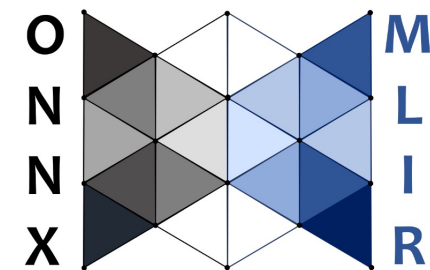


Dynamic Dimension Analysis in onnx-mlir Compiler

Tung D. Le (Speaker), Alexandre E Eichenberger, and Tong Chen

IBM Research



<https://github.com/onnx/onnx-mlir>

ONNX Community Meetup Day 2023



Dynamic Dimensions in Shape Inference

Dynamic dimension is an unknown value
at compile time

```
"onnx.Add"(%A, %B) : (tensor<3x?x?xf32>, tensor<3x?x?xf32>) -> tensor<3x?x?xf32>
```

- Information about shapes is key to high performance
 - Ruling out broadcasting,
 - Enabling Fusion / optimized SIMD code patterns
 - Determining if an operation is suitable for hardware accelerator,...
- ONNX-provided shape inference is not enough
 - Compiler has shape inference anyway to generate code that computes dynamic shapes
 - Compiler changes ONNX operation patterns, has multiple dialects in addition to ONNX



Dynamic Dimension Analysis in onnx-mlir

- Purpose: to explore relations among dynamic dimensions
 - Current focus on discovering dimension equivalences
 - For shapes with arbitrary mixtures of static and dynamic dimensions
- Two phases

Shape-related Operator
Canonicalization



Analysis

- Represent a shape by scalar constants and dimensions.
- Propagate scalars through shape-manipulating operators: Shape, Slice, Squeeze, etc.
- Explore dynamic dimension relationship using the existing shape inference infrastructure in onnx-mlir
- APIs to query information: `sameDim()`, `sameShape()`



Shape-related Operator Canonicalization (1/3)

- onnx.Shape is often the starting point of shape calculation

Shape is a tensor of integers, which is unknown at compile time

```
%0 = "onnx.Shape"(%arg0)      : (tensor<128x?x?xf32>) -> tensor<3xi64>  
%1 = "onnx.Reshape"(%arg1, %0) : (tensor<?x256xf32>, tensor<3xi64>)  
                                -> tensor<?x?x?xf32>
```

The consuming operator has no dimension information from the tensor shape

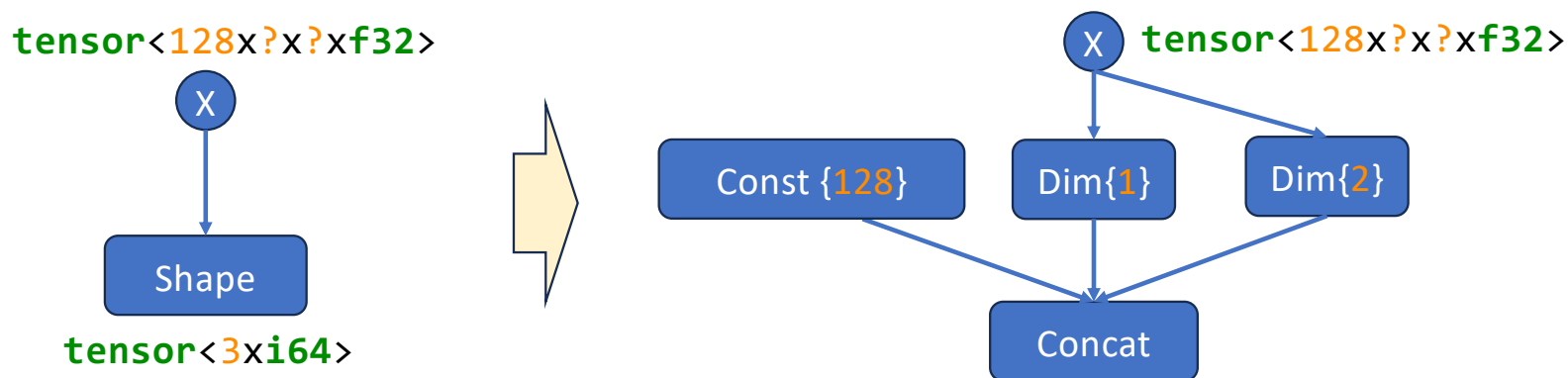
128 was not propagated to this output

- For complete shape inference, we must start keeping track of shape values that were saved in tensors.



