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

Djordje Todorovic

August 2025

## Contents

# 1 Introduction: Compilers as Representation and Transformation Engines

## 1.1 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

## 1.2 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

## 1.3 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.

## 1.4 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

# 2 Chapter 1: Understanding Intermediate Representations (IR)

## 2.1 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)

## 2.2 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:

### 2.2.1 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

### 2.2.2 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

### 2.2.3 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

## 2.3 Learning from Giants: LLVM and MLIR

### 2.3.1 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

### 2.3.2 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

## 2.4 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

## 2.5 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

## 2.6 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

## 2.7 Understanding IR Structure

The IR in xDSL follows a hierarchical structure:

Module

contains

Function

has body

Block

contains

Operation

## 2.8 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)

# 3 Chapter 2: Operations and Operators

## 3.1 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)

## 3.2 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

## 3.3 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

## 3.4 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

# 4  Chapter 3: Analysis Passes

## 4.1  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)

## 4.2 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

## 4.3 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:
```

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

Listing 10: Use-Def Analysis

## 4.4 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

## 4.5 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
```
<center>Listing 12: Constant Propagation Analysis</center>

# 5 Chapter 4: Transformations and Optimizations

## 5.1 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

## 5.2 Pattern-Based Rewriting

xDSL provides powerful pattern matching for transformations:

```python
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])
```
<center>Listing 13: Simple Algebraic Optimization</center>

## 5.3 Common Optimizations

### 5.3.1 Dead Code Elimination

Remove operations whose results are never used:

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

<center>12</center>

```
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

### 5.3.2 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

### 5.3.3 Loop Optimizations

Optimize loop structures for better performance:

```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

## 5.4 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

# 6 Chapter 5: Practical Examples

## 6.1 Building a Simple Expression Compiler

Let's build a complete compiler for arithmetic expressions:

```
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

## 6.2 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

## 6.3 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)
```

```
16
17  @op_def
18  class ReluOp(IRDLOperation):
19      """ReLU activation"""
20      name = "ml.relu"
21
22      input = operand_def(TensorType)
23      output = result_def(TensorType)
24
25  # Optimization: Fuse Linear + ReLU
26  class FuseLinearRelu(RewritePattern):
27      @op_type_rewrite_pattern
28      def match_and_rewrite(self, relu: ReluOp,
29                            rewriter: PatternRewriter):
30          # Check if input is from linear
31          if isinstance(relu.input.owner, LinearOp):
32              linear = relu.input.owner
33
34              # Create fused operation
35              fused = LinearReluOp(
36                  linear.input,
37                  linear.weight,
38                  linear.bias
39              )
40
41              # Replace both ops with fused version
42              rewriter.replace_op(relu, [fused.output])
43              rewriter.erase_op(linear)
44
45  # Lower to hardware operations
46  class LowerMLToHardware(Pass):
47      """Lower ML ops to target hardware"""
48
49      def apply(self, module: builtin.ModuleOp):
50          target = self.get_target()  # CPU, GPU, TPU, etc.
51
52          for op in module.walk():
53              if isinstance(op, LinearOp):
54                  if target == "gpu":
55                      self.lower_linear_to_cublas(op)
56                  elif target == "cpu":
57                      self.lower_linear_to_loops(op)
```

Listing 20: Mini ML Compiler

# 7 Chapter 6: Advanced Topics

## 7.1 Multi-Level IR and Progressive Lowering

Real compilers use multiple levels of abstraction:

```
1  class CompilationPipeline:
2      """Progressive lowering through multiple IR levels"""
3
4      def compile(self, source_code: str) -> str:
5          # Level 1: High-level ML operations
6          ml_ir = parse_to_ml_dialect(source_code)
7          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

## 7.2   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
```

Listing 22: Cost-Based Optimization

## 7.3  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

# 8  Chapter 7: Hands-On Exercises

## 8.1  Exercise 1: Build Your Own Optimizer

Implement an optimization that eliminates redundant memory operations:

```
1 class StoreLoadForwarding(Pass):
2     """Forward stored values to subsequent loads"""
3
4     def apply(self, module: builtin.ModuleOp):
5         # TODO: Your implementation here
6         # Hint: Track recent stores to each memory location
7         # Replace loads with stored values when possible
8         pass
```

