# MLIR Tutorial

Jaeho Lee

Corelab @ Yonsei University

# Course Information

**Introduction to MLIR**

**2 Chapters 2 Practices**

**Based on LLVM 15.0.0**

# Course Outline

- Basic Concept of MLIR
  - Goal: Learn the basic concept of LLVM IR

- Key components of MLIR
  - Goal: Learn the key components and essentials of MLIR

# **List of Practices**

- Basic Concept of MLIR


- Key components of MLIR
  - Practice 1: First Compilation
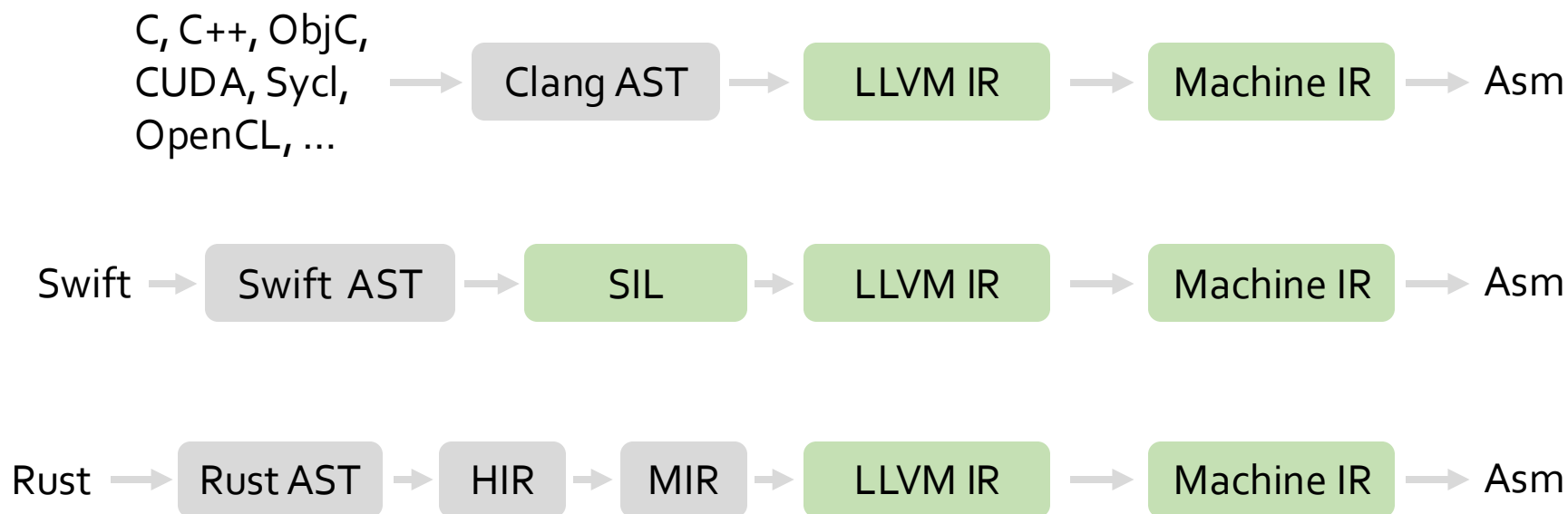  - Practice 2: Create transform pass

# Basic Concept of MLIR

# MLIR: Concept

- **M**ulti-**L**evel **I**ntermediate **R**epresentation
- Compiler Infrastructure
  - Domain-specific intermediate representation
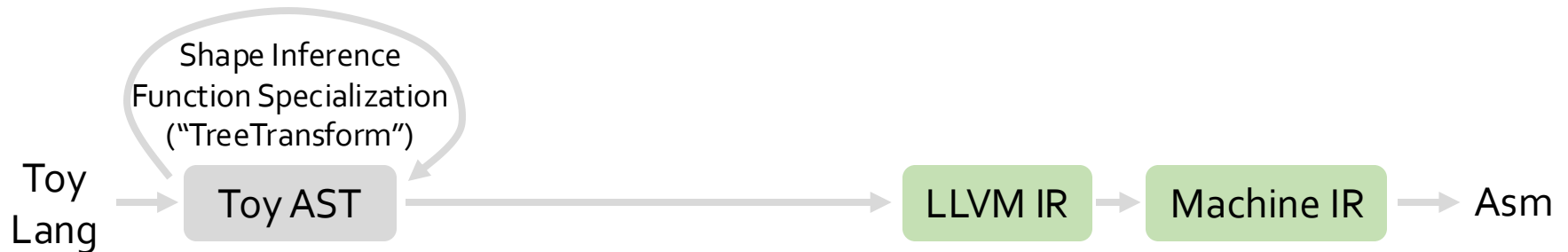  - High-level optimizations and portability

