# An Introduction to Torch-MLIR
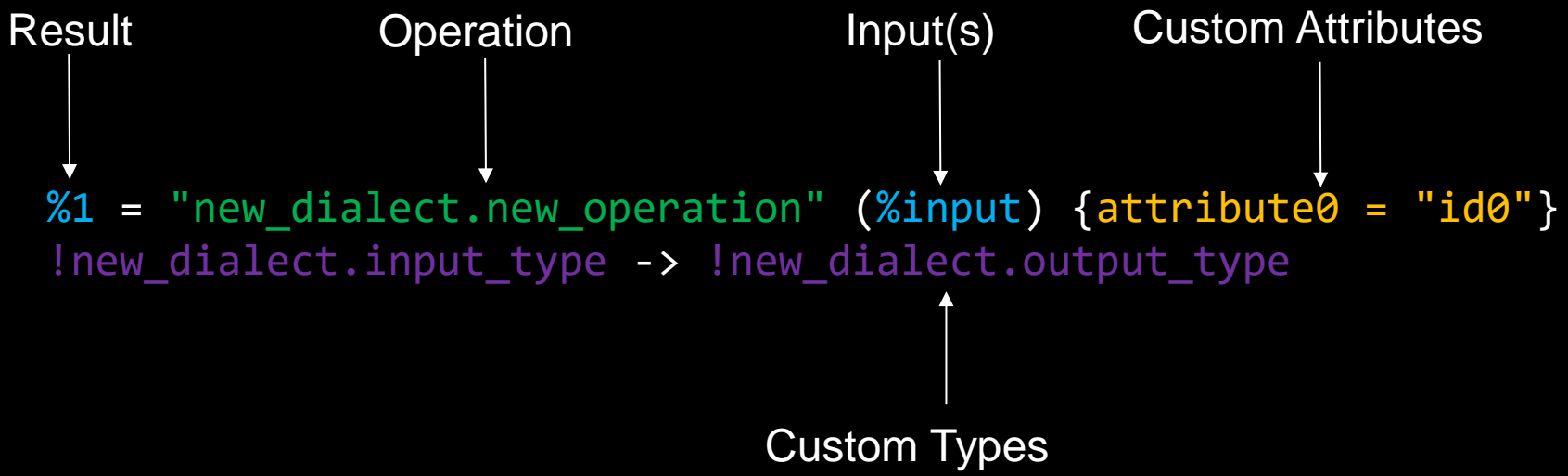
**Marius Brehler**

**AMD**
together we advance_

# Outline

- What is MLIR?

- What is Torch-MLIR?

