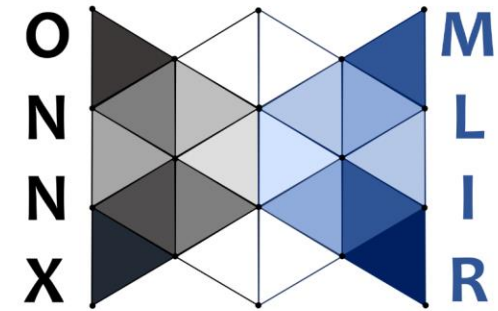# Onnx-mlir:
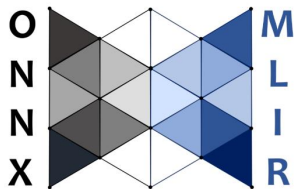# MLIR-based Compiler
# for ONNX Models
# The Latest Status

Tung D. Le, Tong Chen, Ettore Tiotto, Haruki
Imai, Yasushi Negishi, Kevin O'Brien, Kiyokuni
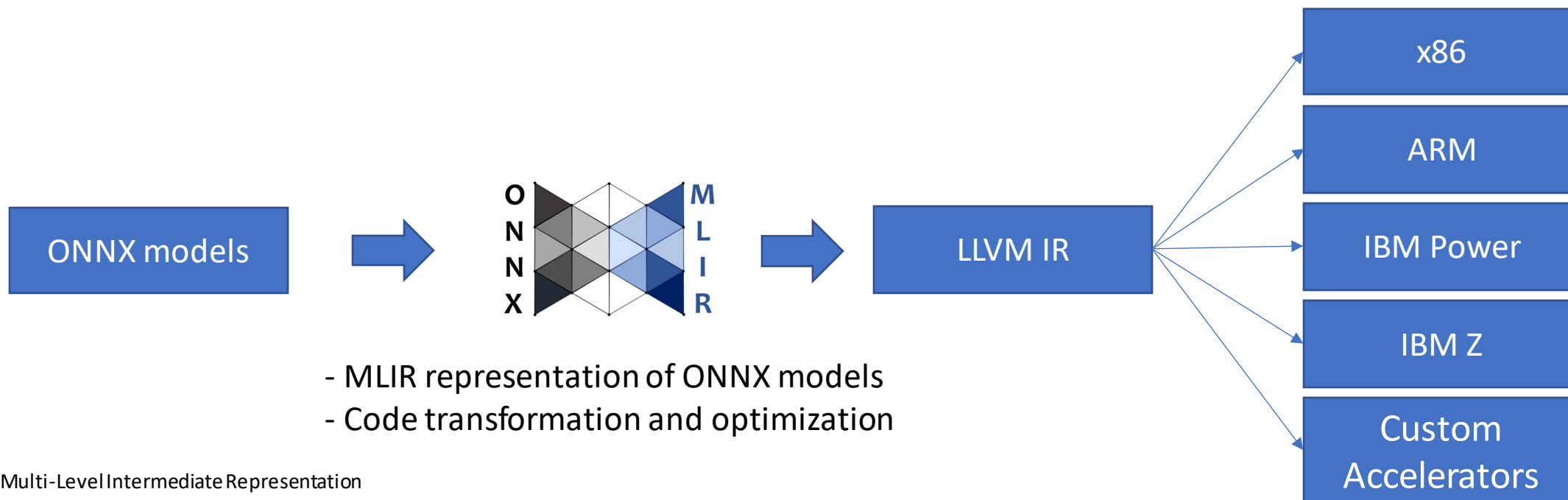Kawachiya, Alexandre E Eichenberger