# Why we need MLIR rather LLVM
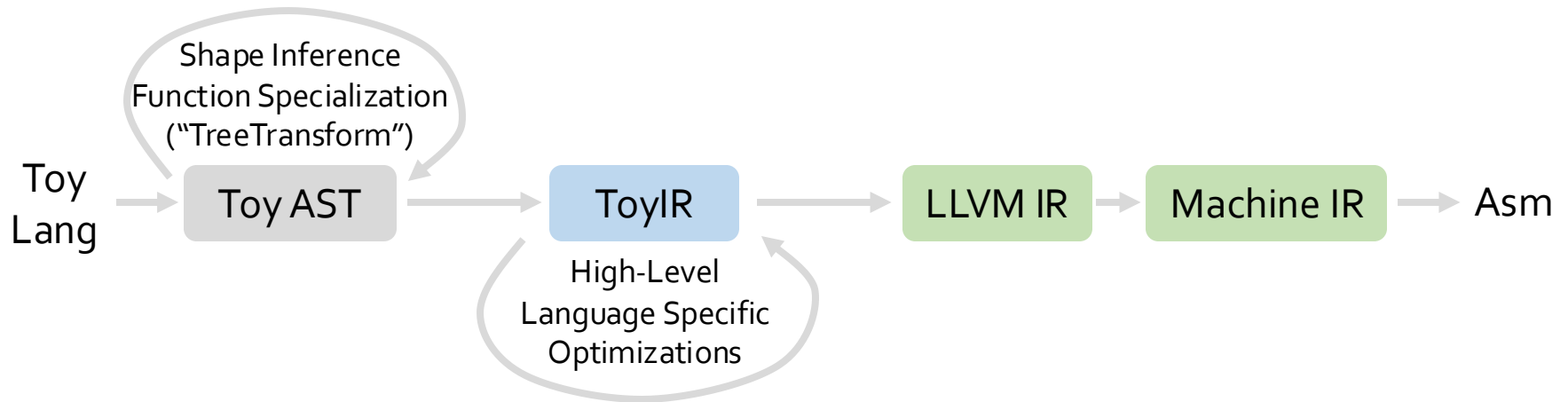
- Existing Successful Compilation Models

C, C++, ObjC, CUDA, Sycl, OpenCL, … → Clang AST → LLVM IR → Machine IR → Asm

Swift → Swift AST → SIL → LLVM IR → Machine IR → Asm

Rust → Rust AST → HIR → MIR → LLVM IR → Machine IR → Asm

# Why we need MLIR rather LLVM

- **The Toy Compiler:** the "Simpler" Approach of Clang



Shape Inference
Function Specialization
("TreeTransform")

Toy Lang → Toy AST → LLVM IR → Machine IR → Asm

# Why we need MLIR rather LLVM

- **The Toy Compiler:** With Language Specific Optimizations

Shape Inference
Function Specialization
("TreeTransform")

Toy
Lang → Toy AST → ToyIR → LLVM IR → Machine IR → Asm
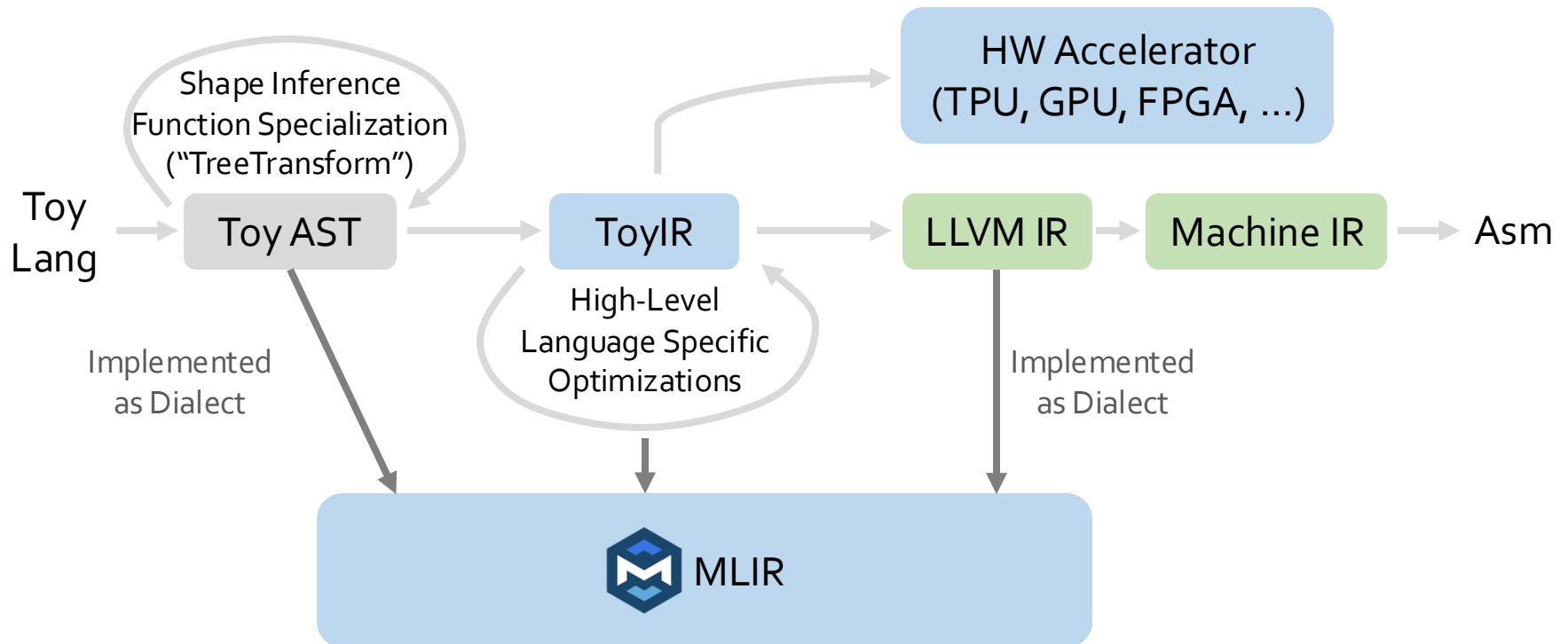
High-Level
Language Specific
Optimizations