- The `torch` Dialect

- Building the Project

- Importing, Converting and Running Models

- How to Get Involved

**AMD**
together we advance_

# What is MLIR?

- (Novel) Approach to build a reusable and extensible compiler infrastructure

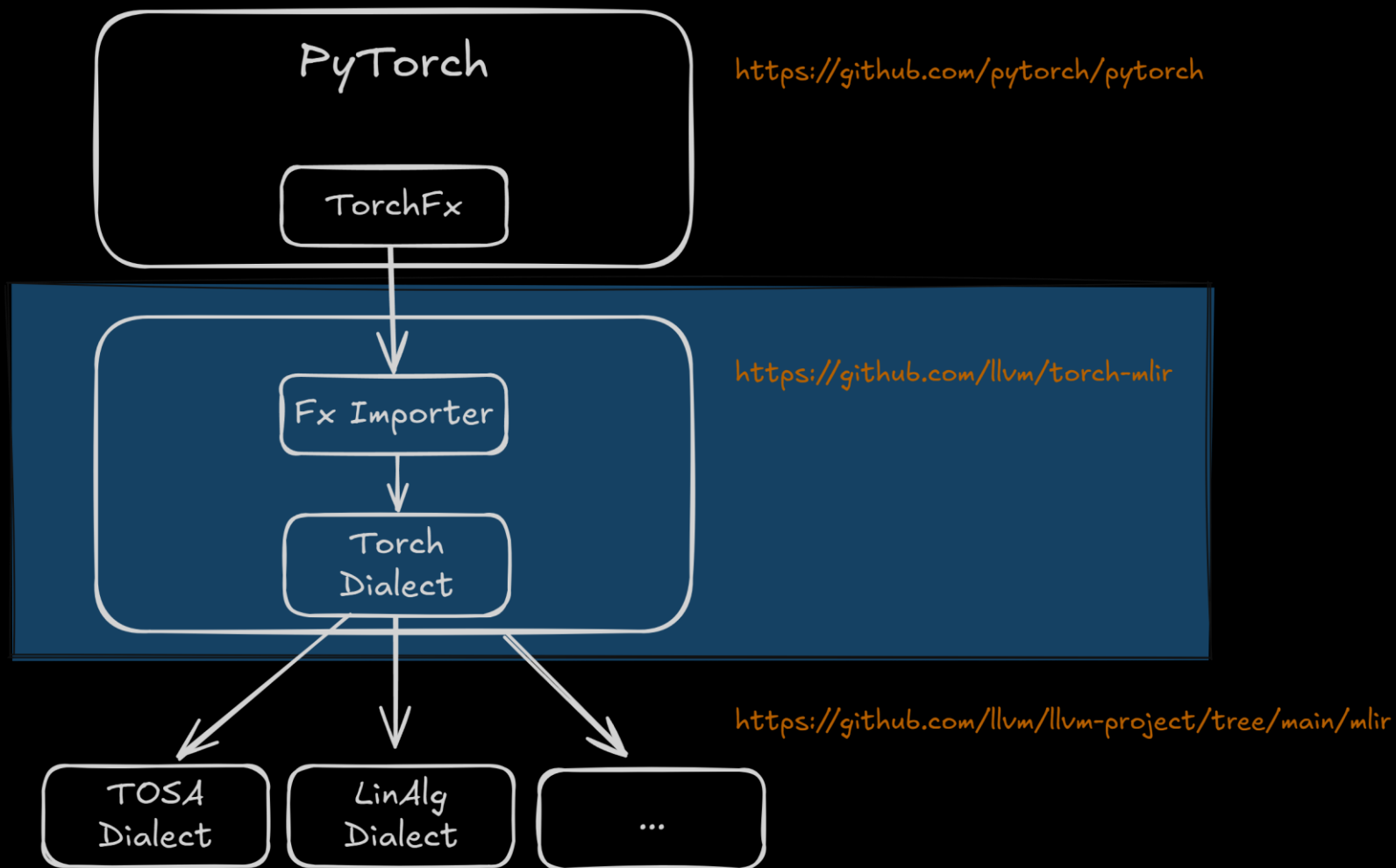- Provides an extensible intermediate representation (IR)



- MLIR is a project of the of the LLVM Foundation

# What is Torch-MLIR?

- Aims to provide first class compiler support from the PyTorch ecosystem to the MLIR ecosystem

- Participating in the LLVM Incubator process

- Dual-Licensed under the Apache License 2.0 with LLVM Exceptions and a BSD-style License

**AMD**
together we advance_

# What is Torch-MLIR?



PyTorch

TorchFx

https://github.com/pytorch/pytorch

Fx Importer

https://github.com/llvm/torch-mlir

Torch Dialect

TOSA Dialect

LinAlg Dialect

...

https://github.com/llvm/llvm-project/tree/main/mlir

**AMD**
together we advance_

# The `torch` Dialect

- The central MLIR abstraction is the `torch` dialect

- Mostly auto-generated based on the PyTorch JIT IR operator registry

- Some manually implemented ops:

  - Used for modeling PyTorch IValue object graphs e.g. `torch.nn_module, torch.class_type`)

  - `torch.global_slot` and related ops to model an incremental lowering of the IValue object graphs

  - Ops supported in the JIT interpreter directly and therefore without a corresponding op in the registry

  - `torch.operator` to represent ops from the registry which haven't been auto-generated

