

Compiling ONNX Neural Network Models Using MLIR

Author: Tian Jin, Gheorghe-Teodor Bercea, Tung D. Le, Tong Chen, Gong Su,
Haruki Imai, Yasushi Negishi, Anh Leu, Kevin O'Brien, Kiyokuni Kawachiya,
Alexandre E. Eichenberger

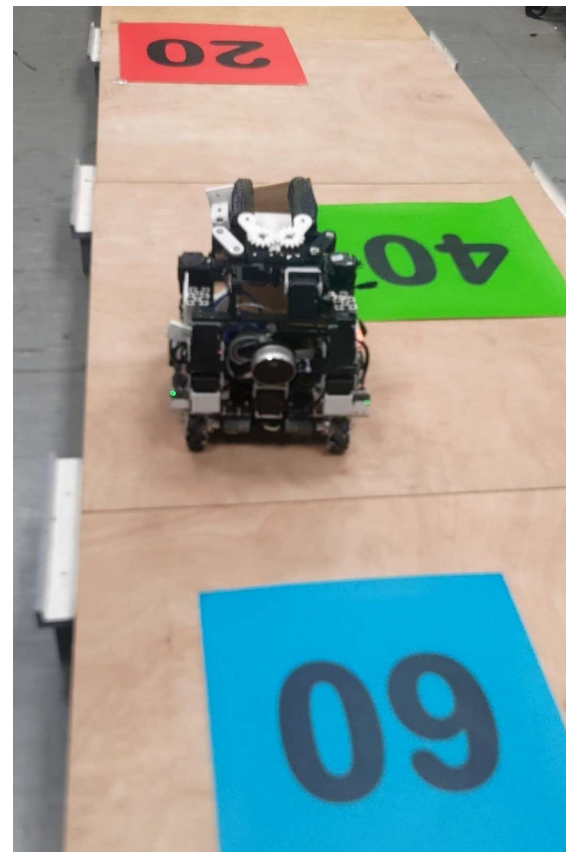
Advisor: 陳鵬升 教授
Presenter: 洪祐鈞

What Does this project do?

- Objective: Transform ONNX graphs into code
- Approach: Utilizes LLVM/MLIR compiler technology
- Focus: Inference
- Goal:
 - Implementation with stand-alone runtime support.
 - Integration to other MLIR compilers.

What is ONNX?

- ONNX (Open Neural Network Exchange)
- A common format for deep learning model.
- Supported by Pytorch, TensorFlow, Caffe,...
 - `torch.onnx.export(model, example_input, onnx_path)`
- Real life example: TensorRt
 - Deploy onnx model to NVIDIA GPUs.
 - Optimization and acceleration on inference:
 - Kernel auto-tuning.
 - Precision quantization.



What is ONNX?

- LeakyRelu Operation Example:

Listing 1: ONNX model for LeakyRelu operator (printed using 'protoc' command).

```
1 ir_version: 3
2 producer_name: "backend-test"
3 graph {
4   node {
5     input: "x"
6     output: "y"
7     op_type: "LeakyRelu"
8     attribute {
9       name: "alpha"
10      f: 0.1
11      type: FLOAT
12    }
13  }
14  name: "test_leakyrelu"
15  input {
16    name: "x"
17    type {
18      tensor_type {
19        elem_type: 1
20        shape {
21          dim {
22            dim_value: 3
23          }
24          dim {
25            dim_value: 4
26          }
27          dim {
28            dim_value: 5
29          }

```

```
34   output {
35     name: "y"
36     type {
37       tensor_type {
38         elem_type: 1
39         shape {
40           dim {
41             dim_value: 3
42           }
43           dim {
44             dim_value: 4
45           }
46           dim {
47             dim_value: 5
48           }
49         }
50       }
51     }
52   }
53 }
54 opset_import {
55   version: 9
56 }
```

What is ONNX?

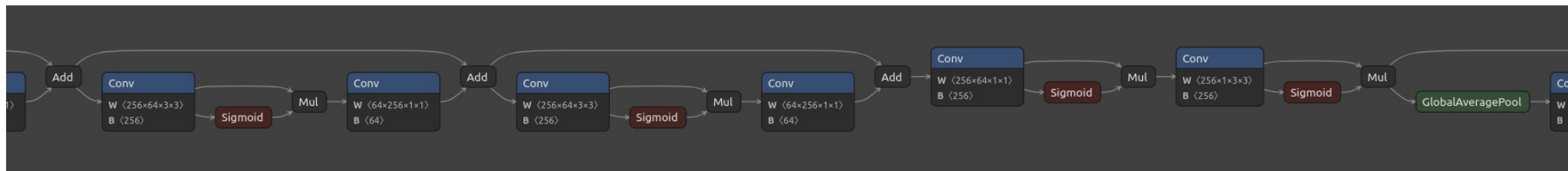
- Conv Operation Example:

```
onnx_model = onnx.load(onnx_file_path)
print(onnx_model)
```

```
ir_version: 8
producer_name: "pytorch"
producer_version: "2.1.1"
graph {
  node {
    input: "input.1"
    input: "onnx::Conv_1396"
    input: "onnx::Conv_1397"
    output: "/features/features.0/features.0.0/Conv_output_0"
    name: "/features/features.0/features.0.0/Conv"
    op_type: "Conv"
    attribute {
      name: "dilations"
      ints: 1
      ints: 1
      type: INTS
    }
    attribute {
      name: "group"
      i: 1
      type: INT
    }
  }
  attribute {
    name: "kernel_shape"
    ints: 3
    ints: 3
    type: INTS
  }
  attribute {
    name: "pads"
    ints: 1
    ints: 1
    ints: 1
    type: INTS
  }
  attribute {
    name: "strides"
    ints: 2
    ints: 2
    type: INTS
  }
}
```

What is ONNX?

- EfficientNet v1, visualized by Netron.
- Input, output, dimension, weight, bias, operation,...



What is MLIR?

