

# MLIR-PyTACO: An End-to-End Use Case For the Sparse Tensor Compiler

Bixia Zheng ([bixia@google.com](mailto:bixia@google.com))

Aart Bik ([ajcbik@google.com](mailto:ajcbik@google.com))

For: MLIR ODM, 2/10/2022



# Outline

- What is MLIR-PyTACO
- Sparse Tensor Algebra Support in MLIR
- Implementing PyTACO with MLIR
- Source Code Organization
- Ongoing/Future Work



# MLIR-PyTACO

- TACO: The Tensor Algebra Compiler
- PyTACO: Python-based DSL for using TACO
- MLIR-PyTACO
  - Use MLIR to implement the DSL
    - Index notation based tensor algebra syntax
    - Tensor I/O



# Index Notation Based Tensor Algebra

- Tensor:  $A = \text{tensor}([2,3], [\text{compressed}, \text{dense}])$
- Index\_var:  $i, j = \text{get\_index\_var}(2)$ 
  - Represents dimensions in the iteration space
- Tensor access expression:  $A[i, j]$
- Tensor assignment:  $C[i, j] = A[i, j] + B[i, j]$ 
  - The same index\_var corresponds to dimensions of the same size
  - Element-wise operations:  $+, -, *, /$
  - Implicit reduction operations:  $+$



# Tensor I/O

- Coordinate format
  - (coordinates, non-zero-value)
- Support variants of coordinate formats
  - Matrix Market format (.mtx)
  - FROSTT format (.tns)
- Tensor input:  **$A = \text{read}(\text{"input.mtx"}, [\text{compressed}, \text{dense}])$**
- Tensor output:  **$\text{write}(\text{"output.tns"}, C)$**
- Create dense tensor from numpy array: `from_array`



# SpMV via MLIR-PyTACO

```
import mlir_pytaco_api as pt
import numpy as np

csr = pt.format([pt.dense, pt.compressed])
dv = pt.format([pt.dense])

A = pt.read("pwtk.mtx", csr)
X = pt.from_array(np.random.uniform(size=A.shape[1]))
Z = pt.from_array(np.random.uniform(size=A.shape[0]))
Y = pt.tensor([A.shape[0]], dv)

i, j = pt.get_index_vars(2)

Y[i] = A[i, j] * X[j] + Z[i]
pt.write("Y.tns", Y)
```



# Outline

- What Is MLIR-PyTACO
- Sparse Tensor Algebra Support in MLIR



# Tensor Types And Sparsity

```
tensor<3x4xf64>
```

```
#CSR = #sparse_tensor.encoding<{ dimLevelType = [ "dense", "compressed" ] }>  
tensor<3x4xf64, #CSR>
```

```
#DCSR = #sparse_tensor.encoding<{ dimLevelType = [ "compressed", "compressed" ] }>  
tensor<3x4xf64, #DCSR>
```





# Tensor Algebra via linalg.generic

```
C[i] = A[i, j] * B[j, i]
```

```
#trait = {  
  indexing_maps = [  
    affine_map<(i,j) -> (i,j)>, // A  
    affine_map<(i,j) -> (j,i)>, // B  
    affine_map<(i,j) -> (i)> // C  
  ],  
  iterator_types = ["parallel", "reduction"], // Can we derive this from the result tensor affine_map?  
}  
%0 = linalg.generic #trait  
  ins(%argA, %argB : !typeA, !typeB)  
  outs(%argC : !typeC) {  
    ^bb0(%A: f64, %B: f64, %C: f64):  
      %1 = arith.mulf %A, %B : f64  
      %2 = arith.addf %C, %1 : f64  
      linalg.yield %2 : f64  
  } -> !typeC
```