Shape-related Operator Canonicalization (2/3)

- Expose individual dimensions to compiler



- A compile time shape value is represented by an onnx.Const
- A runtime shape value is represented by an onnx.Dim *

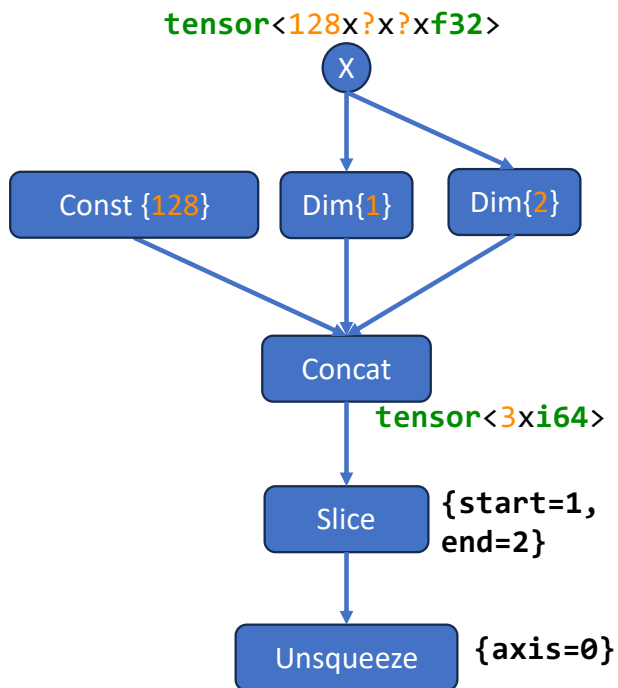
(*) onnx.Dim is an additional ONNX operator in onnx-mlir

Tung, Alexandre and Tong, Dynamic Dimension Analysis in onnx-mlir Compiler

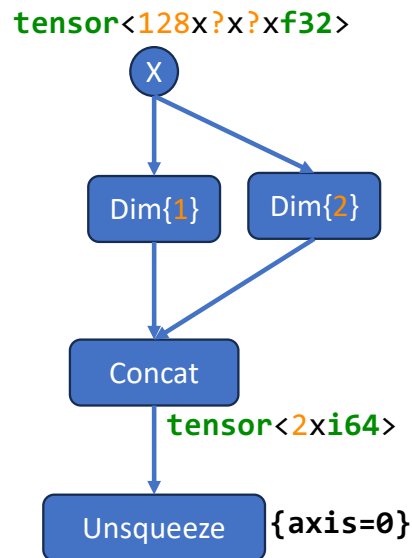


Shape-related Operator Canonicalization (3/3)

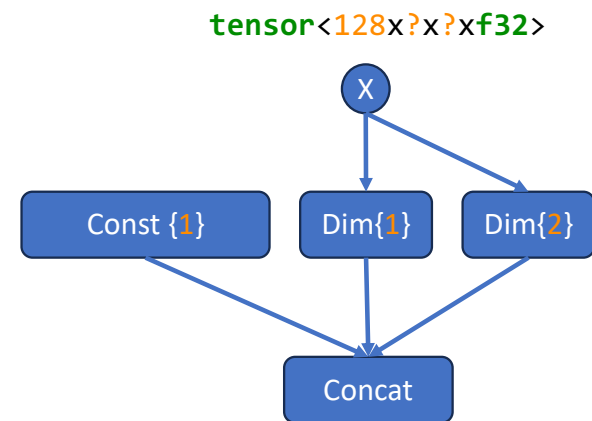
- Propagate scalar dimensions through shape-manipulating operators