- MLIR (Multi-Level Intermediate Representation)
- Motivation:
 - Various kind of hardware for domain specific accelerator. (DSA)
 - Historically, there was scalar and static vector types, now MLIR has tensor type.
 - Lowering the development cost of domain specific compiler.
 - Devs usually build their own high-level IR before going to lower level.
- Chris Lattner: “LLVM is a subset of MLIR”

What is MLIR?

- MLIR Operation (Op):
 - Op is central semantic unit of mlir.
 - Instruction, functions, modules,...
 - Encourage users to define custom Ops
 - Ops contains a list of regions.
 - Regions contain list of blocks.
 - In a sense of control flow graphs.
 - Blocks can contain list of ops
 - Allow recursive structures.

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

Figure 3: Operation (Op) is a main entity in MLIR; operations contain a list of regions, regions contain a list of blocks, blocks contains a list of ops, enabling recursive structures

What is MLIR?

- MLIR Dialects:
 - A group of operations into a namespace.
- MLIR Optimization Passes:
 - Can be expressed to DRRs (Declarative rewriting Rules)
 - Tablegen, C++ code
 - Conversion:
 - KrnlToAffine, KrnlToLLVM, ONNXToKrnl, ONNXToTOSA, ONNXToHLO
 - Transformation:
 - Decompose, shape inference, ConvOpt
 - Translation: Involve external representation.

What is MLIR?

- MLIR ecosystem:
 - Tensorflow: where MLIR originated
 - mhlo: Part of TF, Dynamic scaled XLA (for accelerating linear algebra)
 - tfirt: TF runtime
 - torch-mlir: connecting pytorch and mlir ecosystem
 - cirt: Hardware/Software Co-design
 - iree: Deep Learning E2E compiler. Kinda like TVM

What does the project provide?

- ONNX Dialect that can be integrated to other projects.
- Compiler interfaces:
 - Lowers ONNX graphs to MLIR, LLVM, C/C++, Python and Java.
 - Graph-to-code transformation, stand-alone runtime environment.
- Target:
 - Generic CPUs. (Linux, Windows, macOS)
 - IBM's Telum AI accelerator. (IBM z16)
 - A server that can handle huge amount of queries.

What does the project provide?

- What does it means?
 - Easy to write optimization of CPU and custom accelerators
 - MLIR
 - Easy to deploy
 - No libraries need, only stand-alone driver and runtime support.
 - Integrate ONNX to MLIR eco-system

Motivation:

- Many deep learning frameworks use their own optimized libraries for specific accelerators during inference.
- Drawbacks:
 - Number of supported accelerator of model is limited to libraries.
 - Needed to install various libraries.
 - The implementation varies for the same operation.
- Proposed Solution: Developing a compiler that rewrites a trained model to native code for a target hardware.

Architecture:

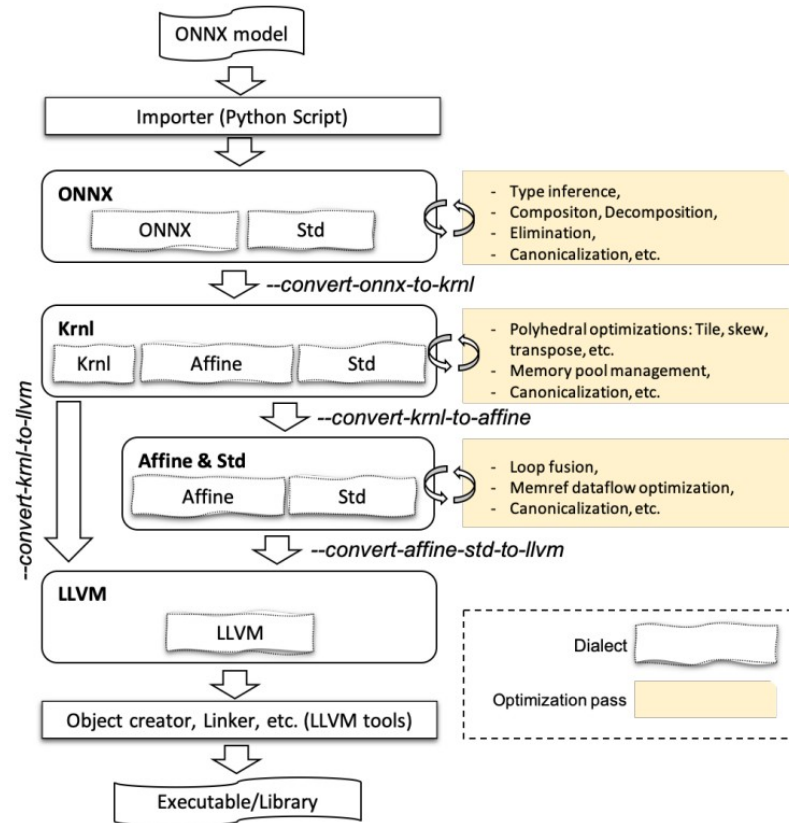


Fig. 2: Architecture of onnx-mlir. Names prefixed with '--' are passes.

Architecture:

- 4 abstraction layer:
- 5 main Dialect:
 - ONNX
 - Krnl
 - Std
 - Affine
 - LLVM