# Linalg Tensor Algebra Semantics

- Implicit broadcasting: broadcast along “missing dimensions”
  - Expression:  $C[i, j] = A[i, j] + B[i]$
  - Semantics:  $C[i, j] = A[i, j] + \text{broadcast}(j, B[i])$
- Reduction: op-reduce over “reduction dimension”
  - Expression:  $C[i] = A[i, j] + B[i, j]$
  - Semantics:  $C[i] = \text{op}(j, A[i, j] + B[i, j])$
- Reduction is performed on the outer expression
  - Expression:  $D[i] = A[i, j] + B[i, j] + C[i]$
  - Semantics:  $D[i] = \text{sum}(j, A[i, j] + B[i, j], \text{broadcast}(j, C[i]))$
  - PyTACO Semantics:  $D[i] = \text{sum}(j, A[i, j] + B[i, j]) + C[i]$

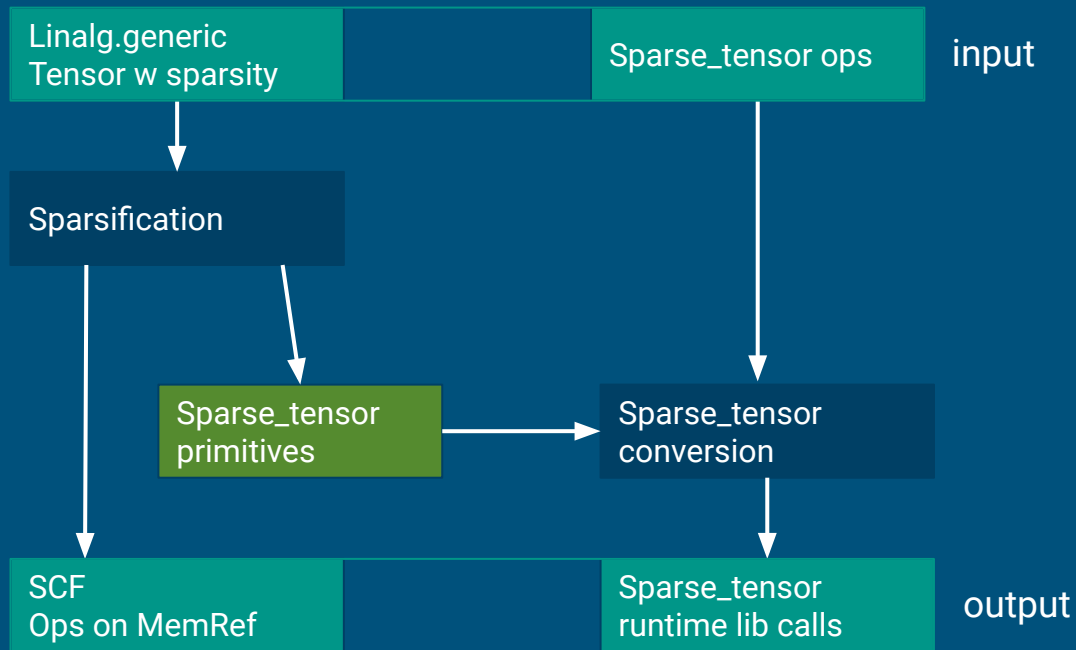
Same as  
PyTACO

Similar to  
PyTACO

Different from  
PyTACO



# Sparse Tensor Compiler



# Outline

- What Is MLIR-PyTACO
- Sparse Tensor Algebra Support in MLIR
- Implementing PyTACO With MLIR



# MLIR Python Bindings

- Thanks for the easy to use Python bindings for
  - Constructing instructions, functions
    - `builtin.FuncOp.from_py_func`
  - Constructing compilation pipelines
  - Building a JIT engine for a compiled module
  - JIT engine execution
- `Linalg.generic` through `linalg.OpDsl`