# Why we need MLIR rather LLVM

- **The Toy Compiler:** Compilers in a Heterogenous World

# Why we need MLIR rather LLVM
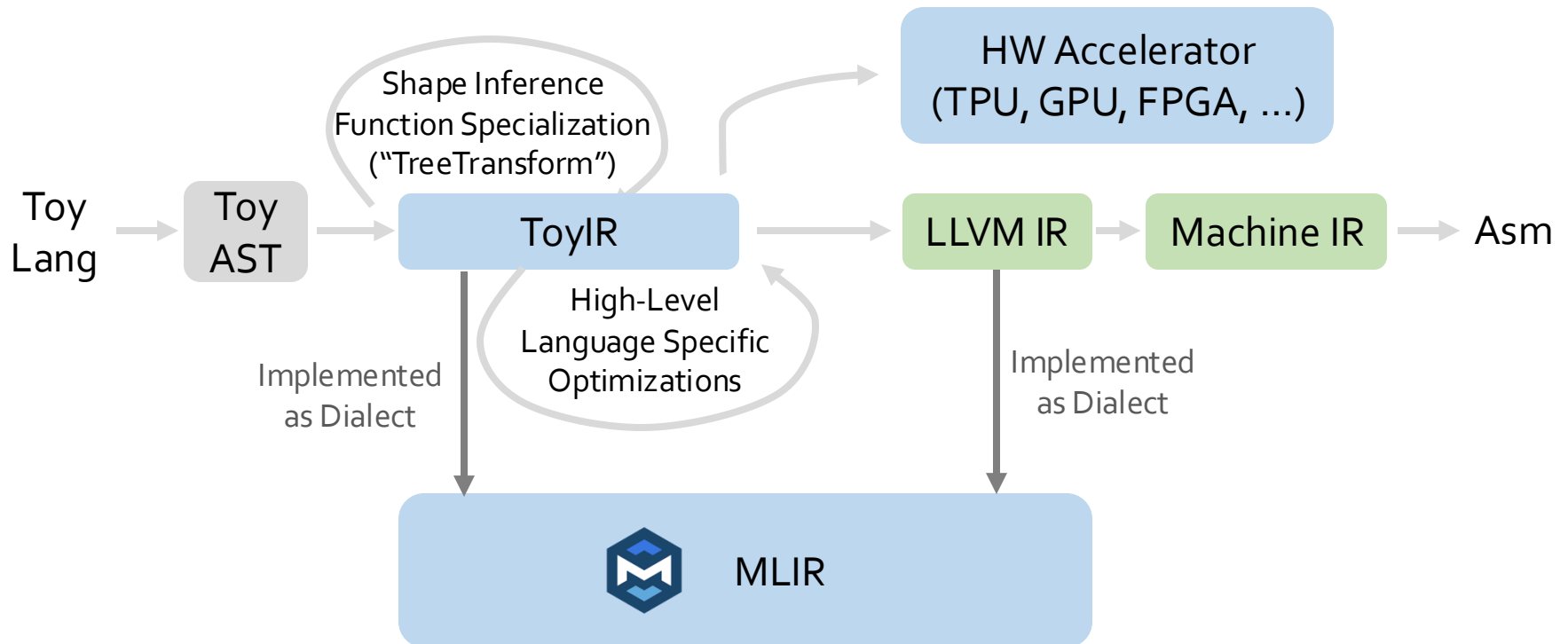
- **It's all about Dialects**

Shape Inference
Function Specialization
("TreeTransform")

Toy
Lang

Toy AST

ToyIR

HW Accelerator
(TPU, GPU, FPGA, ...)

LLVM IR

Machine IR

Asm

Implemented
as Dialect

High-Level
Language Specific
Optimizations

Implemented
as Dialect

MLIR

# Why we need MLIR rather LLVM

• **Dialect for Toy IR:** still <u>flexible enough</u> to perform shape inference and some high-level optimizations

# MLIR design principles

**Little built-in, everything customizable**

- fully customizable IR
- express many different abstractions with Common abstractions
  - machine learning graphs
  - mathematical abstractions
  - instruction-level intermediate representations
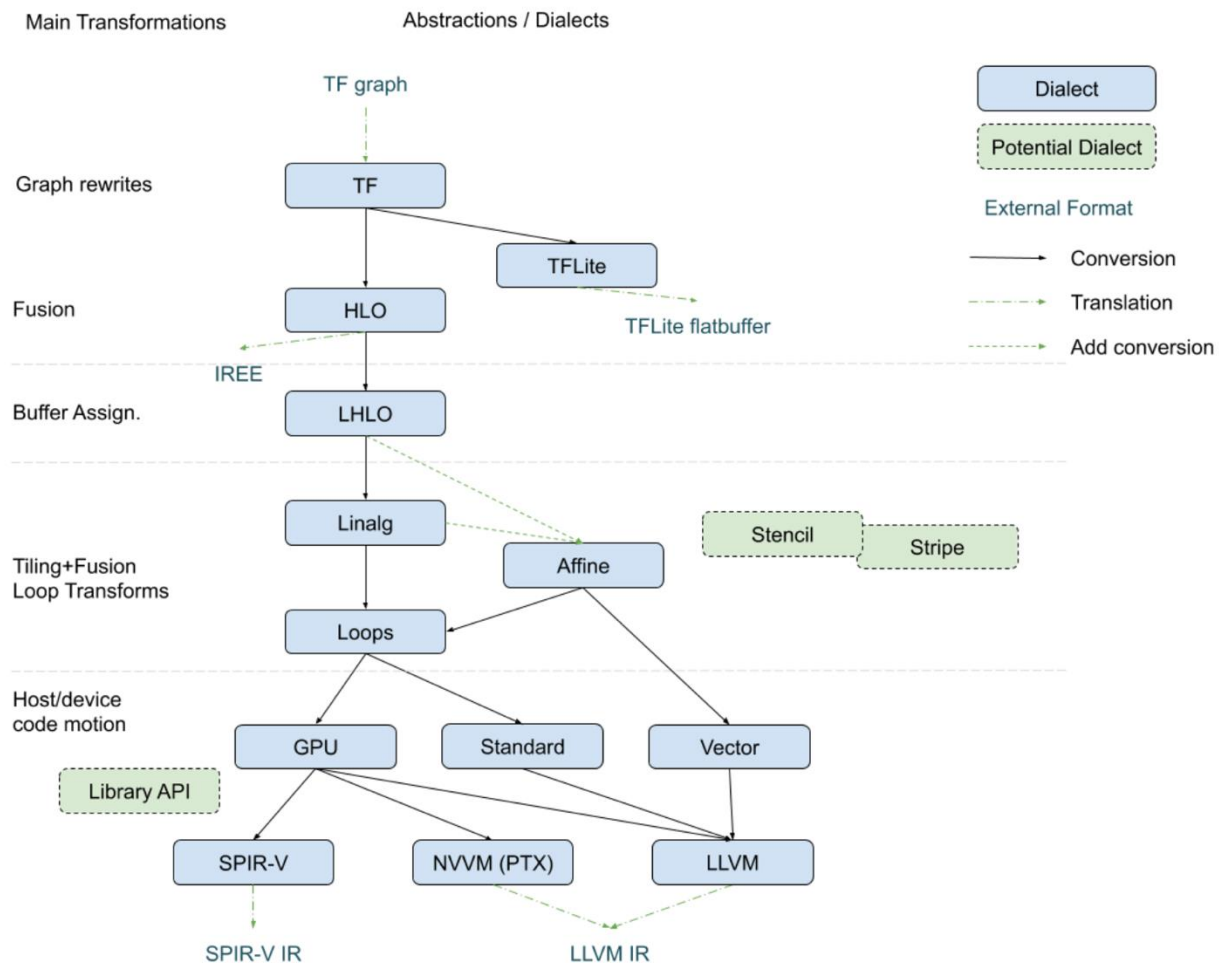  - etc