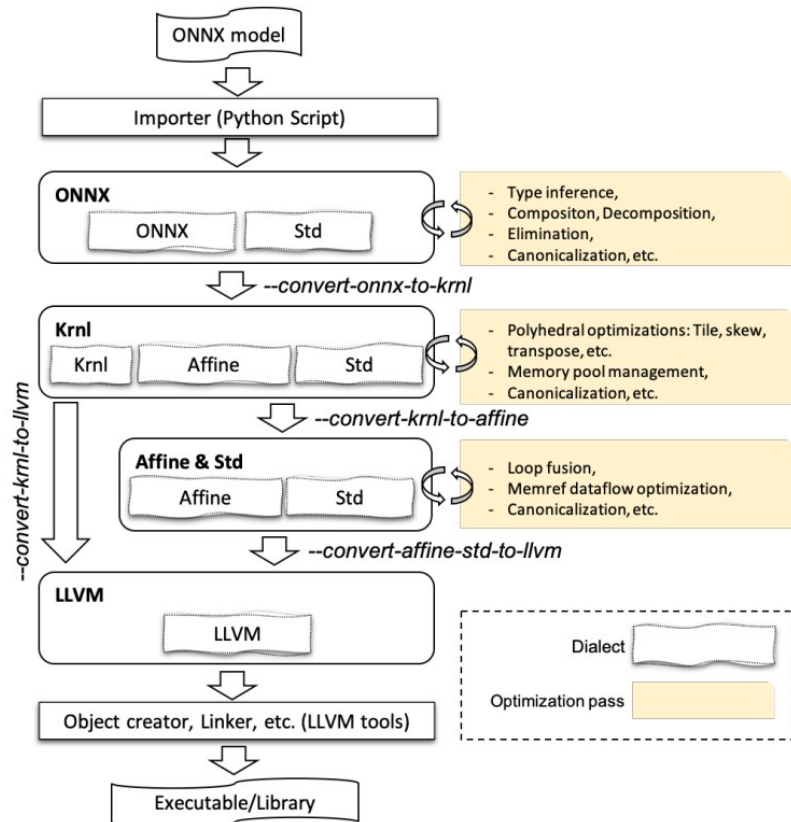


Fig. 2: Architecture of onnx-mlir. Names prefixed with '--' are passes.

Architecture:

- ONNX Dialect:
 - Represent ONNX with mlir
 - Tablegen, Generated by Python script
 - In MLIR, usually you write .td file yourself.
 - You can define the following in .td file:
 - Dialect
 - Operation
 - Attribute
 - Type

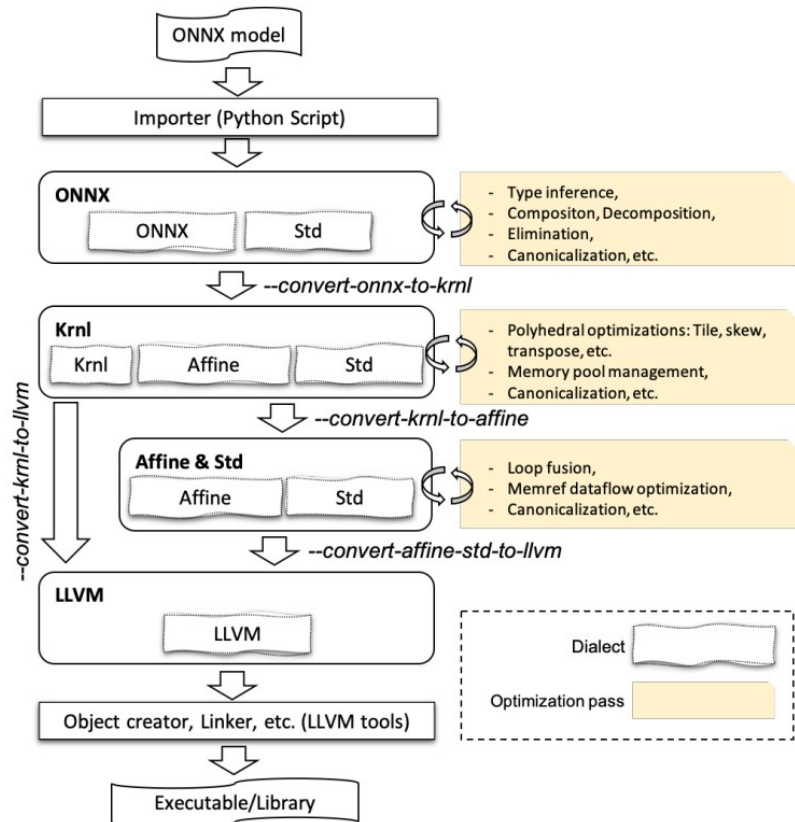


Fig. 2: Architecture of onnx-mlir. Names prefixed with '--' are passes.

Architecture

- ONNX Dialect:
 - Tablegen allow to represent information to human-readable form:
 - Arguments: input & attribute
 - Results: output
 - The ONNX and Tablegen definition are quite similar.

Listing 6: Tablegen-based definition for operation relu.

```
1 def ONNXLeakyReluOp:ONNX_Op<"LeakyRelu",  
2   [NoSideEffect, DeclareOpInterfaceMethods<ShapeInferenceOpInterface>]> {  
3   let summary = "ONNX LeakyRelu operation";  
4   let description = [{"LeakyRelu takes ... "}]  
5   let arguments = (ins AnyTypeOf<[TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$X, ←  
6     DefaultValuedAttr<F32Attr, "0.01">:$alpha);  
7   let results = (outs AnyTypeOf<[TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$Y);  
8   let extraClassDeclaration = [{ ... }];  
9 }
```

Architecture

- ONNX Dialect (LeakyRelu):
 - Real Example: (src/Dialect/ONNX/ONNXOPs.td.inc)
 - Generated by: (utils/gen_onnx_mlir.py)

```
def ONNXLeakyReluOp:ONNX_Op<"LeakyRelu",
  [Pure, DeclareOpInterfaceMethods<ShapeInferenceOpInterface>, DeclareOpInterfaceMethods<ShapeHelperOpInterface>]> {
  let summary = "ONNX LeakyRelu operation";
  let description = [{
    LeakyRelu takes input data (Tensor<T>) and an argument alpha, and produces one
    output data (Tensor<T>) where the function `f(x) = alpha * x for x < 0`,
    `f(x) = x for x >= 0`, is applied to the data tensor elementwise.
  }];
  let arguments = (ins AnyTypeOf<[TensorOf<[BF16]>, TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$X,
    DefaultValuedAttr<F32Attr, "0.01">:$alpha);
  let results = (outs AnyTypeOf<[TensorOf<[BF16]>, TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$Y);
  let extraClassDeclaration = [{
    static int getNumberOfOperands() {
      return 1;
    }
  }]
```