IBM Research

https://github.com/onnx/onnx-mlir
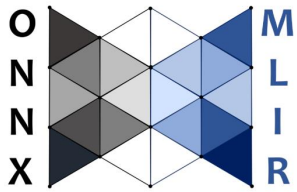
Presenting the work of many people!
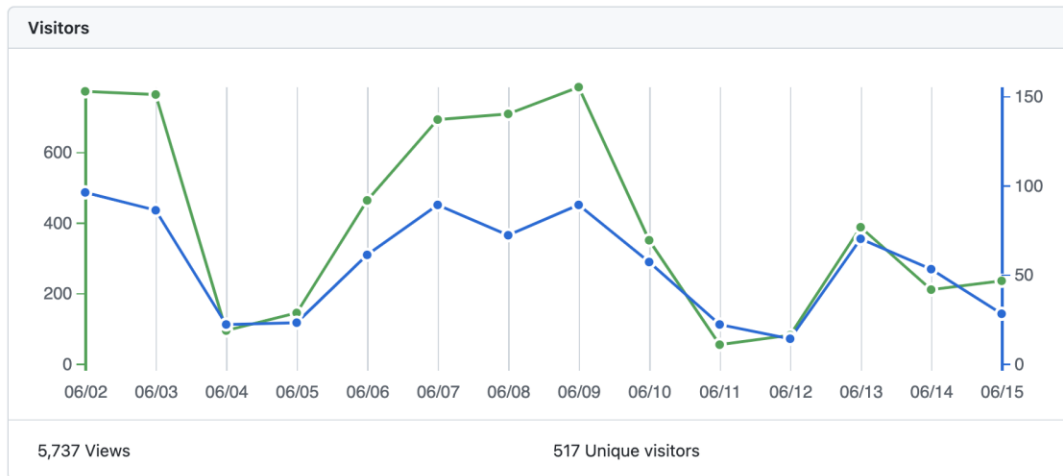
# What is onnx-mlir?

- Compile an ONNX model to an optimized binary using
  - MLIR* to perform high-level optimization transformations
  - LLVM to perform low-level optimizations and code generation



- MLIR representation of ONNX models
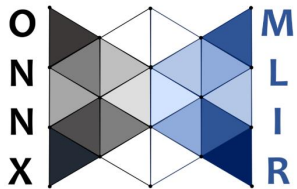- Code transformation and optimization

*MLIR: Multi-Level Intermediate Representation

# Young but active project



Commits per week over the last year



- Contributions and support from IBM, Microsoft, Arm, Facebook, and others

# Design goals

- A reference ONNX dialect in MLIR
- Easy to write optimizations for CPU and custom accelerators
  - From high-level (e.g., graph level) to low-level (e.g., instruction level)
- Easy to deploy
  - Stand-alone driver and runtime support in Python/C/C++/Java
  - Integration into other MLIR-based compilers
- Continuously tested
  - Unit tests and ONNX model zoo
  - x86, Power, z/Architecture
  - Windows, Linux, z/OS, macOS
  - Python/C++/Java

# Onnx-mlir in practice

- Tested with many models in the ONNX model zoo
- Deployed in IBM Watson Machine Learning for z/OS (WMLz)
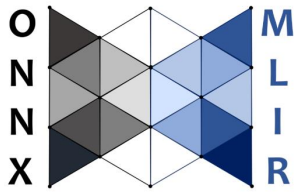- One of the inference engines in BentoML framework

| System | Build Status | Model Zoo Status |
|---|---|---|
| s390x-Linux | Jenkins CI  passing | Models  Total:159 Skipped:33 Passed:110 Failed:16 |
| ppc64le-Linux | Jenkins CI  passing | Models  Total:159 Skipped:33 Passed:110 Failed:16 |
| amd64-Linux | Jenkins CI  passing | Models  Total:159 Skipped:33 Passed:110 Failed:16 |
| amd64-Windows | Azure Pipelines  succeeded | |
| amd64-macOS | GitHub Action MacOS amd64  passing | |
| | openssf best practices  passing | |

Online scoring services in IBM WMLz

BENTOML

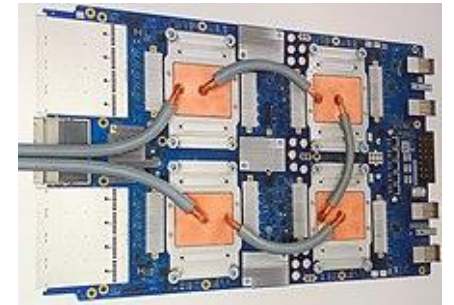The Unified Model Serving Framework
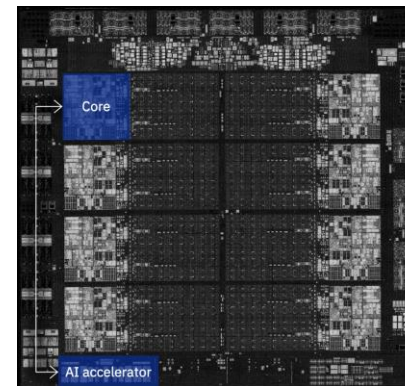
# Recent developments in onnx-mlir