**Progressive Lowering**

- lowering in small steps along multiple abstraction levels
- using multiple levels of abstractions
  - support variety of platforms and programming models
  - flexible design of compilation pipeline

# Key components of MLIR

# MLIR Terminology

- Dialect
  - Intermediate representation of specific abstraction level

- Conversion
  - Compile to different dialects

- Transform
  - Optimization inside such dialect

- Translate
  - converts operations of dialect into some other language (specifically LLVM)
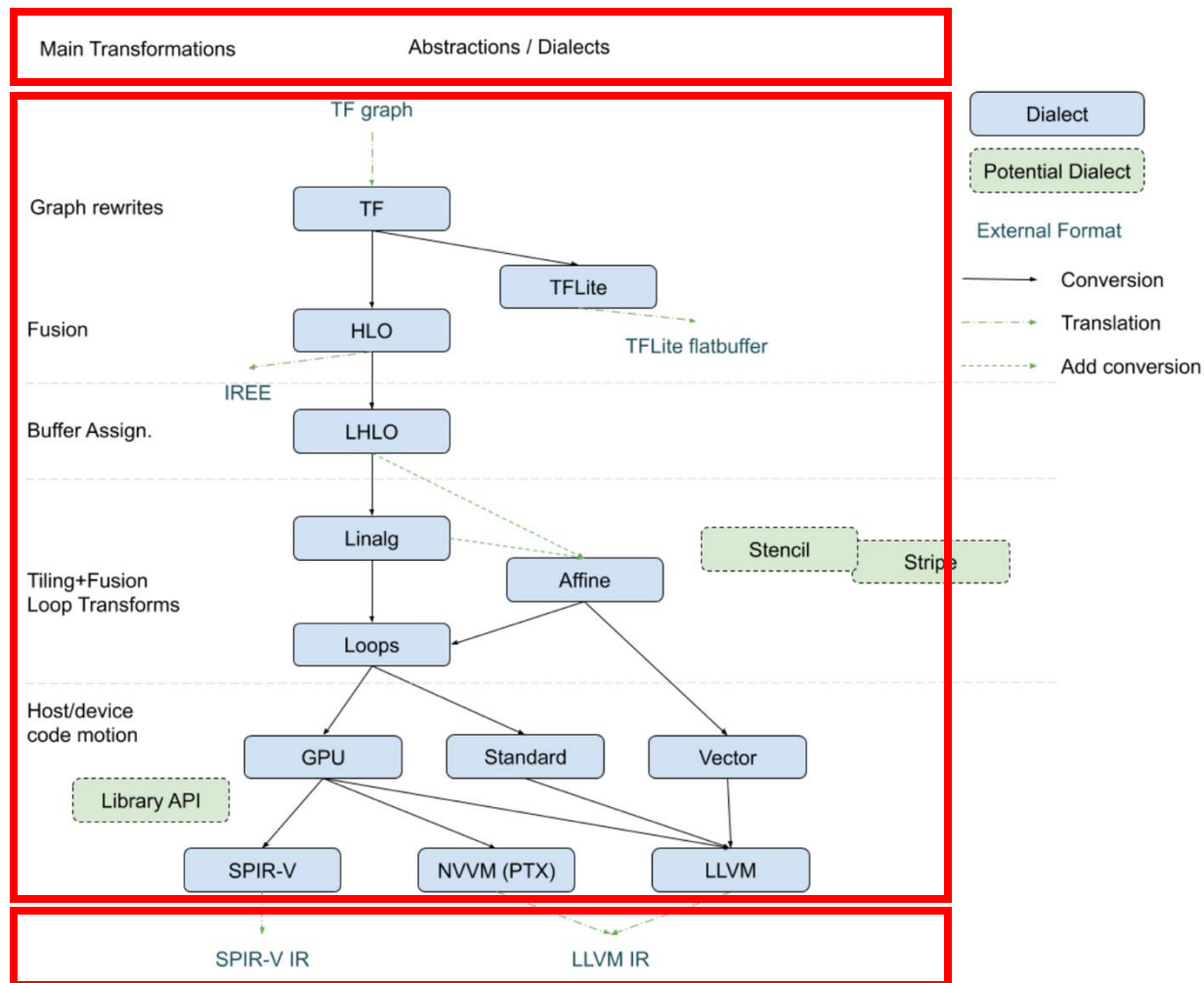
# Structure of MLIR

# Structure of MLIR

Frontend

Midend

Backend

# Let's set Docker environment

- run_docker.sh

docker run -td --workdir /root --name 나만의이름 --gpus all corelabyonsei/mlir-tutorial:latest
docker start 나만의이름

- exec_docker.sh

docker exec -it 나만의이름 bash

# Practice 1: First Compilation

- Goal
  - Learn how to use MLIR opt

- Steps
  1) Write a simple high-level language program (**example.toy**)

```
1 def multiply_transpose(a, b) {
2     return transpose(a) * transpose(b);
3 }
4
5 def main() {
6     var a<2, 3> = [[1, 2, 3], [4, 5, 6]];
7     var b<2, 3> = [1, 2, 3, 4, 5, 6];
8     var c = multiply_transpose(a, b);
9     var d = multiply_transpose(b, a);
10    print(d);
11 }
```

# Practice 1: First Compilation

- Goal
  - Learn how to use MLIR opt

- Steps
  1) Write a simple high-level language program (**example.toy**)
  2) Generate **AST (ast)** with command **toyc-ch6**
  3) Generate **MLIR (mlir)** with command **toyc-ch6**
  4) Generate **LLVM dialect (mlir-llvm)** with command **toyc-ch6**
  5) Generate **LLVM IR (llvm)** with command **toyc-ch6**

- Command
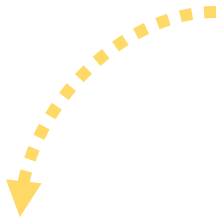  - Ex) `toyc-ch6` **`-emit=ast`** `example.toy`

# MLIR: Example

High-level lang

```
1  def multiply_transpose(a, b) {
2    return transpose(a) * transpose(b);
3  }
4
5  def main() {
6    var a<2, 3> = [[1, 2, 3], [4, 5, 6]];
7    var b<2, 3> = [1, 2, 3, 4, 5, 6];
8    var c = multiply_transpose(a, b);
9    var d = multiply_transpose(b, a);
10   print(d);
11 }
```