Architecture

- ONNX Dialect (Conv):
 - Real Example: (src/Dialect/ONNX/ONNXOPs.td.inc)

```
def ONNXConvOp:ONNX_Op<"Conv",
[Pure, DeclareOpInterfaceMethods<ShapeInferenceOpInterface>, DeclareOpInterfaceMethods<ShapeHelperOpInterface>]> {
  let summary = "ONNX Conv operation";
  let description = [{
    The convolution operator consumes an input tensor and a filter, and
    computes the output.
  }];
  let arguments = (ins AnyTypeOf<[TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$X,
    AnyTypeOf<[TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$W,
    AnyTypeOf<[TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>, NoneType]>:$B,
    DefaultValuedStrAttr<StrAttr, "NOTSET">:$auto_pad,
    OptionalAttr<I64ArrayAttr>:$dilations,
    DefaultValuedAttr<SI64Attr, "1">:$group,
    OptionalAttr<I64ArrayAttr>:$kernel_shape,
    OptionalAttr<I64ArrayAttr>:$pads,
    OptionalAttr<I64ArrayAttr>:$strides);
  let results = (outs AnyTypeOf<[TensorOf<[F16]>, TensorOf<[F32]>, TensorOf<[F64]>]>:$Y);
  let builders = [
    OpBuilder<(ins "Value":$X, "Value":$W, "Value":$B, "StringAttr":$auto_pad, "ArrayAttr":$dilations, "IntegerAttr":$group,
      "ArrayAttr":$kernel_shape, "ArrayAttr":$pads, "ArrayAttr":$strides), [{
        auto resultType = UnrankedTensorType::get(X.getType().cast<ShapedType>().getElementType());
        build($_builder, $_state, resultType, X, W, B, auto_pad, dilations, group, kernel_shape, pads, strides);
      }]>,
    OpBuilder<(ins "ValueRange":$operands, "ArrayRef<NamedAttribute>":$attributes), [{
        auto resultType = UnrankedTensorType::get(operands[0].getType().cast<ShapedType>().getElementType());
        build($_builder, $_state, {resultType}, operands, attributes);
      }]>
  ];
};
```

Architecture:

- krnl Dialect:
 - Bridge between onnx & affine dialect.
 - Affine dialect is ready to use.
 - Most computation are loop nest.
 - Represent loops in polyhedral model.
 - Providing optimization chances:
 - Parallelism
 - Cache locality.
 - SIMD, vectorization

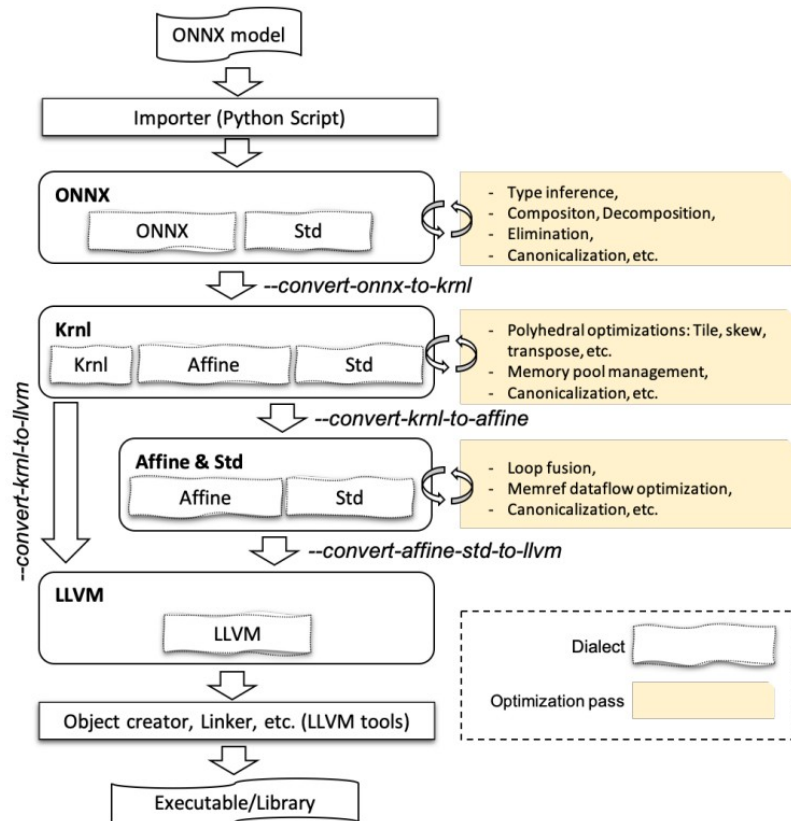


Fig. 2: Architecture of onnx-mlir. Names prefixed with '--' are passes.

Architecture:

- krnl Dialect:
 - At left, Imagine 2 10*10 array added to each other with same array.
 - Now we want to do it with tiling size 2, with a single loop.
 - krnl.iterate (schedule_loop) with (original_loop).
 - Separate program semantics and schedule.

```
1 %i1, %jj = krnl.define_loops 2
2 krnl.iterate(%i1, %jj) with (%i1 -> %i = 0
   to 10, %jj -> %j = 0 to 10) {
3   %foo = std.addi %i, %j : index
4 }
```

```
1 %i1 = krnl.define_loops 1
2 %ib, %il = krnl.block %i1 2 : ←
   (!krnl.loop)->(!krnl.loop, !krnl.loop)
3 krnl.iterate(%ib, %il) with (%i1 -> %i = 0
   to 10) {
4   %foo = std.addi %i, %i : index
5 }
```

Architecture:

- krnl Dialect:
 - --convert-krnl-to-affine pass generates optimized affine.for based loops.
 - The inner loop affine.for is iterating tile size 2.
 - Skew and permutation are applied similarly, and are composable.