AMD△
together we advance_

# The `torch` Dialect - `TorchOps.td`

≡ TorchOps.td  ✕

include > torch-mlir > Dialect > Torch > IR > ≡ TorchOps.td

```
21
22    class Torch_Op<string mnemonic, list<Trait> traits = []>
23      │ │   : Op<Torch_Dialect, mnemonic, traits> {
24    }
25
26    include "torch-mlir/Dialect/Torch/IR/GeneratedTorchOps.td"
27
28    //===----------------------------------------------------------------------===//
29    // TorchScript `torch.nn.Module` object instantiation ops.
30    //===----------------------------------------------------------------------===//
31
32    def Torch_NnModuleOp : Torch_Op<"nn_module", [
33      │ │   DeclareOpInterfaceMethods<SymbolUserOpInterface>,
34      │ │   SingleBlockImplicitTerminator<"::mlir::torch::Torch::NnModuleTerminatorOp">]> {
35      │   let summary = "Constructs a torch.nn.Module";
```

AMD
together we advance_

# Setup the Project

```
# Clone the project and initialize the submodules

git clone https://github.com/llvm/torch-mlir
cd torch-mlir
git submodule update --init --progress --depth=1



# Create a virtual environment and install dependencies

python3 -m venv mlir_venv
source mlir_venv/bin/activate

python -m pip install --upgrade pip

python -m pip install -r requirements.txt
python -m pip install -r torchvision-requirements.txt
```

**AMD**
together we advance_

# Build the Project – Monolithic Build

```
cmake -GNinja -Bbuild-external \
  externals/llvm-project/llvm \
  -DCMAKE_BUILD_TYPE=Release \
  -DCMAKE_C_COMPILER=clang -DCMAKE_CXX_COMPILER=clang++ \
  -DCMAKE_LINKER_TYPE=LLD \
  -DCMAKE_C_COMPILER_LAUNCHER=ccache -DCMAKE_CXX_COMPILER_LAUNCHER=ccache \
  -DLLVM_ENABLE_ASSERTIONS=ON -DPython3_FIND_VIRTUALENV=ONLY \
  -DLLVM_ENABLE_PROJECTS=mlir \
  -DMLIR_ENABLE_BINDINGS_PYTHON=ON -DLLVM_TARGETS_TO_BUILD=host \
  -DLLVM_EXTERNAL_PROJECTS="torch-mlir" \
  -DLLVM_EXTERNAL_TORCH_MLIR_SOURCE_DIR="$PWD" \
  -DLIBTORCH_CACHE=ON -DLIBTORCH_SRC_BUILD=ON -DLIBTORCH_VARIANT=shared \
  -DTORCH_MLIR_ENABLE_PYTORCH_EXTENSIONS=ON \
  -DTORCH_MLIR_ENABLE_JIT_IR_IMPORTER=ON

cmake --build build-external/ --target check-torch-mlir
```

AMD

together we advance_

# Build the Project – Component Build (1)

```
# Build and install MLIR
export LLVM_INSTALL_DIR="`realpath install-llvm/`"

cmake -GNinja –Bbuild-llvm \
   externals/llvm-project/llvm \
   -DCMAKE_BUILD_TYPE=Release \
   -DCMAKE_C_COMPILER=clang -DCMAKE_CXX_COMPILER=clang++ \
   -DCMAKE_LINKER_TYPE=LLD \
   -DCMAKE_C_COMPILER_LAUNCHER=ccache -DCMAKE_CXX_COMPILER_LAUNCHER=ccache \
   -DLLVM_ENABLE_ASSERTIONS=ON -DPython3_FIND_VIRTUALENV=ONLY \
   -DLLVM_ENABLE_PROJECTS=mlir \
   -DMLIR_ENABLE_BINDINGS_PYTHON=ON -DLLVM_TARGETS_TO_BUILD=host \
   -DLLVM_INSTALL_UTILS=ON \
   -DCMAKE_INSTALL_PREFIX="${LLVM_INSTALL_DIR}"

cmake --build build-llvm/ --target install
```