- A framework for supporting custom accelerators
    - Easy to offload ONNX operators to accelerators
    - Custom optimizations for custom accelerators



NVIDIA GPU

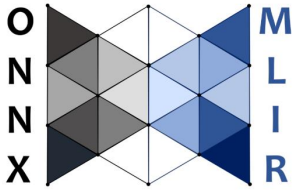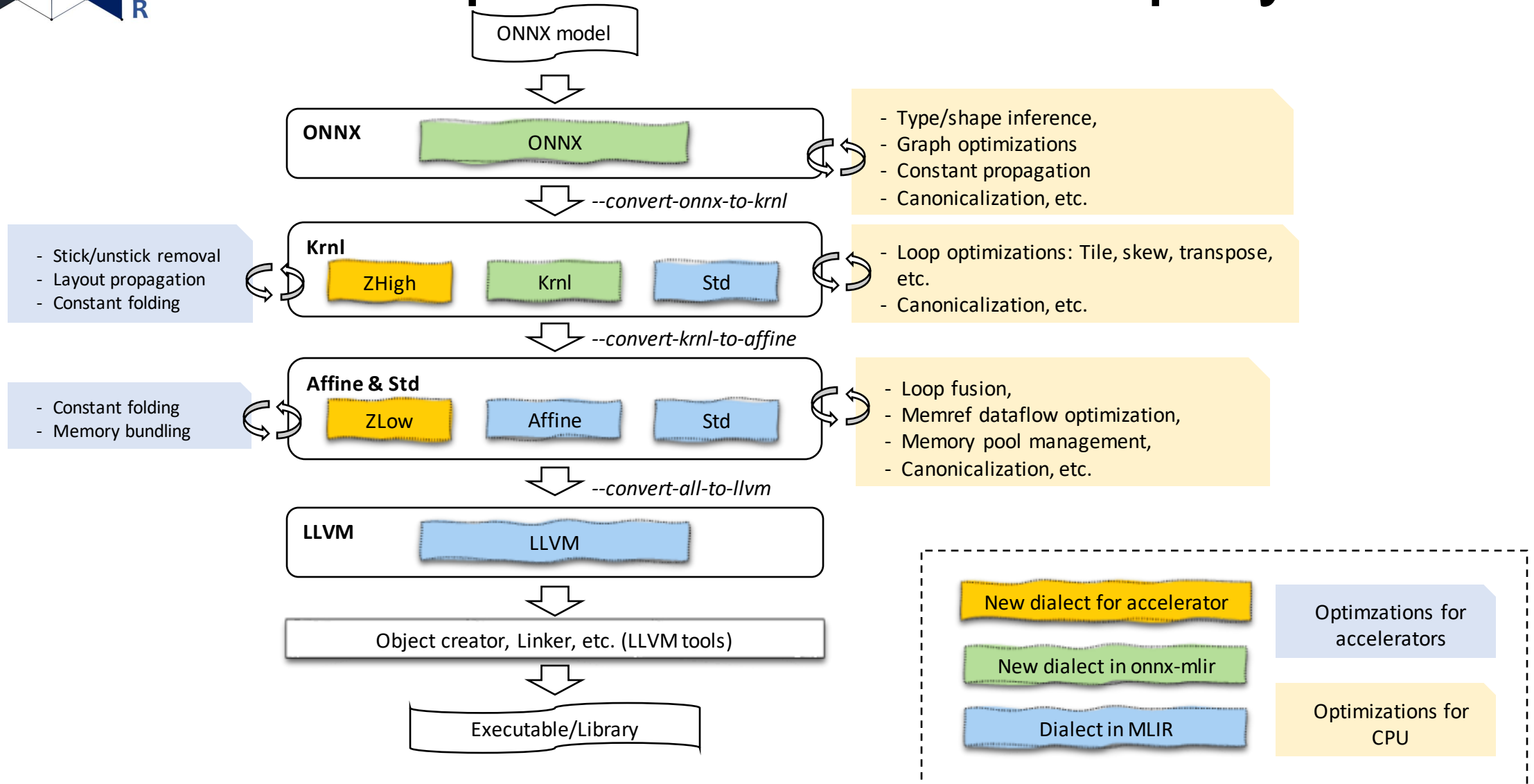- We have demonstrated the framework for IBM on-chip low-latency AI accelerator introduced in IBM z16.
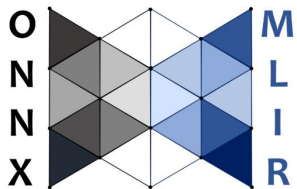