- Example code (-emit=milr)

MLIR

```
1  module {
2    toy.func private @multiply_transpose(%arg0: tensor<*xf64>, %arg1: tensor<*xf64>) ->
3      %0 = toy.transpose(%arg0 : tensor<*xf64>) to tensor<*xf64>
4      %1 = toy.transpose(%arg1 : tensor<*xf64>) to tensor<*xf64>
5      %2 = toy.mul %0, %1 : tensor<*xf64>
6      toy.return %2 : tensor<*xf64>
7    }
8    toy.func @main() {
9      %0 = toy.constant dense<[[1.000000e+00, 2.000000e+00, 3.000000e+00], [4.000000e+0
10     %1 = toy.reshape(%0 : tensor<2x3xf64>) to tensor<2x3xf64>
11     %2 = toy.constant dense<[1.000000e+00, 2.000000e+00, 3.000000e+00, 4.000000e+00,
12     %3 = toy.reshape(%2 : tensor<6xf64>) to tensor<2x3xf64>
13     %4 = toy.generic_call @multiply_transpose(%1, %3) : (tensor<2x3xf64>, tensor<2x3
14     %5 = toy.generic_call @multiply_transpose(%3, %1) : (tensor<2x3xf64>, tensor<2x3
15     toy.print %5 : tensor<*xf64>
16     toy.return
17   }
18 }
```

# MLIR: Example

High-level lang

```
1  def multiply_transpose(a, b) {
2    return transpose(a) * transpose(b);
3  }
4
5  def main() {
6    var a<2, 3> = [[1, 2, 3], [4, 5, 6]];
7    var b<2, 3> = [1, 2, 3, 4, 5, 6];
8    var c = multiply_transpose(a, b);
9    var d = multiply_transpose(b, a);
10   print(d);
11 }
```

- Example code (-emit=milr)

Operation

Region

```
1  module {
2    toy.func private @multiply_transpose(%arg0: tensor<*xf64>, %arg1: tensor<*xf64>) ->
3      %0 = toy.transpose(%arg0 : tensor<*xf64>) to tensor<*xf64>
4      %1 = toy.transpose(%arg1 : tensor<*xf64>) to tensor<*xf64>
5      %2 = toy.mul %0, %1 : tensor<*xf64>
6      toy.return %2 : tensor<*xf64>
7    }
8    toy.func @main() {
9      %0 = toy.constant dense<[[1.000000e+00, 2.000000e+00, 3.000000e+00], [4.000000e+0
10     %1 = toy.reshape(%0 : tensor<2x3xf64>) to tensor<2x3xf64>
11     %2 = toy.constant dense<[1.000000e+00, 2.000000e+00, 3.000000e+00, 4.000000e+00,
12     %3 = toy.reshape(%2 : tensor<6xf64>) to tensor<2x3xf64>
13     %4 = toy.generic_call @multiply_transpose(%1, %3) : (tensor<2x3xf64>, tensor<2x3
14     %5 = toy.generic_call @multiply_transpose(%3, %1) : (tensor<2x3xf64>, tensor<2x3
15     toy.print %5 : tensor<*xf64>
16     toy.return
17   }
18 }
```
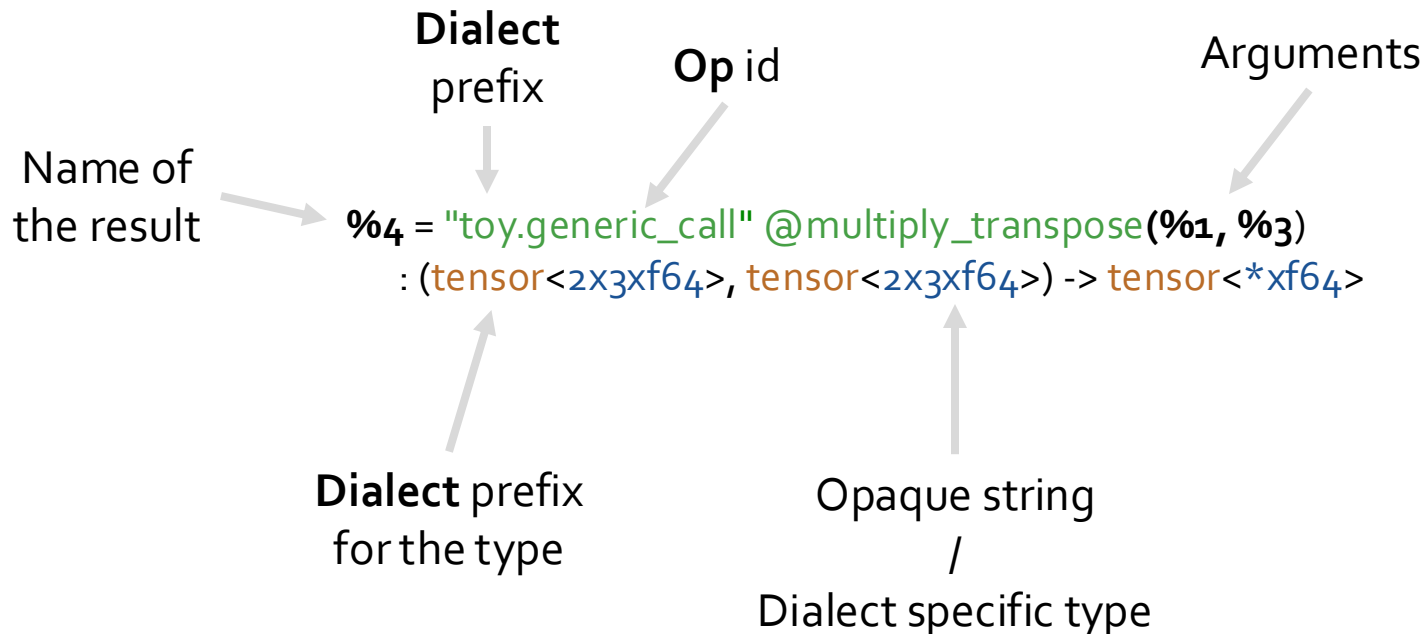
# MLIR: Components