AMD
together we advance_

# Build the Project – Component Build (2)

```
# Build Torch-MLIR
export LLVM_BUILD_DIR="`realpath build-llvm/`"

cmake -GNinja -Bbuild-component \
  -DCMAKE_BUILD_TYPE=Release \
  -DCMAKE_C_COMPILER=clang -DCMAKE_CXX_COMPILER=clang++ \
  -DCMAKE_LINKER_TYPE=LLD \
  -DCMAKE_C_COMPILER_LAUNCHER=ccache -DCMAKE_CXX_COMPILER_LAUNCHER=ccache \
  -DLLVM_ENABLE_ASSERTIONS=ON -DPython3_FIND_VIRTUALENV=ONLY \
  -DMLIR_DIR="${LLVM_INSTALL_DIR}/lib/cmake/mlir/" \
  -DLLVM_DIR="${LLVM_INSTALL_DIR}/lib/cmake/llvm/" \
  -DLLVM_EXTERNAL_LIT="${LLVM_BUILD_DIR}/bin/llvm-lit" \
  -DMLIR_ENABLE_BINDINGS_PYTHON=ON -DLLVM_TARGETS_TO_BUILD=host \
  -DLIBTORCH_CACHE=ON -DLIBTORCH_SRC_BUILD=ON -DLIBTORCH_VARIANT=shared \
  -DTORCH_MLIR_ENABLE_PYTORCH_EXTENSIONS=ON \
  -DTORCH_MLIR_ENABLE_JIT_IR_IMPORTER=ON

cmake --build build-external/ --target check-torch-mlir
```

AMD together we advance_

# Alternative: Install Python Wheels

```
# Install PyTorch
pip install --index-url https://download.pytorch.org/whl/nightly/cpu --pre \
  torch==2.7.0.dev20250120


# Install ONNX
pip install onnx


# Install Torch-MLIR
pip install --find-links https://github.com/llvm/torch-mlir-
release/releases/expanded_assets/dev-wheels torch-mlir
```

AMD
together we advance_

# A Simple Torch Model

```python
from pathlib import Path
import torch

class SimpleNN(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = torch.nn.Linear(16, 10, bias=False)
        self.linear.weight = torch.nn.Parameter(torch.ones(10, 16))
        self.relu = torch.nn.ReLU()
        self.train(False)

    def forward(self, input):
        return self.relu(self.linear(input))

input = torch.randn(2, 16)
model = SimpleNN()
```

AMD△
together we advance_

# A Simple Torch Model Imported to MLIR

```python
from pathlib import Path
import torch

class SimpleNN(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = torch.nn.Linear(16, 10, bias=False)
        self.linear.weight = torch.nn.Parameter(torch.ones(10, 16))
        self.relu = torch.nn.ReLU()
        self.train(False)

    def forward(self, input):
        return self.relu(self.linear(input))

input = torch.randn(2, 16)
model = SimpleNN()

from torch_mlir.fx import export_and_import
mlir_module = export_and_import(model, input, output_type="torch")
```

AMD⌐
together we advance_

# Importing an ONNX Model

```
# Download an ONNX model
wget
https://github.com/onnx/models/raw/main/validated/vision/classification/mnist/
model/mnist-12.onnx


# Import the model
torch-mlir-import-onnx mnist-12.onnx --opset-version 17 -o mnist.mlir
```

**AMD**
together we advance_

# mnist.mlir (Snippet)

…

```
%9 = torch.operator "onnx.Add"(%8, %3) : (!torch.vtensor<[1,8,28,28],f32>,
!torch.vtensor<[8,1,1],f32>) -> !torch.vtensor<[1,8,28,28],f32>

%10 = torch.operator "onnx.Relu"(%9) : (!torch.vtensor<[1,8,28,28],f32>) ->
!torch.vtensor<[1,8,28,28],f32>

%11 = torch.operator "onnx.MaxPool"(%10) {torch.onnx.auto_pad = "NOTSET",
torch.onnx.kernel_shape = [2 : si64, 2 : si64], torch.onnx.pads = [0 : si64, 0
: si64, 0 : si64, 0 : si64], torch.onnx.strides = [2 : si64, 2 : si64]} :
(!torch.vtensor<[1,8,28,28],f32>) -> !torch.vtensor<[1,8,14,14],f32>
```