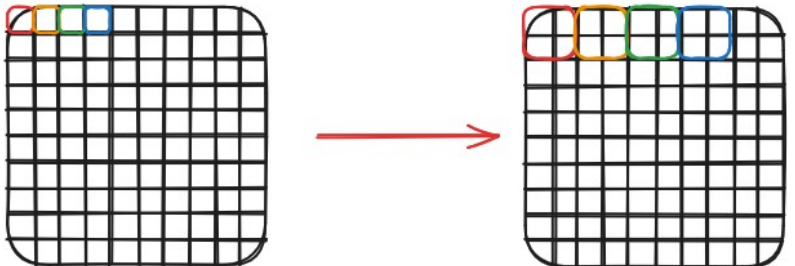
```
1 %ii = krnl.define_loops 1
2 %ib, %il = krnl.block %ii 2 : ↵
   (!krnl.loop)->(!krnl.loop, !krnl.loop)
3 krnl.iterate(%ib, %il) with (%ii -> %i = 0
   to 10) {
4   %foo = std.addi %i, %i : index
5 }
```

```
1 #map0 = affine_map<(d0) -> (d0)>
2 #map1 = affine_map<(d0) -> (d0 + 2)>
3 affine.for %arg0 = 0 to 10 step 2 {
4   affine.for %arg1 = #map0(%arg0) to
     #map1(%arg0) {
5     %0 = addi %arg1, %arg1 : index
6   }
7 }
```

Architecture:

- But, I still don't know what just happened?
 - Me, too.
 - This is my speculation.

```
for (int i = 0; i < 10; ++i){
    for (int j = 0; j < 10; ++j)
        result[i][j] = matrix1[i][j] + matrix2[i][j];
}
```



```
for (int i = 0; i < 10; i += 2) {
    for (int j = 0; j < 10; j += 2) {
        for (int ii = i; ii < i + 2; ++ii) {
            for (int jj = j; jj < j + 2; ++jj) {
                result[ii][jj] = matrix1[ii][jj] + matrix2[ii][jj];
            }
        }
    }
}
```

Architecture:

- Affine Dialect:
 - Mainly for:
 - Actual loop transformation.
 - Dependence analysis.
 - Ready-to-use.

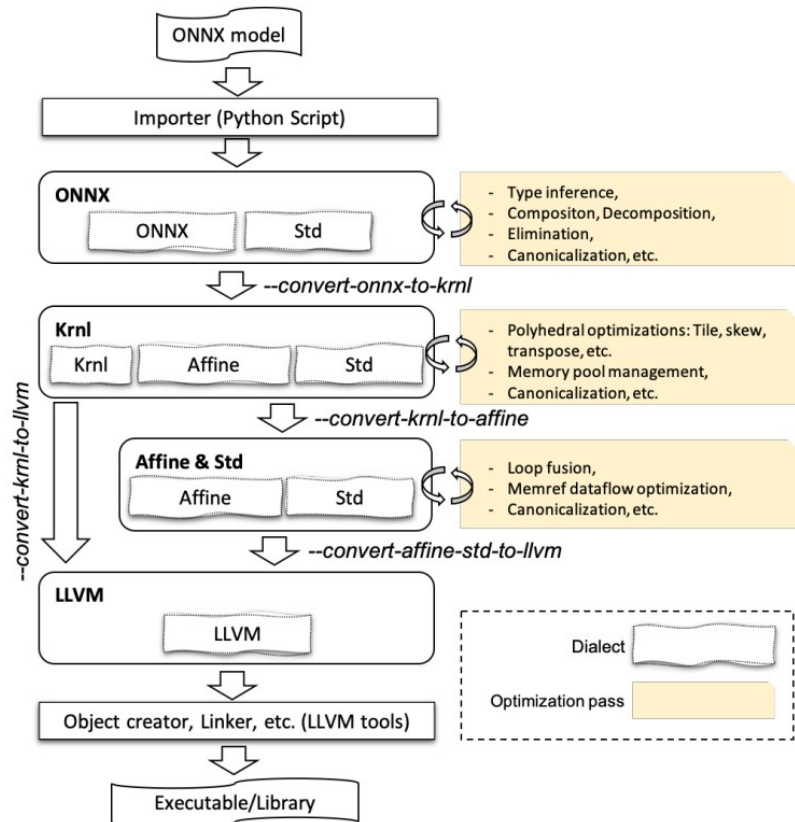


Fig. 2: Architecture of onnx-mlir. Names prefixed with '--' are passes.

Architecture:

- Std Dialect:
 - load, store, addi, addf, absf, call...
 - Obsolete
 - ops are moved to other dialects.

```
%0 = affine.load %arg0[%arg2, %arg3] : memref<3x2xf32>
%1 = affine.load %arg1[%arg2, %arg3] : memref<3x2xf32>
%2 = arith.addf %0, %1 : f32
affine.store %2, %alloc[%arg2, %arg3] : memref<3x2xf32>
```

```
%1 = affine.load %arg0[%arg2, %arg3, %arg4] : memref<3x4x5xf32>
%2 = affine.load %arg1[%arg2, %arg3, %arg4] : memref<3x4x5xf32>
%3 = std.addf %1, %2 : f32
affine.store %3, %0[%arg2, %arg3, %arg4] : memref<3x4x5xf32>
```

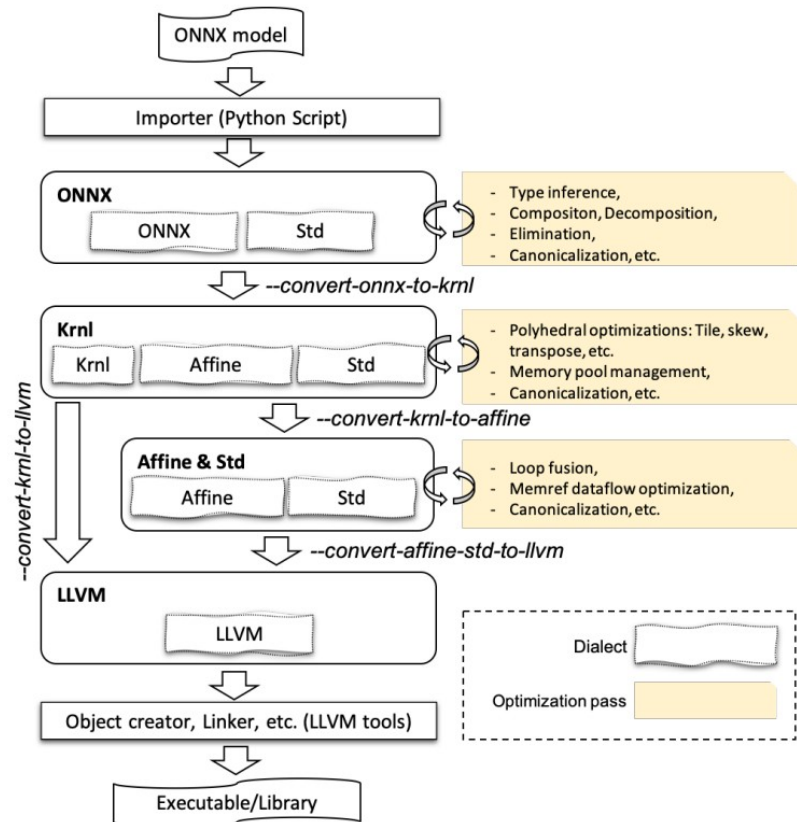


Fig. 2: Architecture of onnx-mlir. Names prefixed with '--' are passes.