- Operation, Region, Block

```
%results:2 = "d.operation"(%arg0, %arg1) ({
  // Regions belong to Ops and can have multiple blocks.        Region
  ^block(%argument: !d.type):                                   Block
    // Ops have function types (expressing mapping).
    %value = "nested.operation"() ({
      // Ops can contain nested regions.                        Region
      "d.op"() : () -> ()
    }) : () -> (!d.other_type)
    "consume.value"(%value) : (!d.other_type) -> ()
  ^other_block:                                                 Block
    "d.terminator"() [^block(%argument : !d.type)] : () -> ()
  })
// Ops can have a list of attributes.
{attribute="value" : !d.type} : () -> (!d.type, !d.other_type)
```
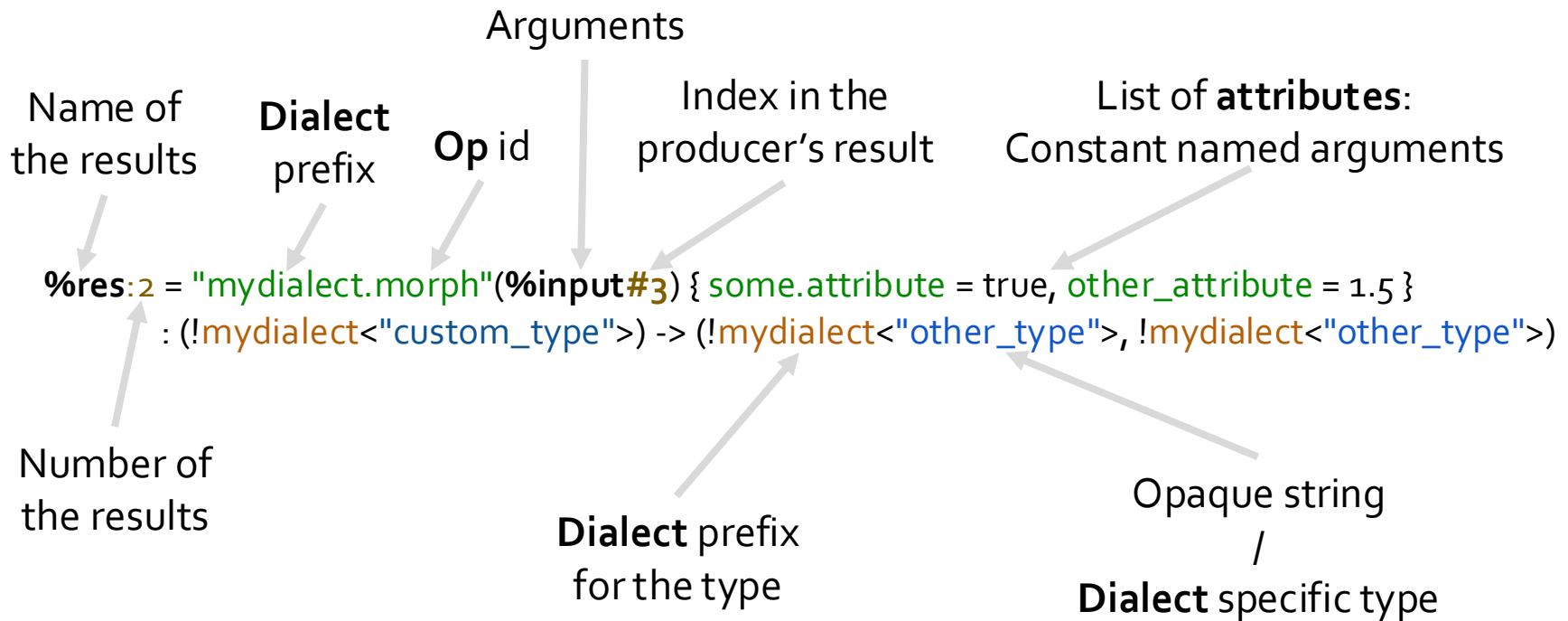
# Operations

- **Operation**! Not instruction (-emit=mlir)
  - No predefined set of instructions
  - Operations are like "opaque functions" to MLIR

**Dialect**
prefix

**Op** id

Arguments

Name of
the result

`%4` = "toy.generic_call" @multiply_transpose**(%1, %3)**
: (tensor<2x3xf64>, tensor<2x3xf64>) -> tensor<*xf64>

**Dialect** prefix
for the type

Opaque string
/
Dialect specific type

# Operations

Arguments

Name of
the results

**Dialect**
prefix

**Op** id

Index in the
producer's result

List of **attributes**:
Constant named arguments

```
%res:2 = "mydialect.morph"(%input#3) { some.attribute = true, other_attribute = 1.5 }
        : (!mydialect<"custom_type">) -> (!mydialect<"other_type">, !mydialect<"other_type">)
```

Number of
the results

**Dialect** prefix
for the type

Opaque string
/
**Dialect** specific type

# Regions

- Container that holds one or more blocks

```
1  module {
2    func @main() {
3      ...
4
5    ^bb1(%105: i64):    // 2 preds: ^bb0
6      %106 = llvm.icmp "slt" %105, %105 : i64
7      llvm.cond_br %106, ^bb2, ^bb6
8    ^bb2:    // pred: ^bb1
9      %107 = llvm.mlir.constant(0 : in
10     %108 = llvm.mlir.constant(2 : index) : i64
11     %109 = llvm.mlir.constant(1 : index) : i64
12     llvm.br ^bb3(%107 : i64)
13
14     ...
15     llvm.return
16   }
17 }
```

Region

Block 1

Block 2

# Regions: SSACFG regions

- SSACFG (Static Single Assignment Control Flow Graph)
  - Describes control flow between blocks

```
1  func @example(%arg: i32) {
2    ^bb0:
3      cond_br %condition, ^bb1, ^bb2
4
5    ^bb1:
6      // do something
7      br ^bb3
8
9    ^bb2:
10     // do something else
11     br ^bb3
12
13   ^bb3:
14     // merge point
15 }
```

# Regions: Graph Regions

- Do not necessarily require control flow between blocks
- Used in representations like Data Flow Graphs or Computational Graphs

```
1  // A custom operation that encapsulates a Graph Region
2  graph_op {
3    // Entry Block (No explicit control flow from here)
4    ^bb0:
5      // Operations without control-flow dependencies
6      %result1 = some_op1
7      %result2 = some_op2
8      graph_terminator
9  }
```