Listing 23: Exercise: Store-Load Forwarding

## 8.2  Exercise 2: Create a Custom Dialect

Design a dialect for quantum computing operations:

```
1 @dialect_def
2 class QuantumDialect(Dialect):
3     name = "quantum"
4
5 # TODO: Define operations like:
6 # - Hadamard gate
7 # - CNOT gate
8 # - Measurement
9 # - Quantum circuit optimization passes
```

Listing 24: Exercise: Quantum Dialect

## 8.3 Exercise 3: Implement Vectorization

Transform scalar operations to vector operations:

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

    def apply(self, module: builtin.ModuleOp):
        # TODO: Find vectorizable loops
        # Transform scalar ops to vector ops
        # Handle loop remainder
        pass
```

Listing 25: Exercise: Auto-Vectorization

## 8.4 Exercise 4: Build a Simple JIT Compiler

Create a JIT compiler for mathematical expressions:

```python
class ExpressionJIT:
    """JIT compile and execute expressions"""

    def compile(self, expr_string: str):
        # Parse expression
        # Generate IR
        # Optimize
        # Generate machine code (or Python bytecode)
        # Return callable function
        pass

# Usage:
jit = ExpressionJIT()
f = jit.compile("x * 2 + y")
result = f(x=5, y=3)  # Should return 13
```

Listing 26: Exercise: JIT Compiler

# 9 Chapter 8: Real-World Applications

## 9.1 Case Study 1: Optimizing Neural Network Inference

Let's optimize a real neural network for deployment:

```python
def optimize_nn_for_deployment(model: NeuralNetwork):
    """Optimize neural network for inference"""

    # Convert to IR
    ir = convert_nn_to_ir(model)

    # Apply optimizations
    optimizations = [
        # Structural optimizations
        FuseLayers(),             # Conv+BN, Linear+ReLU
        RemoveDropout(),          # Not needed for inference

        # Quantization
        QuantizeWeights(bits=8),
        QuantizeActivations(bits=8),

```

```
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 27: Neural Network Optimization Pipeline

## 9.2 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 28: Financial DSL Compiler

## 9.3 Case Study 3: Compiler for Distributed Computing

Compile programs for distributed execution:

```python
@dialect_def
class DistributedDialect(Dialect):
    name = "dist"

    @op_def
    class AllReduceOp(IRDLOperation):
        """All-reduce across distributed nodes"""
        name = "dist.allreduce"

        input = operand_def(TensorType)
        output = result_def(TensorType)
        reduction = attr_def(str)  # "sum", "max", etc.

class DistributionPass(Pass):
    """Distribute computation across nodes"""

    def apply(self, module: builtin.ModuleOp):
        # Analyze data dependencies
        deps = self.analyze_dependencies(module)

        # Partition computation
        partitions = self.partition_computation(module, deps)

        # Insert communication operations
        for partition in partitions:
            self.insert_communication_ops(partition)

        # Generate distributed code
        return self.generate_distributed_code(partitions)
```

Listing 29: Distributed Computing Compiler

# 10 Conclusion: Your Journey Forward

## 10.1 What You've Learned

Congratulations! You now understand:

- How compilers represent programs as structured data (IR)

- The role of operations, attributes, and types

- How to write analysis passes to understand code

- How to implement transformations that optimize programs

- The connection between compiler techniques and ML optimization

## 10.2 Next Steps

Your compiler journey doesn't end here:

1. **Explore Real Compilers**: Study LLVM, GCC, or V8

2. **Contribute to xDSL**: Add new dialects or optimizations

3. **Apply to Your Domain**: Build domain-specific compilers

4. **Learn Advanced Topics**:

- Polyhedral compilation
- Verified compilation
- Quantum compilation
- Neuromorphic compilation

## 10.3 Resources for Further Learning

- **xDSL Documentation**: `https://xdsl.dev`

- **MLIR Tutorial**: `https://mlir.llvm.org/docs/Tutorials/`

- **Engineering a Compiler** (Cooper & Torczon)

- **Modern Compiler Implementation** (Appel)

- **The SSA Book**: `http://ssabook.gforge.inria.fr/`

## 10.4 Final Thoughts

Compilers are the bridge between human intent and machine execution. As AI/ML systems become more complex and heterogeneous hardware becomes more prevalent, understanding compilation becomes crucial. The techniques you've learned - representation, analysis, and transformation - are fundamental to building efficient, scalable systems.

Remember: Every optimization you write makes programs faster, every analysis you implement makes systems more reliable, and every transformation you create pushes the boundaries of what's computationally possible.

Happy compiling!

# A Appendix A: xDSL Quick Reference

## A.1 Common Operations

```
1  # Arithmetic
2  arith.Constant.from_int_and_width(42, 32)   # Integer constant
3  arith.Constant.from_float(3.14, f32)        # Float constant
4  arith.Addi(a, b)   # Integer addition
5  arith.Muli(a, b)   # Integer multiplication
6  arith.Addf(x, y)   # Float addition
7  arith.Mulf(x, y)   # Float multiplication
8
9  # Memory
10 memref.Alloc.get(shape=[10, 20], element_type=f32)
11 memref.Load(mem, indices)
12 memref.Store(value, mem, indices)
13
14 # Control Flow
15 scf.For(start, end, step, init_values)
16 scf.If(condition, has_else=True)
17 scf.While(init_values)
18
19 # Functions
20 func.FuncOp(name, func_type, body)
21 func.Call(func_name, args)
22 func.Return(values)
```

Listing 30: xDSL Operations Quick Reference

## A.2   Common Patterns

```python
# Walking operations
for op in module.walk():
    process(op)

# Pattern matching
if isinstance(op, arith.Addi):
    handle_add(op)

# Building IR
with Builder.implicit_region():
    result = arith.Addi(a, b).result

# Replacing operations
old_op.results[0].replace_by(new_value)
old_op.erase()
```

Listing 31: Common xDSL Patterns

# B   Appendix B: Setting Up Your Environment

## B.1   Installation

```bash
# Install xDSL
pip install xdsl

# For development
git clone https://github.com/xdslproject/xdsl.git
cd xdsl
pip install -e ".[dev]"

# Run tests
pytest tests/
```

Listing 32: Installing xDSL

## B.2   Your First xDSL Program

```python
#!/usr/bin/env python3

from xdsl.context import Context
from xdsl.dialects import builtin, arith, func
from xdsl.printer import Printer

# Create context and load dialects
ctx = Context()
ctx.load_dialect(builtin.Builtin)
ctx.load_dialect(arith.Arith)
ctx.load_dialect(func.Func)

# Build a simple function
func_type = func.FunctionType.from_lists([builtin.i32], [builtin.i32])

with func.FuncOp("double", func_type).body as func_body:
    arg = func_body.block.args[0]
```

```
18    two = arith.Constant.from_int_and_width(2, 32)
19    result = arith.Muli(arg, two)
20    func.Return(result)
21
22 # Create module
23 module = builtin.ModuleOp([func_body.owner])
24
25 # Print the IR
26 printer = Printer()
27 printer.print_op(module)
```

Listing 33: hello$_c$ompiler.py

# C    Appendix C: Glossary

**AST** Abstract Syntax Tree - Tree representation of source code structure

**Basic Block** Sequence of operations with single entry and exit

**CFG** Control Flow Graph - Graph showing possible execution paths

**Dataflow** Analysis tracking how data moves through program

**Dialect** Collection of related operations and types in xDSL

**Dominance** Relationship where one operation must execute before another

**IR** Intermediate Representation - Compiler's internal program representation

**Lowering** Transforming high-level operations to simpler ones

**Operation** Fundamental unit of computation in xDSL

**Operand** Input value to an operation

**Pass** Transformation or analysis applied to IR

**Pattern** Template for matching and replacing IR structures

**Region** Container for blocks in xDSL (e.g., function body, loop body)

**Result** Output value from an operation

**SSA** Single Static Assignment - Form where each variable assigned once

**Type** Description of data format and constraints

**Use-Def Chain** Links between value definitions and uses

**Value** Data flowing between operations (operands and results)

**Verification** Checking that IR is well-formed and valid