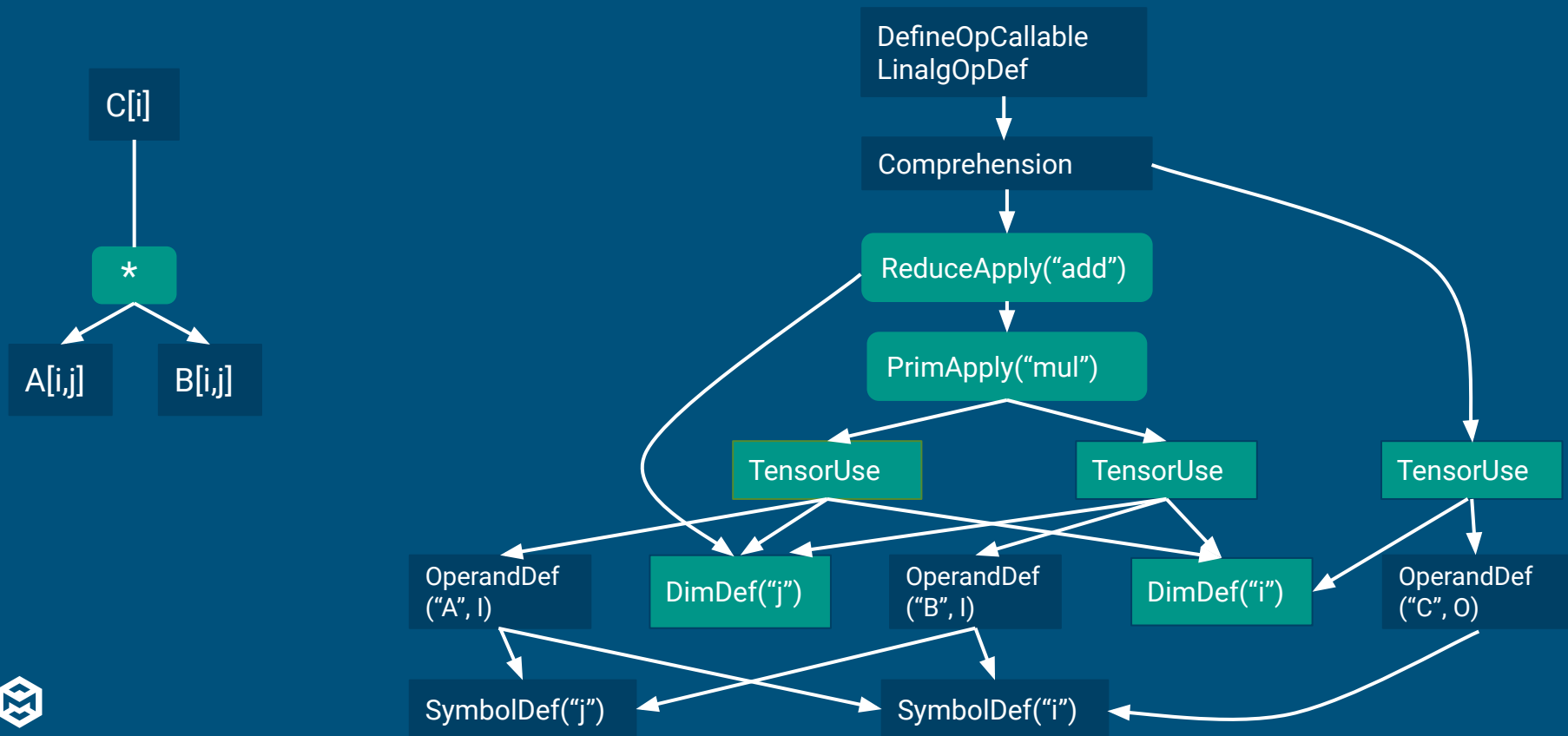


# Linalg.OpDsl Python Bindings

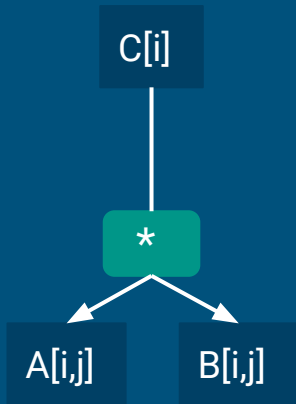
- Construct the AST representation for tensor expressions
  - LinalgOpDef
    - Comprehension: a tensor assignment
  - OperandDef: Input and Output of the function
    - SymbolDef: Dimension in the iteration space
  - TensorUse: represent operands
    - DimDef: Dimension in the operand
  - Operations
    - ReduceApply
    - PrimApply
- Convert the AST to linalg.generic
  - DefineOpCallable



