If both operands are known constants
19                     if (op.lhs in known_constants and
20                         op.rhs in known_constants):
21                         result_value = (known_constants[op.lhs] +
22                                         known_constants[op.rhs])
23                         if op.result not in known_constants:
24                             known_constants[op.result] = result_value
25                             changed = True
26
```

```
27        return known_constants
```
Listing 12: Constant Propagation Analysis

# 6  Chapter 4: Transformations and Optimizations

## 6.1  What is a Transformation?

Transformations modify the IR to improve performance, reduce code size, or prepare for further optimizations. Unlike analysis passes, transformations actually change the program structure.
    Key principles:

- **Correctness**: Must preserve program semantics

- **Profitability**: Should improve some metric (speed, size, power)

- **Composability**: Should work well with other transformations

## 6.2  Pattern-Based Rewriting

xDSL provides powerful pattern matching for transformations:

```
1  from xdsl.pattern_rewriter import (
2      PatternRewriter,
3      RewritePattern,
4      op_type_rewrite_pattern
5  )
6
7  class FoldAddZero(RewritePattern):
8      """Fold x + 0 -> x"""
9
10     @op_type_rewrite_pattern
11     def match_and_rewrite(self, op: arith.Addi,
12                           rewriter: PatternRewriter) -> None:
13         # Check if right operand is zero
14         if isinstance(op.rhs.owner, arith.Constant):
15             if op.rhs.owner.value.value == 0:
16                 # Replace all uses of the add with the left operand
17                 rewriter.replace_op(op, [op.lhs])
18                 return
19
20         # Check if left operand is zero
21         if isinstance(op.lhs.owner, arith.Constant):
22             if op.lhs.owner.value.value == 0:
23                 rewriter.replace_op(op, [op.rhs])
```
Listing 13: Simple Algebraic Optimization

## 6.3  Common Optimizations

### 6.3.1  Dead Code Elimination

Remove operations whose results are never used:

```
1  class DeadCodeElimination(Pass):
2      """Remove unused operations"""
3
4      def apply(self, module: builtin.ModuleOp) -> None:
```

```
5          ops_to_remove = []
6
7          for op in module.walk():
8              # Skip operations with side effects
9              if self.has_side_effects(op):
10                 continue
11
12             # Check if any results are used
13             all_dead = all(not result.uses for result in op.results)
14
15             if all_dead:
16                 ops_to_remove.append(op)
17
18         # Remove dead operations
19         for op in ops_to_remove:
20             op.erase()
21
22     def has_side_effects(self, op: Operation) -> bool:
23         """Check if operation has side effects"""
24         # Memory writes, function calls, etc.
25         return isinstance(op, (memref.Store, func.Call))
```

Listing 14: Dead Code Elimination

### 6.3.2 Common Subexpression Elimination

Identify and eliminate redundant computations:

```
1  class CommonSubexpressionElimination(Pass):
2      """Eliminate redundant computations"""
3
4      def apply(self, module: builtin.ModuleOp) -> None:
5          # Map from (op_type, operands) to result
6          expression_map = {}
7
8          for op in module.walk():
9              if self.is_pure(op):
10                 # Create a key for this expression
11                 key = self.get_expression_key(op)
12
13                 if key in expression_map:
14                     # Found duplicate - replace with existing
15                     existing_result = expression_map[key]
16                     op.results[0].replace_by(existing_result)
17                     op.erase()
18                 else:
19                     # First occurrence - remember it
20                     expression_map[key] = op.results[0]
21
22     def get_expression_key(self, op: Operation):
23         """Create hashable key for expression"""
24         return (op.name, tuple(op.operands))
```

Listing 15: Common Subexpression Elimination

### 6.3.3 Loop Optimizations

Optimize loop structures for better performance:

16

```python
1  class LoopUnrolling(Pass):
2      """Unroll small loops"""
3
4      def apply(self, module: builtin.ModuleOp) -> None:
5          for op in module.walk():
6              if isinstance(op, scf.For):
7                  if self.should_unroll(op):
8                      self.unroll_loop(op)
9
10     def should_unroll(self, loop: scf.For) -> bool:
11         """Decide if loop should be unrolled"""
12         # Get loop bounds if constant
13         if (isinstance(loop.lb.owner, arith.Constant) and
14             isinstance(loop.ub.owner, arith.Constant)):
15             lb = loop.lb.owner.value.value
16             ub = loop.ub.owner.value.value
17             iterations = ub - lb
18
19             # Unroll small loops
20             return iterations <= 4
21         return False
22
23     def unroll_loop(self, loop: scf.For):
24         """Completely unroll a loop"""
25         # Clone loop body for each iteration
26         # Update induction variable references
27         # Remove original loop
28         pass  # Implementation details omitted
```

Listing 16: Loop Unrolling

## 6.4 Lowering Transformations

Lowering transforms high-level operations into simpler ones:

```python
1  class LowerMatrixOps(Pass):
2      """Lower matrix operations to loops"""
3
4      def apply(self, module: builtin.ModuleOp) -> None:
5          for op in module.walk():
6              if isinstance(op, MatMulOp):
7                  self.lower_matmul(op)
8
9      def lower_matmul(self, matmul: MatMulOp):
10         """Lower matmul to three nested loops"""
11         builder = Builder.before(matmul)
12
13         # Get dimensions
14         m, k = matmul.lhs.type.shape
15         _, n = matmul.rhs.type.shape
16
17         # Allocate result
18         result = builder.insert(memref.Alloc.get([m, n], f32))
19
20         # Generate three nested loops
21         with builder.insert(scf.For(0, m, 1)) as i:
22             with builder.insert(scf.For(0, n, 1)) as j:
23                 # Initialize accumulator
```