Propagate through SliceOp



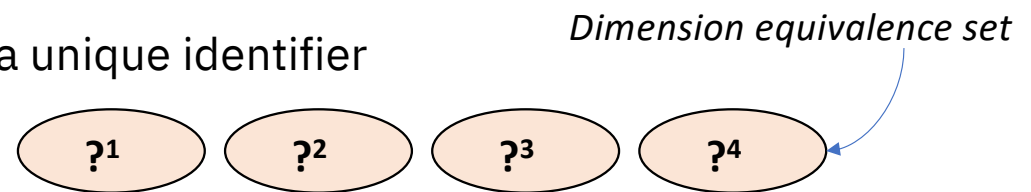
Propagate through UnsqueezeOp



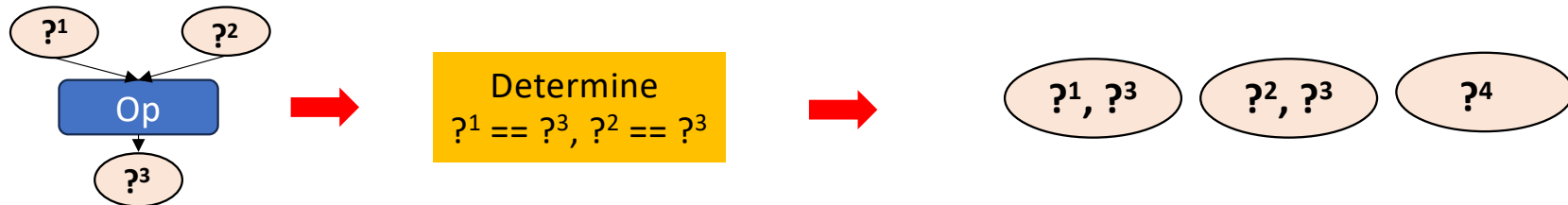
Analysis: Value Numbering for Shapes

- Analysis algorithm:

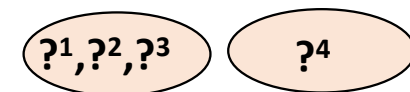
1. Associate each dynamic dimension `?` with a unique identifier



2. Analyze an operation for dimension equivalence



3. Apply transitive closure of dimension equivalence

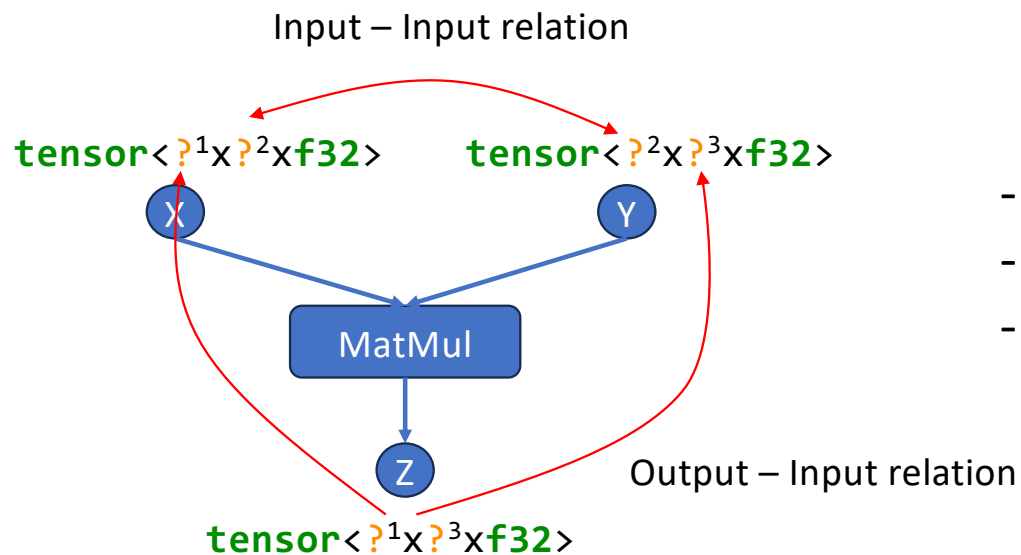


4. Repeat Steps 2. & 3. for operations until steady state is reached



Expand a dimension set (1/2)

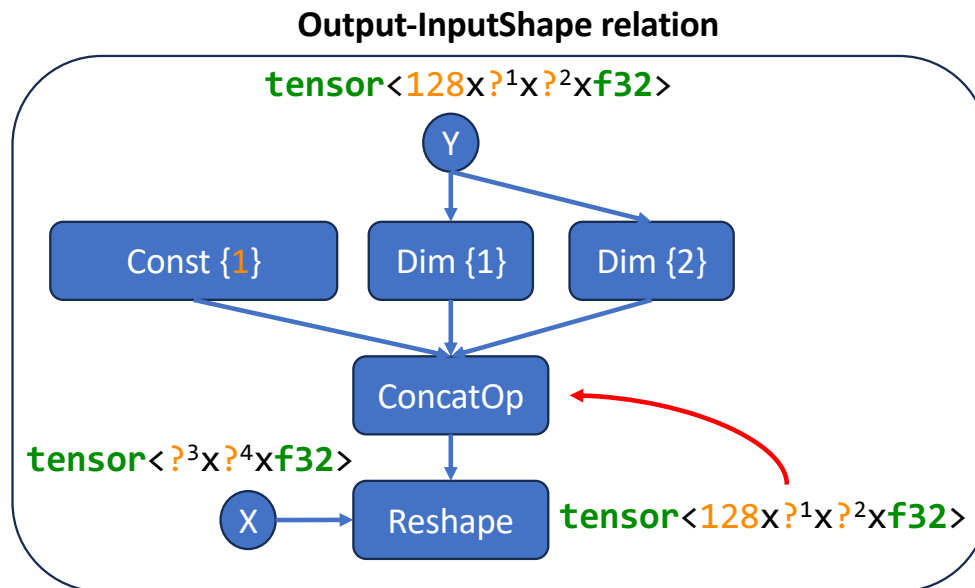
- Given an operator and a dynamic dimension in the output tensor
- Find ALL equivalent dynamic dimensions in the input tensors.