Architecture:

- llvm Dialect:
 - Wrapping the LLVM IR into MLIR.
 - Generate bitcode.

```
%0 = llvm.mlir.constant(2 : i64) : i64
%1 = llvm.mlir.constant(0 : i64) : i64
%2 = llvm.mlir.constant(1 : i64) : i64
%3 = llvm.call @omTensororListGetOmtArray(%arg0) : (!llvm.ptr) -> !llvm.ptr
%4 = llvm.alloca %x2 x !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>) : (i64) -> !llvm.ptr
%5 = llvm.load %3 : !llvm.ptr -> !llvm.ptr
%6 = llvm.alloca %x2 x !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>) : (i64) -> !llvm.ptr
%7 = llvm.mlir.undef : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%8 = llvm.call @omTensororGetDataPtr(%5) : (!llvm.ptr) -> !llvm.ptr
%9 = llvm.insertvalue %8, %7[0] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%10 = llvm.insertvalue %8, %9[1] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%11 = llvm.insertvalue %1, %10[2] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%12 = llvm.call @omTensororGetShape(%5) : (!llvm.ptr) -> !llvm.ptr
%13 = llvm.call @omTensororGetStrides(%5) : (!llvm.ptr) -> !llvm.ptr
%14 = llvm.load %12 : !llvm.ptr -> i64
%15 = llvm.insertvalue %14, %11[3, 0] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%16 = llvm.load %13 : !llvm.ptr -> i64
%17 = llvm.insertvalue %16, %15[4, 0] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%18 = llvm.gemvmliteptr %12[1] : (!llvm.ptr) -> !llvm.ptr, i64
%19 = llvm.load %18 : !llvm.ptr -> i64
%20 = llvm.insertvalue %19, %17[3, 1] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%21 = llvm.gemvmliteptr %13[1] : (!llvm.ptr) -> !llvm.ptr, i64
%22 = llvm.load %21 : !llvm.ptr -> i64
%23 = llvm.insertvalue %22, %20[4, 1] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
llvm.store %23, %6 : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>), !llvm.ptr
%24 = llvm.gemvmliteptr %3[1] : (!llvm.ptr) -> !llvm.ptr, !llvm.ptr
%25 = llvm.load %24 : !llvm.ptr -> !llvm.ptr
%26 = llvm.alloca %x2 x !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>) : (i64) -> !llvm.ptr
%27 = llvm.mlir.undef : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%28 = llvm.call @omTensororGetDataPtr(%25) : (!llvm.ptr) -> !llvm.ptr
%29 = llvm.insertvalue %28, %27[0] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%30 = llvm.insertvalue %28, %29[1] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%31 = llvm.insertvalue %1, %30[2] : !llvm.struct<(ptr, ptr, i64, array<2 x i64>, array<2 x i64>>)
%32 = llvm.call @omTensororGetShape(%25) : (!llvm.ptr) -> !llvm.ptr
```

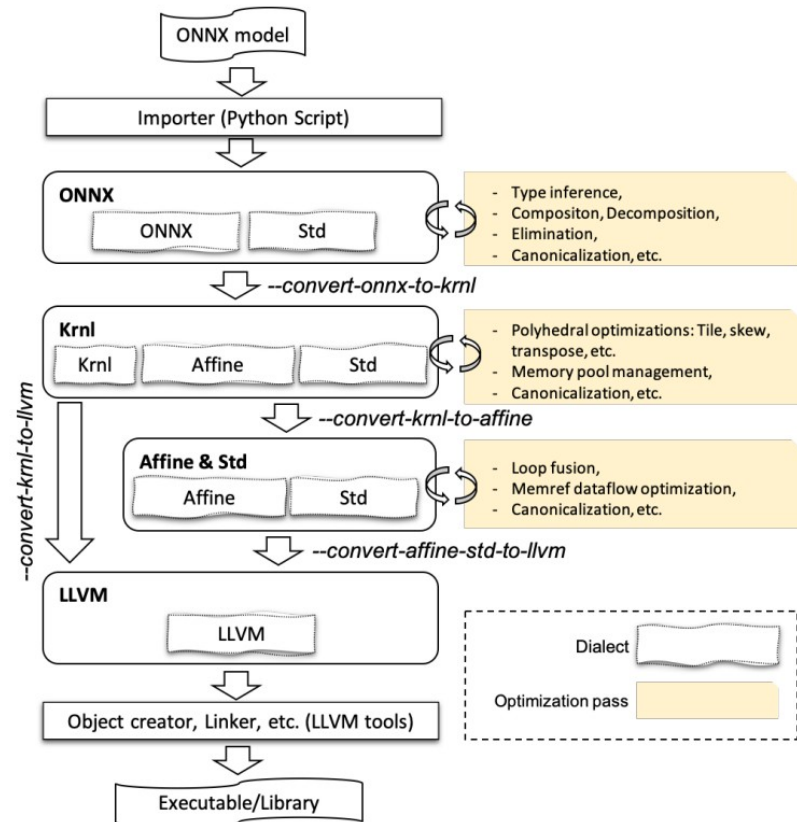


Fig. 2: Architecture of onnx-mlir. Names prefixed with ‘--’ are passes.