# Tensor Expression to Linalg.opdsl AST



# Linalg.opdsl AST to linalg.generic



```
#trait = {  
  indexing_maps = [  
    affine_map<(i,j) -> (i,j)>, // A  
    affine_map<(i,j) -> (i,j)>, // B  
    affine_map<(i,j) -> (i)> // C  
  ],  
  iterator_types = ["parallel", "reduction"],  
}  
%0 = linalg.generic #trait  
  ins(%argA, %argB : !typeA, !typeB)  
  outs(%argC : !typeC) {  
    ^bb0(%A: f64, %B: f64, %C: f64):  
      %1 = arith.mulf %A, %B : f64  
      %2 = arith.addf %C, %1 : f64  
      linalg.yield %2 : f64  
  } -> !typeC
```





# Reduction on “smallest subexpressions”

$$D[i] = A[i, j] * B[i, j] + C[i]$$

Implemented with one `linalg.generic` gives incorrect result:

$$D[i] = \text{sum}(j, A[i, j] * B[i, j] + \text{broadcast}(j, C[i]))$$

Need to implement with two `linalg.generic`:

$$// D[i] = \text{sum}(j, A[i, j] * B[i, j]) + C[i]$$

$$T[i] = A[i, j] * B[i, j]$$

$$D[i] = T[i] + C[i]$$



# Estimate the Sparsity of Temporary Tensors

- $\text{op}(1, 0) = 0$ , computing disjunction of zero values
  - The result dimension is sparse if either source dimensions is sparse
- $\text{op}(1, 0) \neq 0$ , computing conjunction of zero values
  - The result dimension is sparse if both source dimensions are sparse



# Supporting Python Classes

- IndexVar
- ModeFormat
- Tensor
  - `__getitem__`: construct a tensor Access
  - `__setitem__`: construct an IndexExpr assigned to the tensor
  - `evaluate`: evaluate the tensor assignment
- Access
  - `__add__`, `__sub__`, etc: Construct IndexExpr
- IndexExpr



# Evaluating Tensor Assignment

1. Validate the assignment and collect info to support codegen
2. Identify subexpressions corresponding to a `linalg.generic`
3. Create an MLIR module

For each identified subexpression, in reversed topological order  
Emit `linalg.generic` via `DefinedOpCallable`

4. Compile the MLIR module
5. Build an execution engine for the MLIR module
6. Get ctype pointers for the input and output tensor storage
7. Invoke the execution engine and retrieve the result



# Source Code Organization

<https://github.com/llvm/llvm-project/tree/main/mlir/test/Integration/Dialect/SparseTensor/taco>

Aliases to expose PyTACO API  
mlir\_pytaco\_api.py

Tensor  
mlir\_pytaco.py

Tensor I/O  
mlir\_pytaco\_io.py

PyTACO API test  
Test\_simple\_tensor\_algebra.py  
test\_SpMM.py  
test\_SpMV.py  
test\_MTTKRP.py

Unit Test  
unit\_test\_tensor\_utils.py  
unit\_test\_tensor\_io.py



# Ongoing/Future Work

- Support transpose
- Performance tuning
- Support more PyTACO functionality

Paper “**Compiler Support for Sparse Tensor Computations in MLIR**”

Aart J.C. Bik, Penporn Koanantakool, Tatiana Shpeisman, Nicolas Vasilache, Bixia Zheng, and Fredrik Kjolstad

Available as [pre-print on ArXiv](#)



For questions on MLIR-PyTACO, please contact

Aart Bik [ajcbik@google.com](mailto:ajcbik@google.com)

Bixia Zheng [bixia@google.com](mailto:bixia@google.com)