# Blocks

- A Block is a list of operations
- Instructions inside the block are executed in order

```
1  module {
2    func @main() {
3
4
5    ^bb1(%105: i64):   // 2 preds: ^bb0, ^bb5
6      %106 = llvm.icmp "slt" %105, %103 : i64
7      llvm.cond_br %106, ^bb2, ^bb6
8    ^bb2:   // pred: ^bb1
9      %107 = llvm.mlir.constant(0 : index) : i64
10     %108 = llvm.mlir.constant(2 : index) : i64
11     %109 = llvm.mlir.constant(1 : index) : i64
12     llvm.br ^bb3(%107 : i64)
13
14     ...
15     llvm.return
16   }
17 }
```

Block 1

Block 2

# Blocks: Single Static Assignment (SSA)

- Every value must be defined **only once**
  - In other words, every value has a **single** definition

```
%x = 1 + 2
%x = %x + 3
```
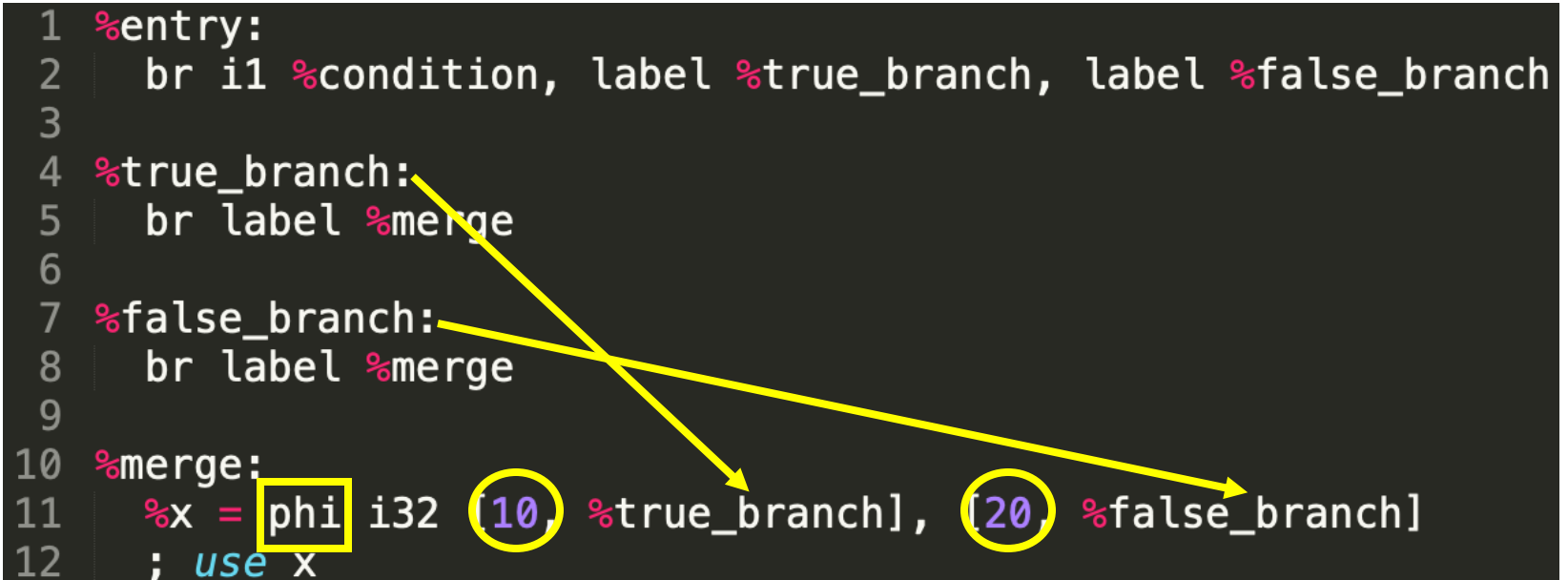❌
```
%x = 1 + 2
%y = %x + 3
```
⭕

- Facilitate program analyses
  - Liveness Analysis: From **DEF** to last **USE**
  - Constant Propagation: If **DEF** is constant, then **USE** is also constant

# Blocks: Block argument

• In LLVM, there is phi-node

```
1  int x;
2  if (condition) {
3      x = 10;
4  } else {
5      x = 20;
6  }
7  // use x
8
```

```
 1  %entry:
 2    br i1 %condition, label %true_branch, label %false_branch
 3
 4  %true_branch:
 5    br label %merge
 6
 7  %false_branch:
 8    br label %merge
 9
10  %merge:
11    %x = phi i32 [10, %true_branch], [20, %false_branch]
12    ; use x
```

# Blocks: Block argument

```
1  int x;
2  if (condition) {
3      x = 10;
4  } else {
5      x = 20;
6  }
7  // use x
8
```

- In LLVM, there is phi-node

```
1  %entry:
2    br i1 %condition, label %true_branch, label %false_branch
3
4  %true_branch:
5    br label %merge
6
7  %false_branch:
8    br label %merge
9
10 %merge:
11   %x = phi i32 [10, %true_branch], [20, %false_branch]
12   ; use x
```

Too low-level

# Blocks: Block argument

```
1  int x;
2  if (condition) {
3      x = 10;
4  } else {
5      x = 20;
6  }
7  // use x
8
```