Optimization Passes:

- Optimization:
 - Operation decomposition
 - Shape inference
 - Graph rewriting
 - Constant propagation

Optimization Passes:

- Operation decomposition:
 - Breaking down complex operations into a sequence of simpler ones
 - e.g. ReduceL1 operation

```
1 def ReduceL1Pattern: Pat<
2   (ReduceL1Op $x, $axes, $keepdims),
3   (ReduceSumOp (AbsOp $x), $axes, $keepdims)
4 >;
```

src > Transform > ONNX >  Decompose.td

```
97  //====-----//
98  // ONNXReduceL1Op %X = ONNXReduceSumOp (ONNXAbsOp %X)
99  //====-----//
100 def ReduceL1OpPattern
101   : Pat<(ONNXReduceL1Op $oprnd, $axes, $keepdims, $noop),
102   (ONNXReduceSumOp(ONNXAbsOp $oprnd), $axes, $keepdims, $noop)>;
103
104 //====-----//
105 // ONNXReduceL2Op %X = ONNXSqrtOp (ONNXReduceSumSquareOp (%X))
106 //====-----//
107 def ReduceL2OpPattern
108   : Pat<(ONNXReduceL2Op $oprnd, $axes, $keepdims, $noop),
109   (ONNXSqrtOp(ONNXReduceSumSquareOp $oprnd, $axes, $keepdims, $noop))>;
```

Optimization Passes:

- Operation decomposition:
 - Breaking down complex operations into a sequence of simpler ones
 - e.g. ReduceL1 operation

L1-norm

$$\|x\|_1 := \sum_{i=1}^n |x_i|$$

- L1-norm 又稱為 Taxicab-norm (計程車幾何範數) 或 Manhattan-norm (曼哈頓距離範數)

L2-norm

$$\|x\|_2 := [\sum_{i=1}^n |x_i|^2]^{1/2} = \sqrt{|x_1|^2 + |x_2|^2 + \dots + |x_n|^2}$$

Optimization Passes:

- Shape inference:
 - Traverses all operations
 - Infers the shapes of tensors with unrank shapes (e.g. `tensor<*xf32>`)
 - Propagates the ranked shapes to consuming operations
 - e.g.: A output to B, A is `tensor<3*4>`, so B is `tensor<3*4>` too.
 - Terminates until all tensors have ranked shapes.

Listing 3: Operation add in onnx dialect, generated using importer.

```
1 module {
2   func @main_graph( %arg0:tensor<3x4x5xf32>, %arg1:tensor<3x4x5xf32>) -> tensor<*xf32> {
```

Listing 4: Operation add in krnl dialect, generated by applying passes `--shape-inference` and `--convert-onnx-to-krnl`.

```
1 module {
2   func @main_graph(%arg0: memref<3x4x5xf32>, %arg1: memref<3x4x5xf32>) -> memref<3x4x5xf32> {
```

Optimization Passes:

- Graph rewriting:
 - Rewriting rules are conveniently represented using DRRs
 - Examples: Fuse add and mul ops together

```
1 def MulAddToGemmPattern : Pat<
2   (AddOp (MatMulOp:$res $m1, $m2), $m3),
3   (GemmOp $m1, $m2, $m3),
4   [(HasOneUse $res)]
5 >;
```

src > Dialect > ONNX > Rewrite.td

```
1 // onnx.add(onnx.matmul(%X, %Y), %Z) = onnx.Gemm(%X, %Y, %Z)
276 def MulAddToGemmOptPattern : Pat<(ONNXAddOp (ONNXMatMulOp:$res $m1, $m2), $m3),
1   (ONNXGemmOp $m1, $m2, $m3, (GemmAlpha), (GemmBeta), (GemmTransA), (GemmTransB)),
2   [(HasOneUse $res), (HasRankOf<2> $m1), (HasRankOf<2> $m2)]>;
```

Optimization Passes:

- Graph rewriting:
 - Rewriting rules are conveniently represented using DRRs
 - Examples: remove identity op

```
1 def IdentityEliminationPattern : Pat<
2   (ONNXIdentityOp $arg),
3   (replaceWithValue $arg)
4 >;
```

src > Dialect > ONNX > Rewrite.td

```
405 // ONNX_Op (onnx.Identity (%X)) = ONNX_Op (%X)
1   def IdentityEliminationPattern : Pat<(ONNXIdentityOp $arg),
2     (replaceWithValue $arg)>;
```


Optimization Passes:

- Constant propagation:
 - Compute the constant together at compile-time.
 - If mixed constant and variable, normalize it:
 - e.g. (onnx.Add)
 - Utilize associativity, group variable and constant respectively.

$$(1) \quad c + x \Rightarrow x + c$$

$$(2) \quad (x + c_1) + c_2 \Rightarrow x + (c_1 + c_2)$$

$$(3) \quad (x + c) + y \Rightarrow (x + y) + c$$

$$(4) \quad x + (y + c) \Rightarrow (x + y) + c$$

$$(5) \quad (x + c_1) + (y + c_2) \Rightarrow (x + y) + (c_1 + c_2)$$

Experiment at the time:

- Support 51/139 ONNX operations.
 - Now: 75/209 not supported.
- Can compile MNIST and Resnet.
 - I can compile efficientnet_v2 (image classification)
- Can run on x86 machine, IBM power systems, system Z.
- Encountered big-endian and little-endian problem.
 - ONNX is little-endian, system Z is big-endian
 - MLIR does not support big-endian well.
 - Patch was able to resolve the problem.

Experiment at the time:

- Benchmark: (2.3-GHz POWER9 processors)
- MNIST: Graph Rewriting is applied.
- ResNet50: No optimization applied.
- Author believe that if polyhedral optimizations, SIMD optimization, and loop fusion is applied, the performance would be improved.

Table 1: Run inferencing with MNIST and ResNet50 on a POWER9 machine. Time in seconds.

Model	Compilation time	Inference time
MNIST	0.237	0.001
ResNet50	7.661	7.540

My Experiment:

- Model: efficientnet_v2
 - 87.9MB
 - Light weight Image Classification model
- Dataset: cifar100
- Compile time: 1m27.788s (single threaded)
- Inference time: 0m1.738s
- What can we get from this information?
 - I don't know.