Google TPU



IBM on-chip AI accelerator
in Telum processor

# How optimizations are deployed?



ONNX model

**ONNX**
ONNX
- Type/shape inference,
- Graph optimizations
- Constant propagation
- Canonicalization, etc.

--convert-onnx-to-krnl

**Krnl**
- Stick/unstick removal
- Layout propagation
- Constant folding

ZHigh | Krnl | Std

- Loop optimizations: Tile, skew, transpose, etc.
- Canonicalization, etc.

--convert-krnl-to-affine

**Affine & Std**
- Constant folding
- Memory bundling

ZLow | Affine | Std

- Loop fusion,
- Memref dataflow optimization,
- Memory pool management,
- Canonicalization, etc.

--convert-all-to-llvm

**LLVM**
LLVM

Object creator, Linker, etc. (LLVM tools)

Executable/Library

New dialect for accelerator — Optimizations for accelerators
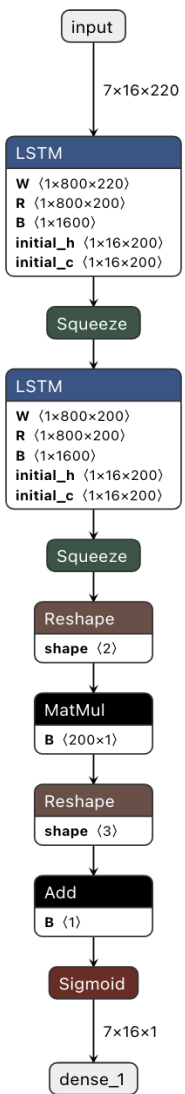
New dialect in onnx-mlir

Dialect in MLIR — Optimizations for CPU

# Example: Credit Card Fraud Detection (CCFD)

CCFD model: https://github.com/IBM/ai-on-z-fraud-detection

input

7×16×220

**LSTM**
W ⟨1×800×220⟩
R ⟨1×800×200⟩
B ⟨1×1600⟩
initial_h ⟨1×16×200⟩
initial_c ⟨1×16×200⟩

Squeeze

**LSTM**
W ⟨1×800×200⟩
R ⟨1×800×200⟩
B ⟨1×1600⟩
initial_h ⟨1×16×200⟩
initial_c ⟨1×16×200⟩

Squeeze

Reshape
shape ⟨2⟩

MatMul
B ⟨200×1⟩

Reshape
shape ⟨3⟩

Add
B ⟨1⟩

Sigmoid

7×16×1

dense_1

Inference speedup for CCFD
with the input size of 7x16x220
(Higher is faster)

11.3 x

1 x
(baseline)

■ zCPU    ■ IBM on-chip AI Accelerator

Real-time detection of fraudulent transactions

Result here is unoffical and does not represent any IBM product

# Summary

- Onnx-mlir is an open-source compiler for ONNX models
  - Easy to do optimizations and support new accelerators

- Call for contribution
  - Cool compiler technologies, e.g., AI for compilers, etc.
  - More architectures of interest

- Some areas of interest in near future
  - Optimize operators: Conv, Pooling, Reduction, etc.
  - Support ONNX machine learning operators
  - Support other accelerators, e.g., GPGPU

**We truly thank all contributors to onnx-mlir!**