- $\text{Dim}(Z, 0)$ is equivalent to $\text{Dim}(X, 0)$
- $\text{Dim}(Z, 1)$ is equivalent to $\text{Dim}(Y, 1)$
- $\text{Dim}(X, 1)$ is equivalent to $\text{Dim}(Y, 0)$



Expand a dimension set (1/2)

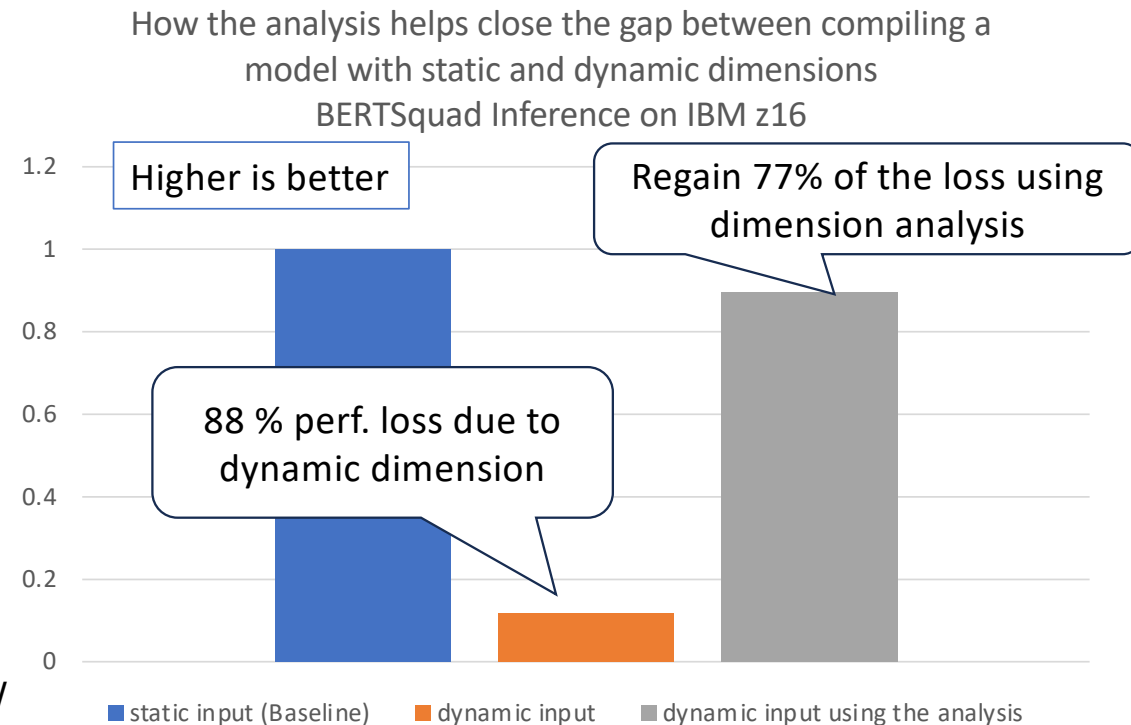


- We have implemented this rule for onnx Reshape, ConstantOfShape, Expand, MaxUnpool, CenterCropPad, Tile.



Apply to the BERTSquad model

- 854 dynamic dimensions are classified into 26 sets.
- Sizes of the 26 sets
 - 1 set of 817 dimensions
 - 12 sets with 2 dimensions
 - 13 sets with 1 dimension
- Future work:
 - Analysis with the presence of +, -, *, /



Thank you for your listening!