…

AMD△
together we advance_

# TorchOnnxToTorch - Abs

lib > Conversion > TorchOnnxToTorch > C++ DefaultDomainAtoF.cpp > ...

```cpp
122    // Simple rewrites for the default domain.
134    void mlir::torch::onnx_c::populateDefaultDomainAtoF(
135        OnnxCustomOpConversionPattern &patterns) {
136    patterns.onOp("Abs", 1,
137                        [](OpBinder binder, ConversionPatternRewriter &rewriter) {
138                            Torch::ValueTensorType resultType;
139                            Value operand;
140                            if (binder.tensorOperand(operand) ||
141                                binder.tensorResultType(resultType))
142                              return failure();
143                            rewriter.replaceOpWithNewOp<Torch::AtenAbsOp>(
144                                binder.op, resultType, operand);
145                            return success();
146                        });
```

AMD
together we advance_

# TorchOnnxToTorch - AveragePool

lib > Conversion > TorchOnnxToTorch > G+ DefaultDomainAtoF.cpp > ⬡ populateDefaultDomainAtoF(OnnxCustomOpConver:

```
135        OnnxCustomOpConversionPattern &patterns) {
456        patterns.onOp(
457            "AveragePool", 11,
458            [](OpBinder binder, ConversionPatternRewriter &rewriter) {
459              std::string autoPad;
460              SmallVector<int64_t> dilations;
461              if (binder.customOpNameStringAttr(autoPad, "auto_pad", "NOTSET"))
462                return failure();
463              if (autoPad != "NOTSET") {
464                // TODO: Add support for `auto_pad` != "NOTSET"
465                return rewriter.notifyMatchFailure(
466                    binder.op, "unsupported conversion: auto_pad != NOTSET");
467              }
468
469            Torch::ValueTensorType resultType;
470            Value operand;
```

18

**AMD**
together we advance_

# DecomposeComplexOps

lib > Dialect > Torch > Transforms > ⊙ DecomposeComplexOps.cpp > {} `anonymous-namespace` > ⫶ DecomposeAt

```cpp
2967    namespace {
2968    class DecomposeAtenHardshrinkOp : public OpRewritePattern<AtenHardshrinkOp> {
2972                              PatternRewriter &rewriter) const override {

2999      Value posMask =
3000          rewriter.create<AtenGtScalarOp>(loc, boolResType, self, poslambd);
3001      Value negMask =
3002          rewriter.create<AtenLtScalarOp>(loc, boolResType, self, neglambd);

3004      Value result = rewriter.create<AtenWhereScalarOtherOp>(loc, resTy, posMask,
3005                                                             self, zero);
3006      result =
3007          rewriter.create<AtenWhereSelfOp>(loc, resTy, negMask, self, result);

3009    rewriter.replaceOp(op, result);
```

AMD △
together we advance_

# Pipelines

- `torch-backend-to-linalg-on-tensors-backend-pipeline`
- `torch-backend-to-stablehlo-backend-pipeline`
- `torch-backend-to-tosa-backend-pipeline`
- `torch-function-to-torch-backend-pipeline`
- `torch-onnx-to-torch-backend-pipeline`
- `torch-shape-refinement-pipeline`
- `torch-simplification-pipeline`
- `torchdynamo-export-to-torch-backend-pipeline`
- `torchscript-module-to-torch-backend-pipeline`

AMD

together we advance_

# RefBackend

```python
# See pt1/examples/fximporter_resnet18.py

from torch_mlir_e2e_test.linalg_on_tensors_backends import refbackend
from torch_mlir_e2e_test.configs.utils import (
    recursively_convert_to_numpy,
)


resnet18 = models.resnet18(pretrained=True).eval()
module = fx.export_and_import(
    resnet18,
    torch.ones(1, 3, 224, 224),
    output_type="linalg-on-tensors",
    func_name=resnet18.__class__.__name__,
)
backend = refbackend.RefBackendLinalgOnTensorsBackend()
compiled = backend.compile(module)
fx_module = backend.load(compiled)
```

**AMD◢**
together we advance_

# How To Get Involved

- Try to import and lower a model by using `torch-mlir-opt`

- Fill issues @ https://github.com/llvm/torch-mlir

- Join discussions @ #torch-mlir (LLVM Discord)

- Implement support for attributes of ONNX operators

- …

AMD
together we advance_

# Disclaimer

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated.  AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS." AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

**AMD**
together we advance_