```
24              zero = builder.insert(arith.Constant.from_float(0.0,
    f32))
25
26              with builder.insert(scf.For(0, k, 1, [zero])) as (k_idx
    , acc):
27                  # Load elements
28                  a_elem = builder.insert(memref.Load(matmul.lhs, [i,
     k_idx]))
29                  b_elem = builder.insert(memref.Load(matmul.rhs, [
    k_idx, j]))
30
31                  # Multiply and accumulate
32                  prod = builder.insert(arith.Mulf(a_elem, b_elem))
33                  new_acc = builder.insert(arith.Addf(acc, prod))
34
35                  scf.Yield(new_acc)
36
37              # Store result
38              builder.insert(memref.Store(acc, result, [i, j]))
39
40          # Replace matmul with lowered result
41          matmul.results[0].replace_by(result)
42          matmul.erase()
```

Listing 17: Lowering High-Level Operations

# 7 Chapter 5: Practical Examples

## 7.1 Building a Simple Expression Compiler

Let's build a complete compiler for arithmetic expressions:

```python
1 from dataclasses import dataclass
2 from typing import Union
3
4 # Define AST nodes
5 @dataclass
6 class Expr:
7     pass
8
9 @dataclass
10 class Number(Expr):
11     value: int
12
13 @dataclass
14 class BinaryOp(Expr):
15     left: Expr
16     op: str
17     right: Expr
18
19 # Parser (simplified)
20 def parse(text: str) -> Expr:
21     """Parse expression like '2 + 3 * 4'"""
22     # Implementation omitted for brevity
23     pass
24
25 # IR Generator
26 class ExprIRGen:
```

```python
27    def __init__(self):
28        self.builder = Builder()
29
30    def compile(self, expr: Expr) -> OpResult:
31        """Compile expression to xDSL IR"""
32        if isinstance(expr, Number):
33            return arith.Constant.from_int_and_width(
34                expr.value, 32).result
35
36        elif isinstance(expr, BinaryOp):
37            left = self.compile(expr.left)
38            right = self.compile(expr.right)
39
40            if expr.op == '+':
41                return arith.Addi(left, right).result
42            elif expr.op == '*':
43                return arith.Muli(left, right).result
44            # ... other operations
45
46 # Optimizer
47 def optimize(module: builtin.ModuleOp):
48     """Apply optimization passes"""
49     passes = [
50         ConstantFolding(),
51         CommonSubexpressionElimination(),
52         DeadCodeElimination(),
53     ]
54
55     for pass_instance in passes:
56         pass_instance.apply(module)
57
58 # Complete pipeline
59 def compile_expression(text: str):
60     # Parse
61     ast = parse(text)
62
63     # Generate IR
64     ir_gen = ExprIRGen()
65     result = ir_gen.compile(ast)
66
67     # Wrap in module
68     module = builtin.ModuleOp([
69         func.FuncOp("compute",
70                     func.FunctionType.from_lists([], [builtin.i32]),
71                     body=[func.Return(result)])
72     ])
73
74     # Optimize
75     optimize(module)
76
77     return module
```

Listing 18: Expression Compiler

## 7.2 Implementing a Domain-Specific Optimization

Let's create an optimization for neural network operations:

```
1   class FuseBatchNormIntoConv(Pass):
2       """Fuse batch normalization into preceding convolution"""
3
4       def apply(self, module: builtin.ModuleOp) -> None:
5           for conv in module.walk():
6               if isinstance(conv, nn.Conv2d):
7                   # Look for batch norm following conv
8                   for use in conv.result.uses:
9                       if isinstance(use.operation, nn.BatchNorm2d):
10                          self.fuse(conv, use.operation)
11
12      def fuse(self, conv: nn.Conv2d, bn: nn.BatchNorm2d):
13          """Fuse batch norm parameters into convolution"""
14          # Extract batch norm parameters
15          gamma = bn.gamma  # scale
16          beta = bn.beta    # shift
17          mean = bn.running_mean
18          var = bn.running_var
19          eps = bn.eps
20
21          # Compute fused weights and bias
22          std = sqrt(var + eps)
23          fused_weight = conv.weight * (gamma / std)
24          fused_bias = gamma * (conv.bias - mean) / std + beta
25
26          # Create new convolution with fused parameters
27          fused_conv = nn.Conv2d(
28              weight=fused_weight,
29              bias=fused_bias,
30              # ... other conv parameters
31          )
32
33          # Replace conv -> bn with fused conv
34          bn.result.replace_by(fused_conv.result)
35          conv.erase()
36          bn.erase()
```

Listing 19: Neural Network Optimization

## 7.3 Building a Mini ML Compiler

Let's create a compiler for a subset of ML operations:

```
1   # Define ML dialect
2   @dialect_def
3   class MLDialect(Dialect):
4       name = "ml"
5
6   @op_def
7   class LinearOp(IRDLOperation):
8       """Linear layer: y = Wx + b"""
9       name = "ml.linear"
10
11      input = operand_def(TensorType)
12      weight = operand_def(TensorType)
13      bias = operand_def(TensorType)
14
15      output = result_def(TensorType)
```

```python
@op_def
class ReluOp(IRDLOperation):
    """ReLU activation"""
    name = "ml.relu"

    input = operand_def(TensorType)
    output = result_def(TensorType)

# Optimization: Fuse Linear + ReLU
class FuseLinearRelu(RewritePattern):
    @op_type_rewrite_pattern
    def match_and_rewrite(self, relu: ReluOp,
                          rewriter: PatternRewriter):
        # Check if input is from linear
        if isinstance(relu.input.owner, LinearOp):
            linear = relu.input.owner

            # Create fused operation
            fused = LinearReluOp(
                linear.input,
                linear.weight,
                linear.bias
            )

            # Replace both ops with fused version
            rewriter.replace_op(relu, [fused.output])
            rewriter.erase_op(linear)

# Lower to hardware operations
class LowerMLToHardware(Pass):
    """Lower ML ops to target hardware"""

    def apply(self, module: builtin.ModuleOp):
        target = self.get_target()  # CPU, GPU, TPU, etc.

        for op in module.walk():
            if isinstance(op, LinearOp):
                if target == "gpu":
                    self.lower_linear_to_cublas(op)
                elif target == "cpu":
                    self.lower_linear_to_loops(op)
```

Listing 20: Mini ML Compiler

# 8 Chapter 6: Advanced Topics

## 8.1 Multi-Level IR and Progressive Lowering

Real compilers use multiple levels of abstraction:

```python
class CompilationPipeline:
    """Progressive lowering through multiple IR levels"""

    def compile(self, source_code: str) -> str:
        # Level 1: High-level ML operations
        ml_ir = parse_to_ml_dialect(source_code)
        ml_ir = optimize_ml_level(ml_ir)
```

```
8
9          # Level 2: Linear algebra operations
10         linalg_ir = lower_ml_to_linalg(ml_ir)
11         linalg_ir = optimize_linalg_level(linalg_ir)
12
13         # Level 3: Loops and memory operations
14         loop_ir = lower_linalg_to_loops(linalg_ir)
15         loop_ir = optimize_loops(loop_ir)
16
17         # Level 4: Target-specific operations
18         target_ir = lower_to_target(loop_ir)
19         target_ir = optimize_target_specific(target_ir)
20
21         # Code generation
22         return generate_code(target_ir)
```

Listing 21: Multi-Level Compilation Pipeline

## 8.2   Cost Models and Autotuning

Choosing the best optimization requires understanding costs:

```
1  class CostModel:
2      """Estimate operation costs"""
3
4      def estimate_cost(self, op: Operation) -> float:
5          if isinstance(op, arith.Mulf):
6              return 1.0  # FP multiply cost
7          elif isinstance(op, memref.Load):
8              return self.memory_latency(op)
9          elif isinstance(op, scf.For):
10             iterations = self.analyze_trip_count(op)
11             body_cost = sum(self.estimate_cost(body_op)
12                             for body_op in op.body.walk())
13             return iterations * body_cost
14
15 class AutotuningPass(Pass):
16     """Try different optimization strategies"""
17
18     def apply(self, module: builtin.ModuleOp):
19         strategies = [
20             self.try_vectorization,
21             self.try_loop_tiling,
22             self.try_parallelization,
23         ]
24
25         best_module = module
26         best_cost = self.cost_model.estimate_cost(module)
27
28         for strategy in strategies:
29             candidate = strategy(module.clone())
30             cost = self.cost_model.estimate_cost(candidate)
31
32             if cost < best_cost:
33                 best_module = candidate
34                 best_cost = cost
35
```

```
36          return best_module
```
<div align="center">Listing 22: Cost-Based Optimization</div>

## 8.3 Connection to ML Compilation

Modern ML frameworks are essentially compilers:

- **PyTorch JIT**: Traces Python code to TorchScript IR

- **TensorFlow XLA**: Compiles TF graphs to optimized kernels

- **JAX**: Uses XLA for JIT compilation

- **Apache TVM**: Optimizes ML models for various hardware

The concepts you've learned apply directly:

- **Graph optimization** = IR transformation

- **Kernel fusion** = Operation merging

- **Memory planning** = Register allocation

- **Quantization** = Type lowering

# 9 Chapter 7: Hands-On Exercises

## 9.1 Introduction: Learning by Doing

The exercises in this chapter are designed to solidify your understanding of compiler concepts through practical implementation. Each exercise builds upon the concepts we've covered, challenging you to apply your knowledge in increasingly sophisticated ways. These aren't just academic exercises - they mirror real-world compiler engineering challenges you'll encounter in production systems.

**Goals of These Exercises:**

- **Reinforce Core Concepts**: Apply IR manipulation, pattern matching, and transformation techniques in practice

- **Develop Compiler Intuition**: Understand trade-offs between different optimization strategies

- **Build Debugging Skills**: Learn to verify correctness and debug transformation passes

- **Connect Theory to Practice**: See how abstract concepts translate to real optimizations

## 9.2 Exercise 1: Build Your Own Optimizer - Store-Load Forwarding

### 9.2.1 Context and Motivation

Memory operations are often bottlenecks in program performance. When a program stores a value to memory and immediately loads it back, we're wasting cycles on unnecessary memory traffic. Store-load forwarding is a crucial optimization in modern compilers and CPUs that eliminates these redundant operations.

**Real-World Impact:** This optimization is so important that modern CPUs implement it in hardware. By doing it in the compiler, we can eliminate the operations entirely, reducing both memory bandwidth and instruction count.

**What You'll Learn:**

- How to track data flow through memory operations

- Alias analysis basics (when is it safe to forward?)

- The importance of maintaining program correctness during optimization

- How to reason about memory dependencies

```python
class StoreLoadForwarding(Pass):
    """Forward stored values to subsequent loads

    Transform patterns like:
        memref.store %value, %mem[%i, %j]
        %loaded = memref.load %mem[%i, %j]
    Into:
        memref.store %value, %mem[%i, %j]
        // Use %value directly instead of %loaded
    """

    def apply(self, module: builtin.ModuleOp):
        # TODO: Your implementation here
        # Step 1: Build a map of recent stores for each memory location
        # Step 2: For each load, check if we have a recent store
        # Step 3: Verify indices match exactly (be conservative!)
        # Step 4: Check for intervening operations that might
    invalidate
        # Step 5: Replace load result with stored value

        # Hints:
        # - Start simple: handle only same-block forwarding
        # - Be conservative: only forward when indices are identical
        # - Watch for: function calls, other stores to same array
        # - Remember: correctness > performance
        pass

    def can_forward(self, store_op, load_op):
        """Check if it's safe to forward from store to load"""
        # TODO: Implement safety checks
        pass
```

Listing 23: Exercise: Store-Load Forwarding

### 9.2.2 Testing Your Implementation

```python
def test_store_load_forwarding():
    """Test your optimization works correctly"""

    # Test case 1: Simple forwarding
    # store 42 to mem[0]
    # load from mem[0]   <-- should be replaced with 42

    # Test case 2: Different indices (no forwarding)
    # store 42 to mem[0]
    # load from mem[1]   <-- should NOT be forwarded

    # Test case 3: Intervening store
    # store 42 to mem[0]
    # store 100 to mem[0]
```

```
15      # load from mem[0]  <-- should get 100, not 42
16
17      # Test case 4: Across basic blocks (advanced)
18      # Requires dominance analysis
```

Listing 24: Test Cases for Store-Load Forwarding

## 9.3 Exercise 2: Create a Custom Dialect - Quantum Computing

### 9.3.1 Context and Motivation

Quantum computing represents a fundamentally different computational model. By creating a quantum dialect, you'll learn how compiler infrastructure can be extended to support entirely new paradigms. This exercise demonstrates xDSL's extensibility and teaches you to think about computation abstractly.

**Why This Matters:** As new computational paradigms emerge (quantum, neuromorphic, photonic), the ability to quickly prototype new compiler infrastructure becomes crucial. This exercise teaches you the skills needed to support future computing platforms.

**What You'll Learn:**

- How to design operations for non-standard computation models

- Type system design for domain-specific constraints

- Verification of domain-specific invariants

- Lowering from high-level domain concepts to implementation

```python
1  from xdsl.irdl import irdl_op_definition, irdl_attr_definition
2  from xdsl.ir import OpResult, SSAValue
3
4  @dialect_def
5  class QuantumDialect(Dialect):
6      name = "quantum"
7
8  # TODO: Define quantum-specific types
9  @irdl_attr_definition
10 class QubitType(TypeAttribute):
11     """Type representing a quantum bit"""
12     name = "quantum.qubit"
13
14 @irdl_attr_definition
15 class QuantumRegisterType(TypeAttribute):
16     """Type for quantum register of n qubits"""
17     name = "quantum.qreg"
18     size: int
19
20 # TODO: Define quantum gates
21 @irdl_op_definition
22 class HadamardOp(IRDLOperation):
23     """Hadamard gate: creates superposition"""
24     name = "quantum.h"
25
26     input_qubit = operand_def(QubitType)
27     output_qubit = result_def(QubitType)
28
29     def verify(self):
30         """Verify operation is well-formed"""
```

```
31          # Hadamard is always valid on a single qubit
32          pass
33
34 @irdl_op_definition
35 class CNOTOp(IRDLOperation):
36     """Controlled-NOT gate"""
37     name = "quantum.cnot"
38
39     control = operand_def(QubitType)
40     target = operand_def(QubitType)
41
42     # CNOT modifies both qubits
43     new_control = result_def(QubitType)
44     new_target = result_def(QubitType)
45
46 @irdl_op_definition
47 class MeasureOp(IRDLOperation):
48     """Measure qubit, collapsing to classical bit"""
49     name = "quantum.measure"
50
51     qubit = operand_def(QubitType)
52     classical_bit = result_def(IntegerType(1))
53
54 # TODO: Optimization passes
55 class QuantumCircuitOptimizer(Pass):
56     """Optimize quantum circuits"""
57
58     def apply(self, module: builtin.ModuleOp):
59         # TODO: Implement optimizations like:
60         # - Gate cancellation (H H = I)
61         # - Gate commutation
62         # - Circuit depth reduction
63         pass
```

Listing 25: Exercise: Quantum Dialect

### 9.3.2 Challenge Extensions

- Implement quantum circuit simulation

- Add error correction operations

- Create a cost model for quantum operations

- Implement decomposition to basic gate sets

## 9.4 Exercise 3: Implement Vectorization

### 9.4.1 Context and Motivation

Modern processors have SIMD (Single Instruction, Multiple Data) units that can process multiple data elements simultaneously. Vectorization transforms scalar operations to leverage these capabilities, often providing 4-8x speedups. This is a cornerstone optimization in high-performance computing.

**Industry Relevance:** Every production compiler implements auto-vectorization. ML frameworks rely heavily on vectorized operations for performance. Understanding vectorization is essential for performance engineering.

**What You'll Learn:**

- How to identify vectorizable patterns

- Dependency analysis for parallel execution

- Handling edge cases (loop remainders, alignment)

- Cost modeling for vectorization decisions

```python
class Vectorization(Pass):
    """Automatically vectorize loops

    Transform:
        for i in range(n):
            c[i] = a[i] + b[i]

    Into:
        for i in range(0, n-4, 4):
            c[i:i+4] = a[i:i+4] + b[i:i+4]  # Vector operation
        for i in range(n-n%4, n):
            c[i] = a[i] + b[i]  # Scalar remainder
    """

    def __init__(self, vector_width=4):
        self.vector_width = vector_width

    def apply(self, module: builtin.ModuleOp):
        for loop in self.find_loops(module):
            if self.is_vectorizable(loop):
                self.vectorize_loop(loop)

    def is_vectorizable(self, loop: scf.For) -> bool:
        """Check if loop can be vectorized"""
        # TODO: Check for:
        # 1. No loop-carried dependencies
        # 2. Contiguous memory access
        # 3. Supported operations (add, mul, etc.)
        # 4. No control flow in loop body
        pass

    def vectorize_loop(self, loop: scf.For):
        """Transform loop to use vector operations"""
        # TODO: Implementation steps:
        # 1. Create vector loop with stride = vector_width
        # 2. Replace scalar ops with vector ops
        # 3. Create remainder loop for n % vector_width iterations
        # 4. Update memory operations to vector loads/stores
        pass

    def analyze_dependencies(self, loop: scf.For) -> bool:
        """Check for loop-carried dependencies"""
        # TODO: Implement dependency analysis
        # Look for patterns like a[i] = a[i-1] + b[i]
        pass
```

Listing 26: Exercise: Auto-Vectorization

### 9.4.2 Advanced Challenges

- Handle non-unit stride access patterns

27

- Implement cost model to decide when vectorization is profitable

- Support reduction operations (sum, max)

- Handle conditional operations with masking

## 9.5 Exercise 4: Build a Simple JIT Compiler

### 9.5.1 Context and Motivation

Just-In-Time compilation bridges the gap between interpretation and ahead-of-time compilation. By compiling code at runtime, JIT compilers can optimize based on actual execution patterns. This exercise gives you hands-on experience with the full compilation pipeline.

**Real-World Applications:** JIT compilation powers JavaScript engines (V8), Java VMs (HotSpot), Python (PyPy), and ML frameworks (JAX, PyTorch). Understanding JIT compilation is crucial for modern system design.

**What You'll Learn:**

- The complete compilation pipeline from source to execution

- Runtime code generation techniques

- Optimization decisions based on runtime information

- Integration between compiled and interpreted code

```python
import ctypes
from typing import Dict, Callable

class ExpressionJIT:
    """JIT compile and execute mathematical expressions

    This JIT compiler will:
    1. Parse expressions into AST
    2. Generate xDSL IR
    3. Optimize the IR
    4. Generate executable code
    5. Return a callable Python function
    """

    def __init__(self):
        self.cache: Dict[str, Callable] = {}
        self.optimization_level = 2

    def compile(self, expr_string: str) -> Callable:
        """Compile expression to native function"""

        # Check cache first
        if expr_string in self.cache:
            return self.cache[expr_string]

        # Step 1: Parse expression
        ast = self.parse_expression(expr_string)

        # Step 2: Generate IR
        ir_module = self.generate_ir(ast)

        # Step 3: Optimize
```

```python
        if self.optimization_level > 0:
            ir_module = self.optimize_ir(ir_module)

        # Step 4: Generate machine code
        machine_code = self.generate_code(ir_module)

        # Step 5: Create callable wrapper
        func = self.create_callable(machine_code)

        # Cache for future use
        self.cache[expr_string] = func

        return func

    def parse_expression(self, expr: str):
        """Parse expression into AST"""
        # TODO: Implement expression parser
        # Support: numbers, variables, +, -, *, /, parentheses
        # Example: "x * 2 + y" -> BinaryOp('+', BinaryOp('*', Var('x'),
    Num(2)), Var('y'))
        pass

    def generate_ir(self, ast) -> builtin.ModuleOp:
        """Generate xDSL IR from AST"""
        # TODO: Create IR module with function
        # Function should take dict of variables, return result
        pass

    def optimize_ir(self, module: builtin.ModuleOp) -> builtin.ModuleOp
    :
        """Apply optimization passes"""
        passes = [
            ConstantFolding(),
            CommonSubexpressionElimination(),
            DeadCodeElimination(),
        ]

        for opt_pass in passes:
            opt_pass.apply(module)

        return module

    def generate_code(self, module: builtin.ModuleOp) -> bytes:
        """Generate machine code or Python bytecode"""
        # TODO: Options:
        # 1. Generate Python bytecode using 'compile()'
        # 2. Use LLVM bindings to generate native code
        # 3. Generate C code and compile with ctypes
        pass

    def create_callable(self, code: bytes) -> Callable:
        """Create Python-callable function from machine code"""
        # TODO: Wrap machine code in Python-callable interface
        pass

# Usage example:
jit = ExpressionJIT()
f = jit.compile("x * 2 + y")
```

```
89  result = f(x=5, y=3)  # Should return 13
90
91  # Advanced: Support for conditionals
92  g = jit.compile("x > 0 ? x * 2 : -x")
93  result2 = g(x=-5)  # Should return 5
```
Listing 27: Exercise: JIT Compiler

### 9.5.2 Implementation Hints

```
1  # Simple approach using Python's compile()
2  def generate_python_code(self, module: builtin.ModuleOp) -> str:
3      """Generate Python source from IR"""
4      code = "def compiled_func(**variables):\n"
5      # Walk IR and generate Python code
6      # ...
7      return code
8
9  def create_callable_simple(self, py_code: str) -> Callable:
10     """Create function from Python code"""
11     exec_globals = {}
12     exec(py_code, exec_globals)
13     return exec_globals['compiled_func']
```
Listing 28: JIT Implementation Starter Code

# 10  Chapter 8: Real-World Applications

## 10.1  Case Study 1: Optimizing Neural Network Inference

Let's optimize a real neural network for deployment:

```
1  def optimize_nn_for_deployment(model: NeuralNetwork):
2      """Optimize neural network for inference"""
3
4      # Convert to IR
5      ir = convert_nn_to_ir(model)
6
7      # Apply optimizations
8      optimizations = [
9          # Structural optimizations
10         FuseLayers(),            # Conv+BN, Linear+ReLU
11         RemoveDropout(),         # Not needed for inference
12
13         # Quantization
14         QuantizeWeights(bits=8),
15         QuantizeActivations(bits=8),
16
17         # Memory optimizations
18         InplaceOperations(),    # Reuse buffers
19         MemoryPlanning(),       # Minimize allocation
20
21         # Target-specific
22         UseTargetIntrinsics(),  # AVX, NEON, etc.
23         AutoTuneKernels(),      # Platform-specific tuning
24     ]
25
```

```
26        for opt in optimizations:
27            ir = opt.apply(ir)
28
29        return ir
```

Listing 29: Neural Network Optimization Pipeline

## 10.2   Case Study 2: Domain-Specific Language Compiler

Build a compiler for financial computations:

```
1  # Define financial operations dialect
2  @dialect_def
3  class FinanceDialect(Dialect):
4      name = "finance"
5
6      @op_def
7      class BlackScholesOp(IRDLOperation):
8          """Black-Scholes option pricing"""
9          name = "finance.black_scholes"
10
11         spot = operand_def(f64)
12         strike = operand_def(f64)
13         time = operand_def(f64)
14         rate = operand_def(f64)
15         volatility = operand_def(f64)
16
17         call_price = result_def(f64)
18         put_price = result_def(f64)
19
20 # Optimize for numerical stability
21 class NumericalStabilityPass(Pass):
22     """Ensure numerical stability in financial calculations"""
23
24     def apply(self, module: builtin.ModuleOp):
25         for op in module.walk():
26             if isinstance(op, math.ExpOp):
27                 # Prevent overflow/underflow
28                 self.add_range_check(op)
29             elif isinstance(op, arith.DivfOp):
30                 # Prevent division by zero
31                 self.add_zero_check(op)
```

Listing 30: Financial DSL Compiler

## 10.3   Case Study 3: Compiler for Distributed Computing

Compile programs for distributed execution:

```
1  @dialect_def
2  class DistributedDialect(Dialect):
3      name = "dist"
4
5      @op_def
6      class AllReduceOp(IRDLOperation):
7          """All-reduce across distributed nodes"""
8          name = "dist.allreduce"
9
10         input = operand_def(TensorType)
```

```
11        output = result_def(TensorType)
12        reduction = attr_def(str)  # "sum", "max", etc.
13
14 class DistributionPass(Pass):
15     """Distribute computation across nodes"""
16
17     def apply(self, module: builtin.ModuleOp):
18         # Analyze data dependencies
19         deps = self.analyze_dependencies(module)
20
21         # Partition computation
22         partitions = self.partition_computation(module, deps)
23
24         # Insert communication operations
25         for partition in partitions:
26             self.insert_communication_ops(partition)
27
28         # Generate distributed code
29         return self.generate_distributed_code(partitions)
```
Listing 31: Distributed Computing Compiler

# 11 Project: Building a Matrix Operations Compiler

## 11.1 Project Overview

In this project, you'll build a complete compiler for a domain-specific language (DSL) focused on matrix operations. This project demonstrates the full compilation pipeline: from parsing high-level Python-like code, through IR generation and optimization, to generating efficient low-level code. The project is inspired by real-world compilers like those used in machine learning frameworks (TensorFlow, PyTorch) that optimize matrix operations for performance. The code is available at https://github.com/djtodoro/matrix-toy-lang.

**Why Matrix Operations?** Matrix operations are fundamental to scientific computing, machine learning, and graphics. They offer rich optimization opportunities like:

- Algebraic simplifications (e.g., $(A^T)^T = A$)

- Operation fusion (combining multiple operations into one)

- Memory layout optimizations (row-major vs column-major)

- Parallelization opportunities

**Learning Objectives:**

- Understand the complete compiler pipeline architecture

- Learn how to leverage existing tools (Python AST) with custom IR

- Implement pattern-matching optimizations

- Experience the design decisions in creating a custom dialect

- Practice incremental development and testing

**Project Structure:** The compiler consists of five main phases:

1. **Frontend**: Parse Python code using the `ast` module

2. **Dialect Design**: Create a custom matrix operations dialect

3. **IR Generation**: Convert Python AST to dialect operations

4. **Optimization**: Apply pattern-based transformations

5. **Lowering**: Generate executable code (optional)

## 11.2 Phase 1: Frontend - Parsing the Input

The compiler begins by parsing the input program. For this project, we'll use Python as our source language and leverage Python's built-in `ast` module to parse the code.

**Example Input Program:**

```python
def matrix_computation(A, B):
    C = A @ B              # Matrix multiplication
    D = C.T.T              # Double transpose (can be optimized to C)
    E = D * 2.0            # Scalar multiplication
    F = E.T                # Single transpose
    G = F.T                # Another transpose (F.T.T = E)
    H = G + B.T            # Matrix addition
    result = H.T.T         # Final double transpose
    return result
```

Listing 32: Example matrix operations with optimization opportunities

This example contains several optimization opportunities:

- Three instances of double transpose that can be eliminated

- Potential for operation fusion

- Opportunities for common subexpression elimination

## 11.3 Phase 2: Building the Parser

The parser traverses the Python AST and extracts matrix operations. We use the Visitor pattern to walk the AST and identify relevant operations.

**Key Parser Tasks:**

1. Identify matrix operations (matmul, transpose, add, scalar mul)

2. Track variable definitions and uses

3. Build a sequence of operations for IR generation

4. Detect optimization patterns (like double transpose)

**Parser Architecture:**

```python
import ast

class MatrixOperationExtractor(ast.NodeVisitor):
    def __init__(self):
        self.operations = []  # List of extracted operations
        self.variables = {}   # Variable tracking

    def visit_Assign(self, node):
        # Extract the operation from the right-hand side
        op_info = self.analyze_expression(node.value)
        target = node.targets[0].id
```

```
12
13        self.operations.append({
14            'target': target,
15            'operation': op_info
16        })
17
18    def analyze_expression(self, node):
19        # Pattern matching for different operations
20        if isinstance(node, ast.BinOp):
21            if isinstance(node.op, ast.MatMult):  # @ operator
22                return {'op': 'matmul', ...}
23            elif isinstance(node.op, ast.Mult):   # * operator
24                return {'op': 'scalar_mul', ...}
25        elif isinstance(node, ast.Attribute):
26            if node.attr == 'T':  # Transpose
27                # Check for double transpose pattern
28                base = self.analyze_expression(node.value)
29                if base.get('op') == 'transpose':
30                    return {'op': 'double_transpose', ...}
31                return {'op': 'transpose', ...}
32        # ... handle other operations
```

<div align="center">Listing 33: Parser structure and key methods</div>

**Key Insights:**

- The parser performs early pattern detection (e.g., double transpose)

- Operations are stored as a linear sequence for easy IR generation

- Variable tracking enables use-def chain analysis

### 11.4  Phase 3: Designing the Matrix Dialect

A dialect in xDSL/MLIR is a collection of operations and types that represent concepts at a specific abstraction level. For matrix operations, we need to design operations that:

- Capture the semantics of matrix operations precisely

- Enable verification of correctness (e.g., dimension checking)

- Facilitate pattern matching for optimizations

- Support gradual lowering to hardware operations

**Core Operations in the Matrix Dialect:**

```
1  from xdsl.irdl import IRDLOperation, irdl_op_definition
2  from xdsl.ir import Dialect
3
4  @irdl_op_definition
5  class MatMulOp(IRDLOperation):
6      """Matrix multiplication: C = A @ B"""
7      name = "matrix.matmul"
8
9      # Shape inference: (m n) @ (n p)     (m p)
10     def verify(self):
11         # Check dimension compatibility
12         assert self.lhs.type.cols == self.rhs.type.rows
13
```

```
14  @irdl_op_definition
15  class TransposeOp(IRDLOperation):
16      """Matrix transpose: B = A^T"""
17      name = "matrix.transpose"
18
19      # Shape transformation: (m  n)      (n  m)
20      # Key for optimization patterns
21
22  @irdl_op_definition
23  class ScalarMulOp(IRDLOperation):
24      """Element-wise scalar multiplication: B =   A  """
25      name = "matrix.scalar_mul"
26
27      # Preserves shape, scales values
28
29  @irdl_op_definition
30  class AddOp(IRDLOperation):
31      """Element-wise addition: C = A + B"""
32      name = "matrix.add"
33
34      # Requires matching dimensions
```

Listing 34: Key dialect operations and their design

**Design Decisions:**

1. **Type System**: Use static shapes when known for better optimization

2. **Verification**: Each operation includes dimension checking

3. **SSA Form**: Operations produce new values (functional style)

4. **Attributes vs Operands**: Constants as attributes, variables as operands

## 11.5   Phase 4: IR Generation

The IR generator transforms the parsed operations into the matrix dialect IR. This phase bridges the gap between the source language (Python) and our optimization framework.

**IR Generation Process:**

1. Create a function wrapper for the operations

2. Map Python variables to SSA values

3. Generate dialect operations in sequence

4. Maintain use-def chains through SSA form

**Example IR Generation:**

```
1  class MatrixIRGenerator:
2      def __init__(self):
3          self.var_map = {}  # Maps Python variables to SSA values
4
5      def generate_operation(self, op_info):
6          op_type = op_info['operation']['op']
7          target = op_info['target']
8
9          if op_type == 'matmul':
10             # C = A @ B becomes:
```

```
11          # %c = matrix.matmul %a, %b : (f32, f32) -> f32
12          left = self.var_map[op_info['operation']['left']]
13          right = self.var_map[op_info['operation']['right']]
14          result = MatMulOp(left, right).result
15          self.var_map[target] = result
16
17      elif op_type == 'double_transpose':
18          # D = C.T.T becomes (before optimization):
19          # %tmp = matrix.transpose %c : (m n) -> (n m)
20          # %d = matrix.transpose %tmp : (n m) -> (m n)
21          matrix = self.var_map[op_info['operation']['matrix']]
22          temp = TransposeOp(matrix).result
23          result = TransposeOp(temp).result
24          self.var_map[target] = result
```

Listing 35: Translating parsed operations to IR

### Generated IR Example:

```
1  func @matrix_computation(%A: matrix<2x3xf32>, %B: matrix<3x2xf32>)
2      -> matrix<2x2xf32> {
3    %C = matrix.matmul %A, %B : (matrix<2x3xf32>, matrix<3x2xf32>)
4                                -> matrix<2x2xf32>
5    %tmp1 = matrix.transpose %C : matrix<2x2xf32> -> matrix<2x2xf32>
6    %D = matrix.transpose %tmp1 : matrix<2x2xf32> -> matrix<2x2xf32>
7    %E = matrix.scalar_mul %D, 2.0 : matrix<2x2xf32> -> matrix<2x2xf32>
8    %F = matrix.transpose %E : matrix<2x2xf32> -> matrix<2x2xf32>
9    %G = matrix.transpose %F : matrix<2x2xf32> -> matrix<2x2xf32>
10   %B_t = matrix.transpose %B : matrix<3x2xf32> -> matrix<2x3xf32>
11   %H = matrix.add %G, %B_t : (matrix<2x2xf32>, matrix<2x3xf32>)
12                                -> matrix<2x2xf32>
13   %tmp2 = matrix.transpose %H : matrix<2x2xf32> -> matrix<2x2xf32>
14   %result = matrix.transpose %tmp2 : matrix<2x2xf32> -> matrix<2x2xf32>
15   return %result : matrix<2x2xf32>
16 }
```

Listing 36: Matrix dialect IR before optimization

## 11.6 Phase 5: Optimization - Double Transpose Elimination

The optimization phase applies pattern-based transformations to improve the IR. Our primary optimization eliminates the mathematical identity $(A^T)^T = A$.

### Pattern Matching Approach:

```
1  from xdsl.pattern_rewriter import PatternRewriter, RewritePattern
2
3  class DoubleTransposeElimination(RewritePattern):
4      """Eliminate pattern: transpose(transpose(X)) = X"""
5
6      def match_and_rewrite(self, op: TransposeOp, rewriter):
7          # Check if input is also a transpose
8          if isinstance(op.input.owner, TransposeOp):
9              # Found the pattern!
10             original = op.input.owner.input
11
12             # Replace double transpose with original value
13             rewriter.replace_op(op, [original])
14
15             # Statistics tracking
```

```
16            self.eliminated_count += 1
```

**How the Optimization Works:**

1. **Pattern Detection**: Walk the IR looking for transpose operations

2. **Match Checking**: Check if the input to a transpose is another transpose

3. **Replacement**: Replace the double transpose with the original matrix

4. **Cleanup**: Remove dead code (unused intermediate transposes)

**Before and After Optimization:**

```
1 // Before optimization:
2 %tmp1 = matrix.transpose %C : matrix<2x2xf32> -> matrix<2x2xf32>
3 %D = matrix.transpose %tmp1 : matrix<2x2xf32> -> matrix<2x2xf32>
4
5 // After optimization:
6 %D = %C  // Direct use, no transposes needed
```

**Additional Optimization Opportunities:**

- **Algebraic Simplifications**: $(A + B)^T = A^T + B^T$

- **Operation Fusion**: Combine scalar_mul operations

- **Common Subexpression Elimination**: Reuse computed transposes

- **Matrix Chain Optimization**: Optimal parenthesization for multiple multiplications

## 11.7   Phase 6: The Complete Compiler Pipeline

The complete compiler integrates all phases into a cohesive pipeline. Here's how the components work together:

**Pipeline Architecture:**

```
1 class MatrixCompiler:
2     def compile(self, source_code: str) -> str:
3         # Phase 1: Parse Python to AST
4         ast_tree = ast.parse(source_code)
5
6         # Phase 2: Extract matrix operations
7         operations = MatrixOperationExtractor().extract(ast_tree)
8
9         # Phase 3: Generate IR
10        module = MatrixIRGenerator().generate(operations)
11
12        # Phase 4: Optimize
13        optimizer = MatrixOptimizationPipeline()
14        optimizer.apply(module)
15
16        # Phase 5: Lower to target (optional)
17        # Could lower to LLVM, C, or other targets
18
19        return module.to_string()
```

**Running the Compiler:**

```
# Compile and optimize
python matrix_compiler.py input.py --optimize

# View IR at different stages
python matrix_compiler.py input.py --emit-ir
python matrix_compiler.py input.py --emit-optimized-ir

# Generate executable (if lowering is implemented)
python matrix_compiler.py input.py -o output
```
Listing 40: Command-line usage

**Example Compilation Flow:**

1. **Input**: Python code with matrix operations

2. **Parse**: Extract operation sequence

3. **Generate**: Create matrix dialect IR

4. **Optimize**: Apply transformations (3 double transposes eliminated)

5. **Output**: Optimized IR or lowered code

## 11.8    Testing and Validation

A robust compiler needs comprehensive testing at each phase:
**Testing Strategy:**

1. **Unit Tests**: Test each phase independently

2. **Integration Tests**: Test the complete pipeline

3. **Correctness Tests**: Verify semantic preservation

4. **Performance Tests**: Measure optimization effectiveness

**Key Test Cases:**

```python
def test_double_transpose_elimination ():
    # Input code with double transpose
    source = "D = C.T.T"

    # Compile without optimization
    ir_before = compile(source, optimize=False)
    assert count_transposes(ir_before) == 2

    # Compile with optimization
    ir_after = compile(source, optimize=True)
    assert count_transposes(ir_after) == 0

    # Verify semantic equivalence
    assert execute(ir_before) == execute(ir_after)
```
Listing 41: Testing optimization correctness

## 11.9   Project Extensions and Exercises

**Recommended Extensions:**

1. **Additional Optimizations**:

   - Implement $(A \cdot B)^T = B^T \cdot A^T$ transformation
   - Add common subexpression elimination
   - Implement operation fusion (combine multiple scalar multiplications)

2. **Advanced Features**:

   - Add support for matrix inverse operations
   - Implement shape inference for unknown dimensions
   - Add support for batched operations

3. **Code Generation**:

   - Lower to LLVM IR for execution
   - Generate optimized C code
   - Target GPU operations (CUDA/OpenCL)

4. **Analysis Passes**:

   - Add cost model for operations
   - Implement data dependency analysis
   - Add memory usage optimization

**Project Structure:**

```
matrix-compiler/
    src/
        frontend/
            parser.py           # AST parsing
            lexer.py            # Tokenization (if needed)
        dialect/
            matrix_ops.py       # Operation definitions
            types.py            # Type definitions
        transforms/
            optimize.py         # Optimization passes
            lower.py            # Lowering passes
        compiler.py             # Main driver
    tests/
        test_parser.py          # Parser tests
        test_dialect.py         # Dialect tests
        test_optimize.py        # Optimization tests
    examples/
        simple.py               # Simple examples
        complex.py              # Complex examples
    README.md                   # Documentation
```

Listing 42: Recommended project organization

## 11.10 Learning Outcomes and Key Takeaways

By completing this project, students will have gained practical experience with:
**Core Compiler Concepts:**

- **Frontend Design**: Parsing and semantic analysis

- **IR Design**: Creating domain-specific abstractions

- **Pattern Matching**: Implementing optimization patterns

- **SSA Form**: Understanding use-def chains

- **Verification**: Ensuring correctness at each phase

**Software Engineering Skills:**

- Designing modular, testable compiler components

- Writing comprehensive test suites

- Documenting complex software systems

- Using modern compiler infrastructure (xDSL/MLIR)

**Mathematical Foundations:**

- Matrix algebra and its optimization opportunities

- Correctness preservation through transformations

- Cost models and performance analysis

## 11.11 Conclusion

This project demonstrates the complete journey of building a domain-specific compiler. Starting from high-level Python code, we've created a custom IR dialect, implemented pattern-based optimizations, and seen how modern compiler infrastructure enables rapid development of specialized compilers.

The matrix operations compiler serves as a foundation for understanding larger systems like TensorFlow's XLA, PyTorch's TorchScript, and other domain-specific compilers that power modern machine learning and scientific computing.
**Next Steps:**

- Explore the xDSL/MLIR ecosystem for more advanced dialects

- Study production compilers like LLVM, GCC, or specialized ML compilers

- Implement more sophisticated optimizations (auto-vectorization, parallelization)

- Build compilers for other domains (graphics, cryptography, databases)

Remember: Compilers are not just about code generation—they're about bridging the gap between human intent and machine execution, enabling both productivity and performance.