- But in MLIR, **Block Argument !**

```
1  func @example(%condition: i1) {
2    cond_br %condition, ^bb1, ^bb2
3
4  ^bb1:
5    br ^bb3(10 : i32)
6
7  ^bb2:
8    br ^bb3(20 : i32)
9
10 ^bb3(%x: i32):
11   // use x
12 }
```
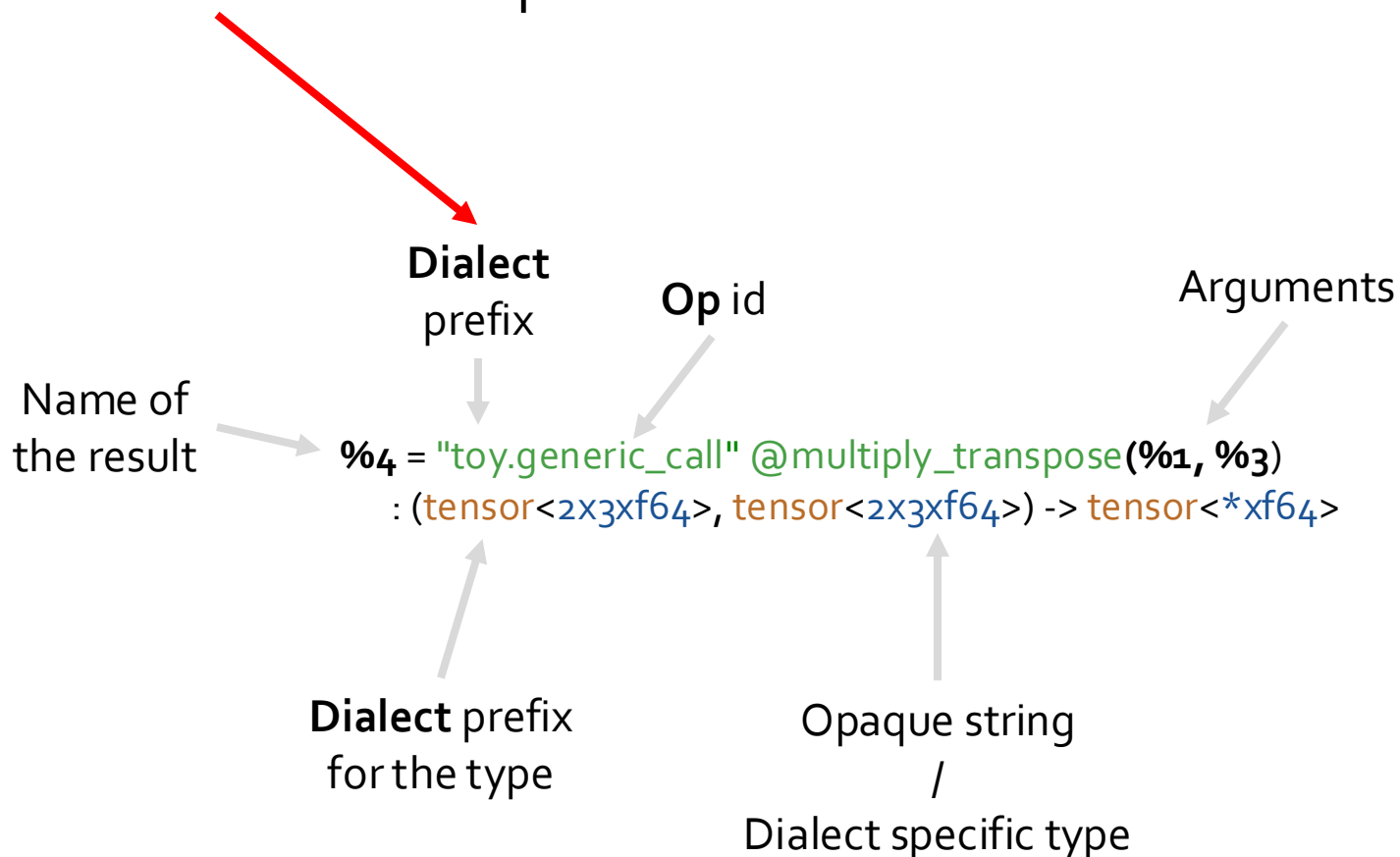
# Blocks: Block argument

```
1   int x;
2   if (condition) {
3       x = 10;
4   } else {
5       x = 20;
6   }
7   // use x
8
```

- But in MLIR, **Block Argument !**

```
1   func @example(%condition: i1) {
2     cond_br %condition, ^bb1, ^bb2
3
4   ^bb1:
5     br ^bb3(10 : i32)
6
7   ^bb2:
8     br ^bb3(20 : i32)
9
10  ^bb3(%x: i32):
11    // use x
12  }
```

similar to using parameters in real programming languages

# Now go back to Operations

- Dialects are the primitive unit of MLIR

**Dialect** prefix

**Op** id

Arguments

Name of the result

%4 = "toy.generic_call" @multiply_transpose**(%1, %3)**
: (tensor<2x3xf64>, tensor<2x3xf64>) -> tensor<*xf64>

**Dialect** prefix for the type

Opaque string / Dialect specific type

# Create Toy Dialect: Directory setup

- {PROJECT_ROOT}**/include/Dialect/**Toy
  - for public include files | *.td, *.h, *.hpp
- {PROJECT_ROOT}**/lib/Dialect/**Toy
  - for sources
- {PROJECT_ROOT}**/lib/Dialect /**Toy**/IR**
  - for operations | *.cpp
- {PROJECT_ROOT}**/lib/Dialect/**Toy**/Transforms**
  - for transforms
- {PROJECT_ROOT}**/test/Dialect/**Toy
  - for tests

# Create Toy Dialect: Cmake configuration

- TableGen Targets
  - {PROJECT_ROOT}/include/Dialect/{Dialect Name}/IR/CMakeLists.txt

```
add_mlir_dialect({Dialect_Name}Ops {tablegen_target_name})
add_mlir_doc({Dialect_Name}Ops {Dialect_Name}Dialect Dialects/ -gen-dialect-doc)
```

- Library Targets

```
add_mlir_dialect_library({dialect-library-target}
        DEPENDS
        <tablegen-targets>

        LINK_COMPONENTS
        Core

        LINK_LIBS PUBLIC
        <some-other-library>
)
```

# Create Toy Dialect: Dialect Implementation

- Dialect definition in ODS format

```
def Toy_Dialect : Dialect {
  let name = "toy";

  let summary = "...";

  let description = [{ ...}];

  let cppNamespace = "mlir::toy";
}

class Toy_Op<string mnemonic, list<OpTrait> traits = []> :
          Op<Toy_Dialect, mnemonic, traits>;
```

# Create Toy Dialect: Details

- Let's go to Docker

# Practice 2: Create transform pass

- Let's go to Docker

```
1  module {
2    func @main(%arg0: tensor<2x3xf64>, %arg1: tensor<2x3xf64>) {
3      %0 = "toy.add"(%arg0, %arg1) : (tensor<2x3xf64>, tensor<2x3xf64>) -> tensor<2x3xf64>
4      %1 = "toy.add"(%arg0, %arg1) : (tensor<2x3xf64>, tensor<2x3xf64>) -> tensor<2x3xf64>
5      toy.return %1 : tensor<2x3xf64>
6    }
7  }
```

- %0 is redundant again!
- Let's remove the "No Use" operations in Module