My Experiment:

- How did I experiment?
 - Use ImageNet-1k pretrained weight to do transfer learning.
 - Convert the pytorch model to .onnx
 - Compile model with onnx-mlir to generate .so file
 - Rename the PyRuntime.xxx.so to PyRuntime.so
 - Link and export don't work.
 - Implement the driver.
- The actually used command:

```
torch.onnx.export(model, example_input, onnx_path, verbose=True)  
~/onnx-mlir/build/Debug/bin/onnx-mlir -O3 --EmitLib /home/sylvex/dl_cnn_hw1/efficientnet_32.onnx
```

My Experiment:

- How did I implement driver?
 - Python Runtime.
 - docs/mnist_example would help you to write c/python/java driver

```
import PyRuntime
session = PyRuntime.OMExecutionSession("./efficientnet_32.so")
outputs = session.run([image_array])
prediction = outputs[0]

for i in range(0, 100):
    score = prediction[0, i]
    if score > max_score:
        candidate = i
        max_score = score
    print(f"candidate:{coarse_label[i]} has score:{score}")

print(f"candidate:{coarse_label[candidate]} has max score:{max_score}")
```

My Experiment:

- Result:

Input:



Output:

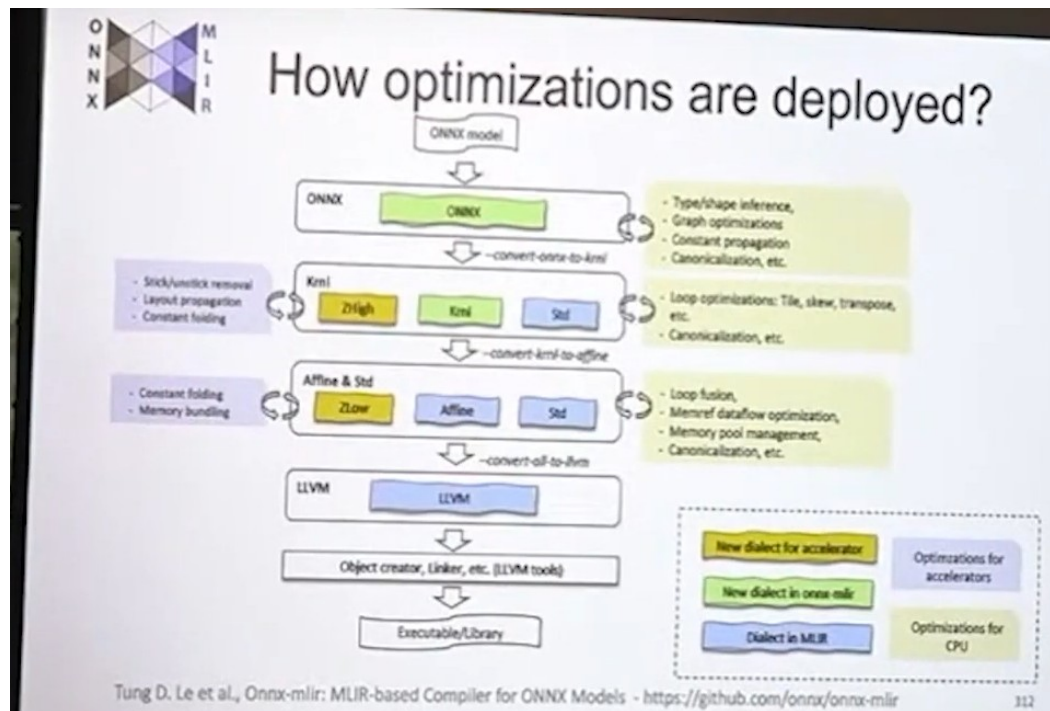
```
candidate:tractor has score:-4.397192001342773   candidate:tractor has score:-0.8939839005470276
candidate:train has score:-4.818047523498535     candidate:train has score:-2.0960440635681152
candidate:trout has score:-2.1939566135406494     candidate:trout has score:-2.1811110973358154
candidate:tulip has score:0.8204006552696228      candidate:tulip has score:-1.3569916486740112
candidate:turtle has score:-5.017601490020752     candidate:turtle has score:-1.5122387409210205
candidate:wardrobe has score:-1.3327709436416626  candidate:wardrobe has score:0.15867894887924194
candidate:whale has score:-3.0263237953186035     candidate:whale has score:-1.0275558233261108
candidate:willow_tree has score:2.586937427520752 candidate:willow_tree has score:-2.600288867950439
candidate:wolf has score:-3.4177582263946533     candidate:wolf has score:6.0673394203186035
candidate:woman has score:-2.024754762649536     candidate:woman has score:7.935925483703613
candidate:worm has score:1.5632660388946533      candidate:worm has score:-0.6879879236221313
candidate:orchid has max score:9.20013427734375  candidate:girl has max score:16.172103881835938
```

Paper Conclusion:

- The project transform onnx graph to native code with MLIR.
- Proposed Dialects: onnx, krnl
- Discussed Optimizations Implemented:
 - Operation decomposition, graph rewriting, constant propagation, etc
- Easily integrates new optimizations thanks to MLIR
- Future Optimizations:
 - Polyhedral optimization, Loop fusion, SIMD optimization.
 - Code generation for accelerators.

Research Opportunity:

- How to add accelerator?



Onnx-mlir: an MLIR-based Compiler for ONNX Models - The Latest Status

Personal thought:

- Prove-of-concept works, project have long way to go.
- Are there optimization chance lost when:
 - Conversion between dialect.
 - Dialect coverage.
 - Weird hardware like 1 operation complete convolution?
- MLIR provide template for building compiler faster, but the core is how you redirect the problem.
- Have to read more about the ecosystem if I research for MLIR.

Reference:

- I forgot.
- I've seemed too much stuff, forget to note them respectively.
- Sorry.
- Definitely not lazy.
- Trust me.

Thank you



喂，老師嗎？

我現在在宿舍樓頂

我真的看不完文獻了

這裡風好大
我好害怕

爸爸:期末考考得怎麼樣？我:



成績出來之前的我:

