

DialEgg

Dialect-Agnostic MLIR Optimizer using Equality Saturation with Egglog

CGO 2025

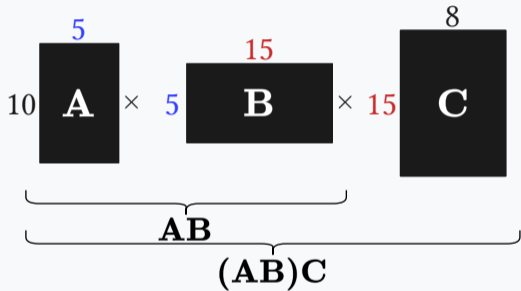
March 3, 2025

Abd-El-Aziz Zayed and Christophe Dubach

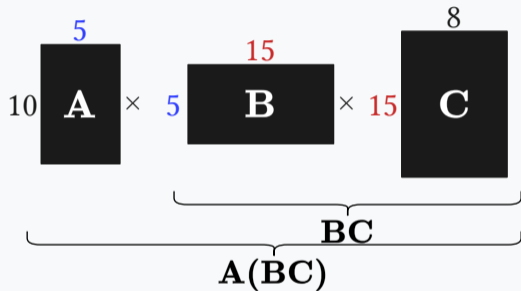
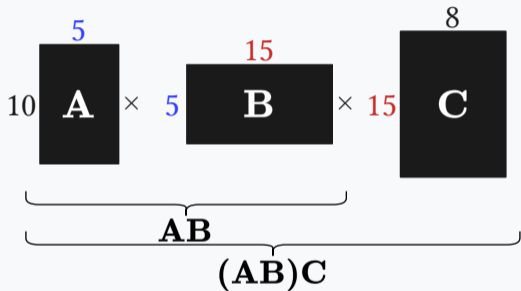


McGill

MatMul is Associative: $(AB)C = A(BC)$

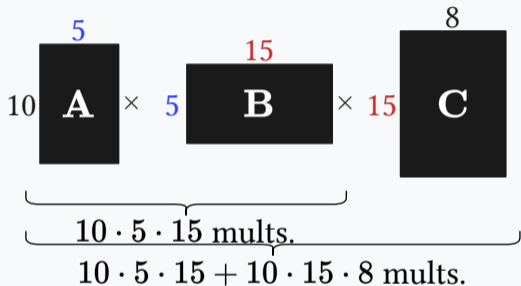


MatMul is Associative: $(AB)C = A(BC)$



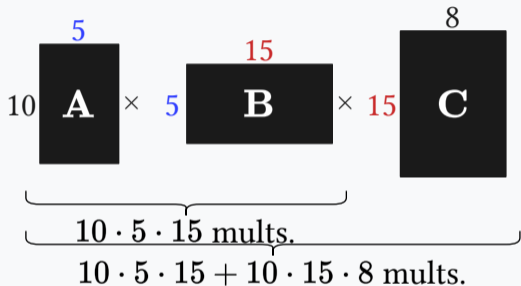
How do we find out which one is more efficient?

The Order of Operations Matters

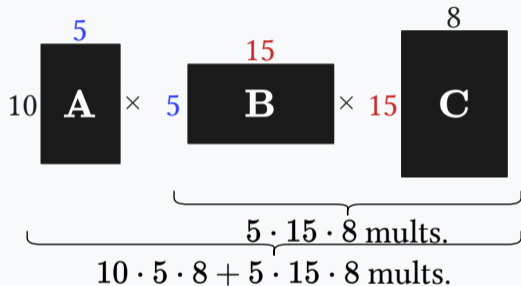


cost $(\text{AB})\text{C} = 1950$
multiplications.

The Order of Operations Matters

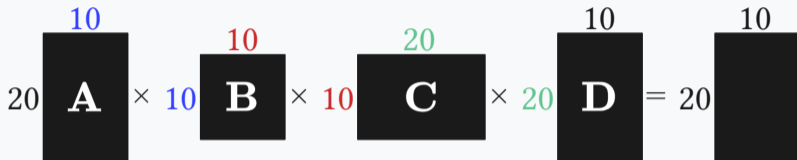


cost $(AB)C = 1950$
multiplications.



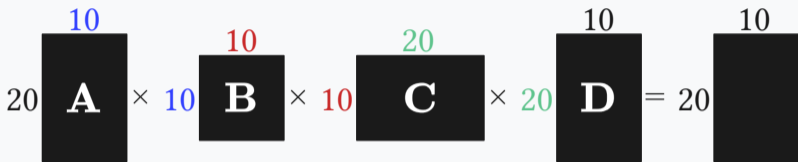
cost $A(BC) = 1000$
multiplications.

Cost of 3MM



- cost $((AB) C) D = 10,000$ multiplications.
- cost $(A (BC)) D = 10,000$ multiplications.
- cost $A ((BC) D) = 6,000$ multiplications.
- **cost $A (B (CD)) = 5,000$ multiplications.**
- cost $(AB) (CD) = 6,000$ multiplications.

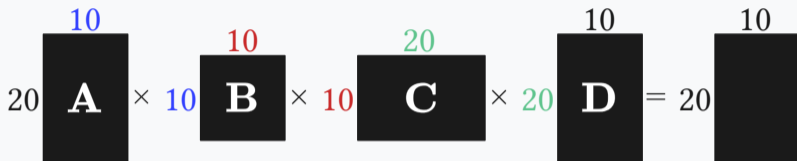
$((AB) C) D$ in MLIR



Why Multi-Level IR?

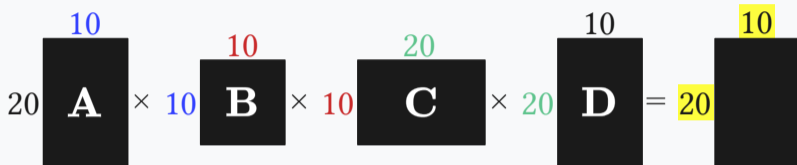
- Perform high-level operations like matrix multiplication.
- Great for domain-specific optimizations.
- Domain-specific operations are grouped into MLIR dialects.
- MatMul in the linear algebra dialect: `linalg.matmul`.

$((AB) C) D$ in Pseudo-MLIR



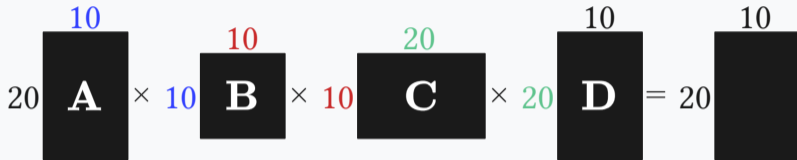
```
1 func.func @3mm(%A: tensor<20x10xi64>, %B: tensor<10x10xi64>,
2             %C: tensor<10x20xi64>, %D: tensor<20x10xi64>)
3             -> tensor<20x10xi64> {
4
5             %AB    = linalg.matmul %A    %B -> tensor<20x10xi64>
6             %ABC   = linalg.matmul %AB   %C -> tensor<20x20xi64>
7             %ABCD  = linalg.matmul %ABC  %D -> tensor<20x10xi64>
8
9             func.return %ABCD
10 }
```


$((AB) C) D$ in Pseudo-MLIR



```
1 func.func @3mm(%A: tensor<20x10xi64>, %B: tensor<10x10xi64>,
2             %C: tensor<10x20xi64>, %D: tensor<20x10xi64>)
3             -> tensor<20x10xi64> {
4
5             %AB = linalg.matmul %A %B -> tensor<20x10xi64>
6             %ABC = linalg.matmul %AB %C -> tensor<20x20xi64>
7             %ABCD = linalg.matmul %ABC %D -> tensor<20x10xi64>
8
9             func.return %ABCD
10 }
```

$((AB) C) D$ in Pseudo-MLIR



```

1 func.func @3mm(%A: tensor<20x10xi64>, %B: tensor<10x10xi64>,
2             %C: tensor<10x20xi64>, %D: tensor<20x10xi64>)
3             -> tensor<20x10xi64> {
4
5             %AB = linalg.matmul %A %B -> tensor<20x10xi64>
6             %ABC = linalg.matmul %AB %C -> tensor<20x20xi64>
7             %ABCD = linalg.matmul %ABC %D -> tensor<20x10xi64>
8
9             func.return %ABCD
10 }

```

Associativity Rewrite in MLIR

How can we express the rewrite $(\mathbf{XY}) \mathbf{Z} \Leftrightarrow \mathbf{X} (\mathbf{YZ})$ in an MLIR pass?

- >120 lines of C++ code.
- Hard to write and maintain.

Likewise for the cost model.

Associativity Rewrite

How we want to express the rewrite $(\mathbf{XY}) \mathbf{Z} \Leftrightarrow \mathbf{X} (\mathbf{YZ})$:

```
1 cost (MatMul X Y) = (nrows X) * (ncols X) * (ncols Y)
2 rewrite (MatMul (MatMul X Y) Z) <=> (MatMul X (MatMul Y Z))
```

Repeatedly apply the rewrite to the AST until a fixed-point is reached.

Associativity Rewrite

How we want to express the rewrite $(\mathbf{XY}) \mathbf{Z} \Leftrightarrow \mathbf{X} (\mathbf{YZ})$:

```
1 cost (MatMul X Y) = (nrows X) * (ncols X) * (ncols Y)
2 rewrite (MatMul (MatMul X Y) Z) <=> (MatMul X (MatMul Y Z))
```

Repeatedly apply the rewrite to the AST until a fixed-point is reached.

Associativity Rewrite

How we want to express the rewrite $(\mathbf{XY}) \mathbf{Z} \Leftrightarrow \mathbf{X} (\mathbf{YZ})$:

```
1 cost (MatMul X Y) = (nrows X) * (ncols X) * (ncols Y)
2 rewrite (MatMul (MatMul X Y) Z) <=> (MatMul X (MatMul Y Z))
```

Repeatedly apply the rewrite to the AST until a fixed-point is reached.

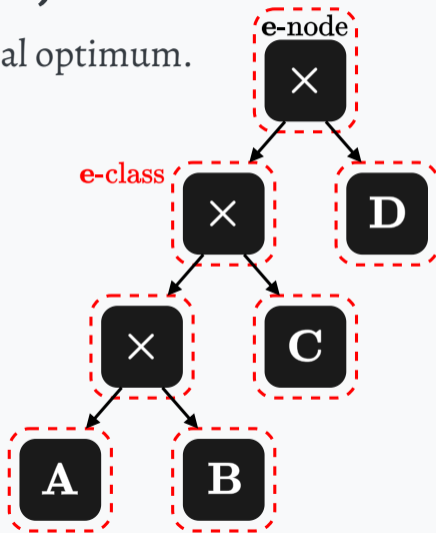
Equality Saturation for $((AB) C) D$

Use **e**quivalence graphs to find the global optimum.

Equality Saturation for ((AB) C) D

Use equivalence graphs to find the global optimum.

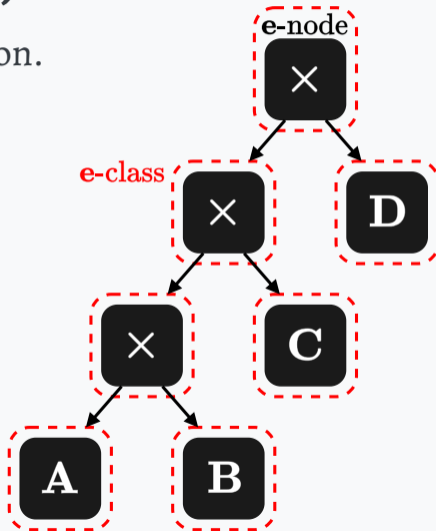
1. Build an **e-graph** of ((**AB**) **C**) **D**
2. Apply the rewrite rule $(\mathbf{XY}) \mathbf{Z} \Leftrightarrow \mathbf{X} (\mathbf{YZ})$ to the e-graph until a fixed-point.
3. Extract the optimized expression via the cost model $\text{cost } \mathbf{XY} = \text{nrows}(\mathbf{X}) \cdot \text{ncols}(\mathbf{X}) \cdot \text{ncols}(\mathbf{Y})$



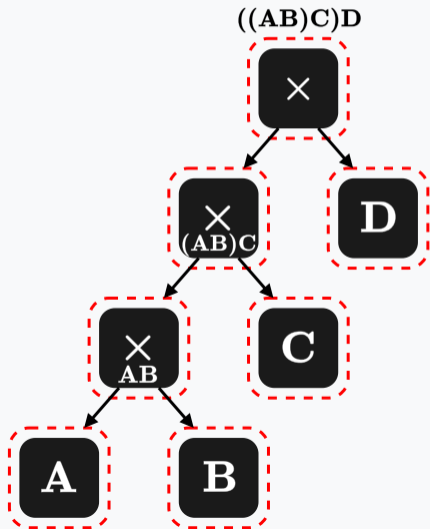
1. Build an e-graph of $((AB) C) D$

AST-like representation of the expression.

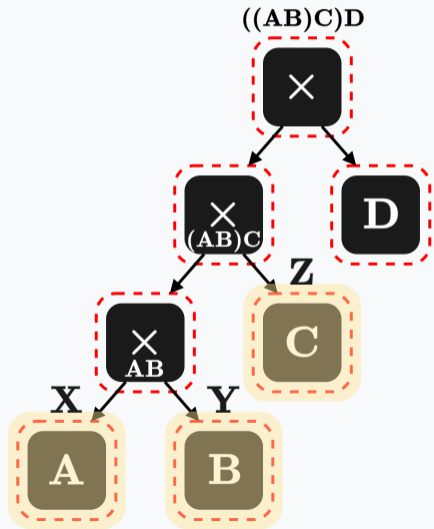
1. **e-node**: Operation node.
2. **e-class**: Set of equivalent e-nodes.
3. **e-graph**: Set of e-classes, capturing equivalent versions of a program compactly.



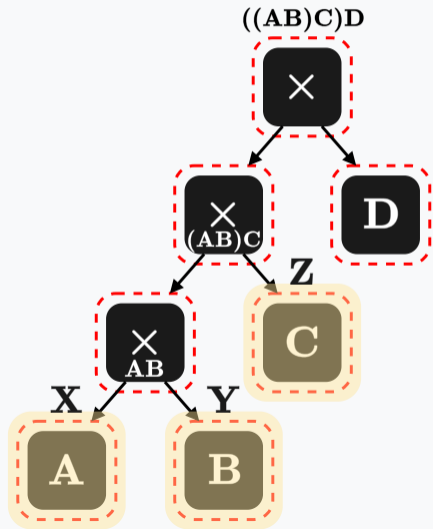
2. Exhaustively Apply the Rewrite $(XY)Z \Leftrightarrow X(YZ)$



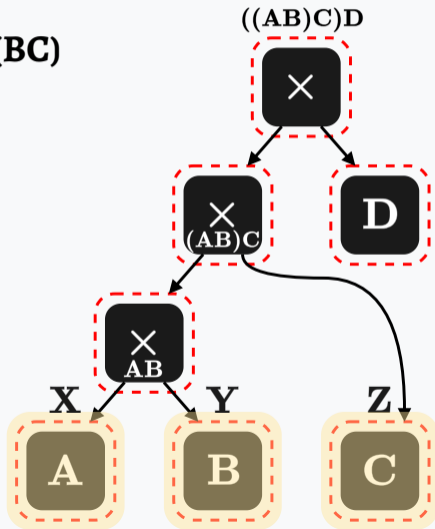
2. Exhaustively Apply the Rewrite $(XY) Z \Leftrightarrow X (YZ)$



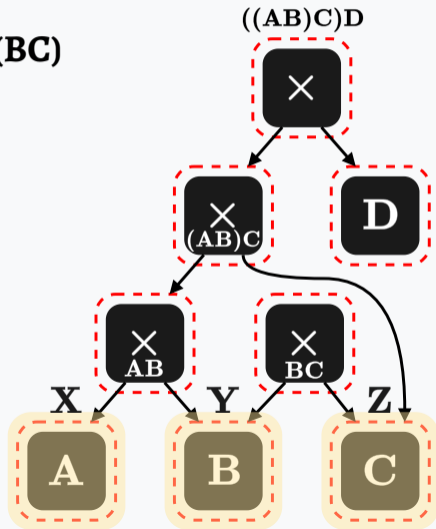
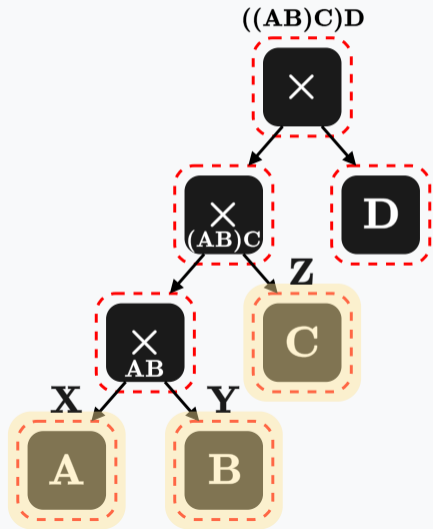
2. Exhaustively Apply the Rewrite $(XY) Z \Leftrightarrow X (YZ)$



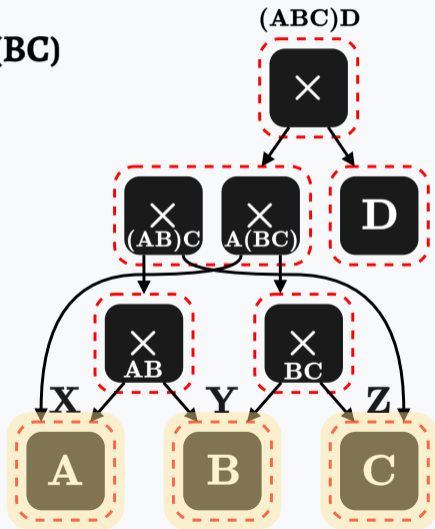
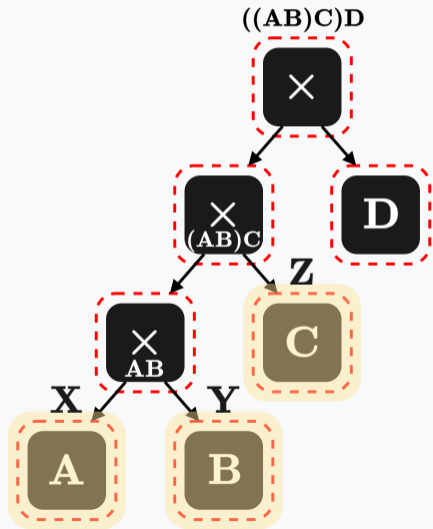
$$(AB) C = A (BC)$$



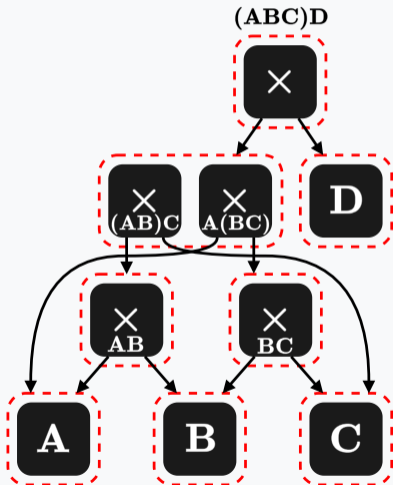
2. Exhaustively Apply the Rewrite $(XY) Z \Leftrightarrow X (YZ)$



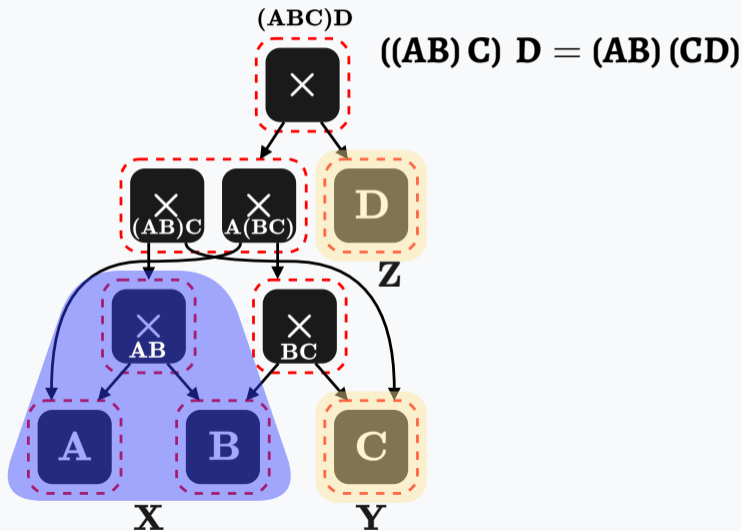
2. Exhaustively Apply the Rewrite $(XY)Z \Leftrightarrow X(YZ)$



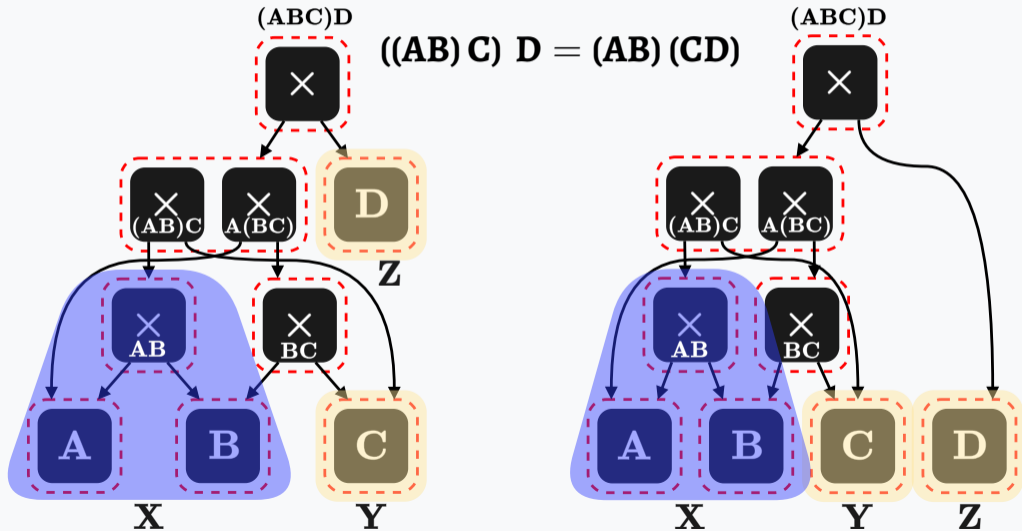
2. Exhaustively Apply the Rewrite $(XY)Z \Leftrightarrow X(YZ)$



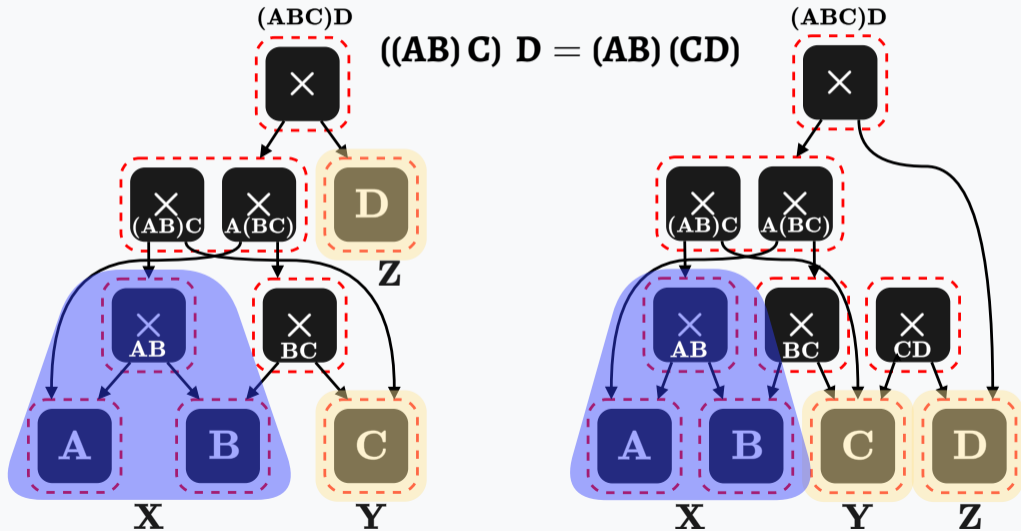
2. Exhaustively Apply the Rewrite $(XY)Z \Leftrightarrow X(YZ)$



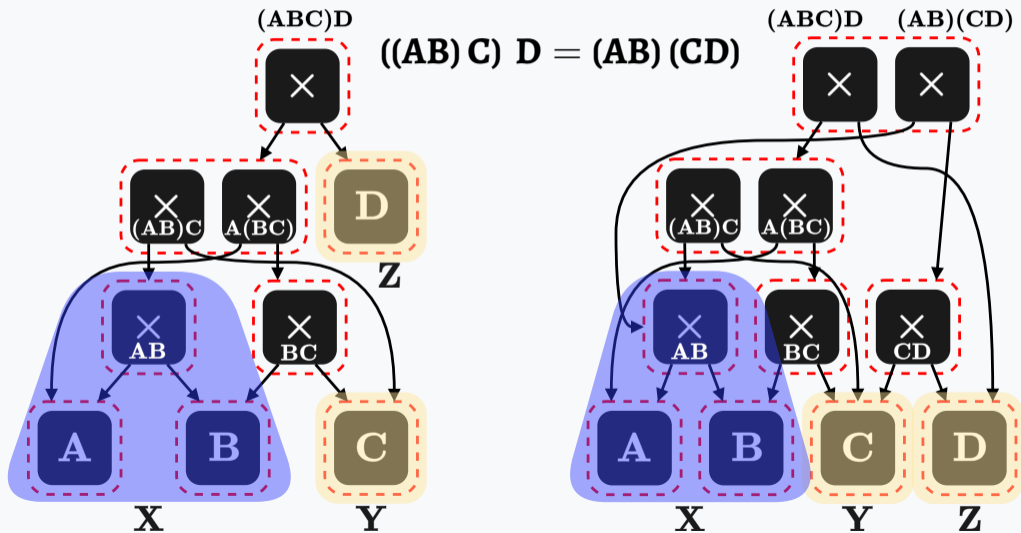
2. Exhaustively Apply the Rewrite $(XY) Z \Leftrightarrow X (YZ)$



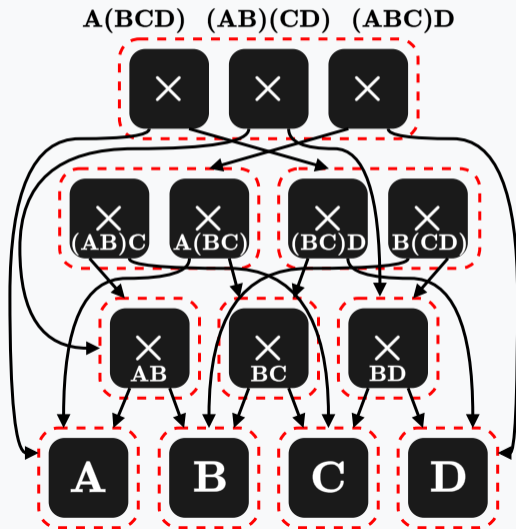
2. Exhaustively Apply the Rewrite $(XY) Z \Leftrightarrow X (YZ)$



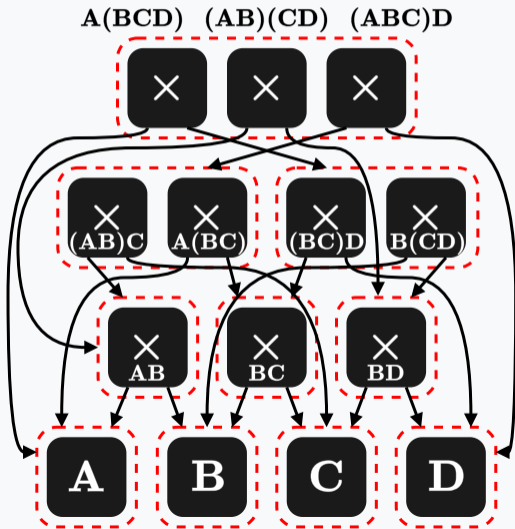
2. Exhaustively Apply the Rewrite $(XY) Z \Leftrightarrow X (YZ)$



Final E-Graph After Saturation

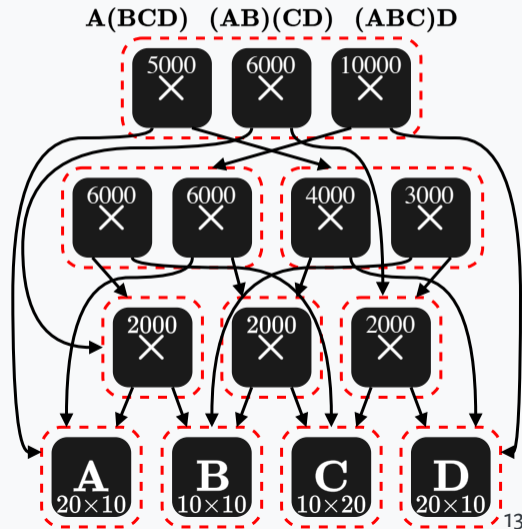
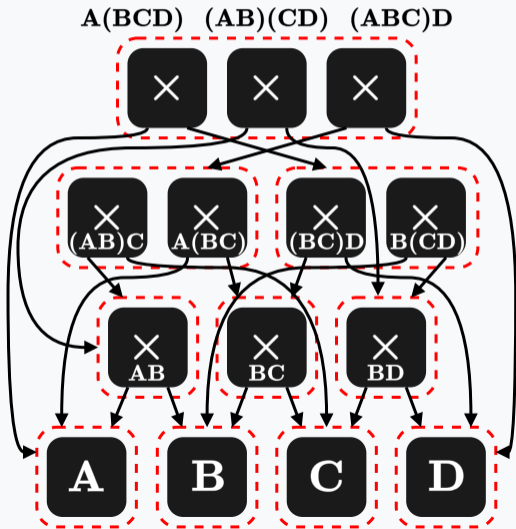


3. Extract the Optimal Expression



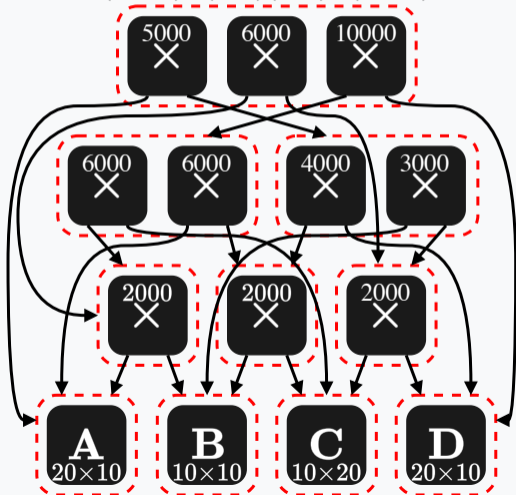
- Root e-class contains equivalent expressions of the program.
- Compute cost of e-nodes in the root e-class.
- Pick the e-node in the root e-class with the lowest cost.

3. Extract the Optimal Expression

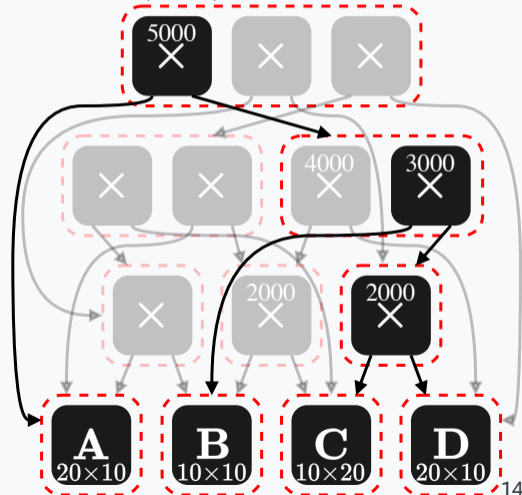


The Optimal Expression is A (B (CD))

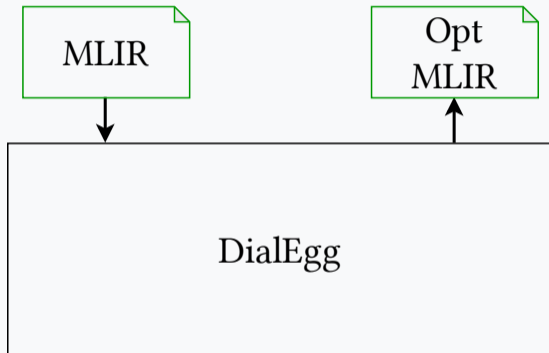
A(BCD) (AB)(CD) (ABC)D



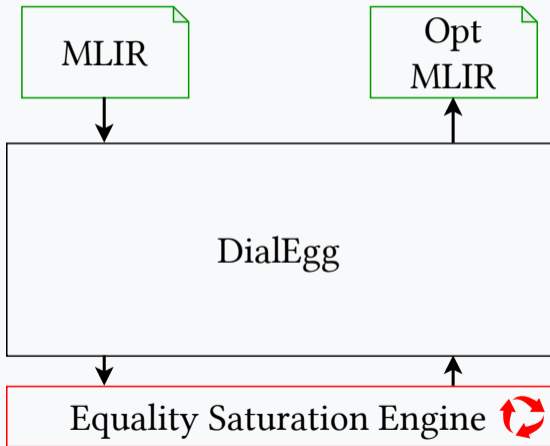
A(BCD)



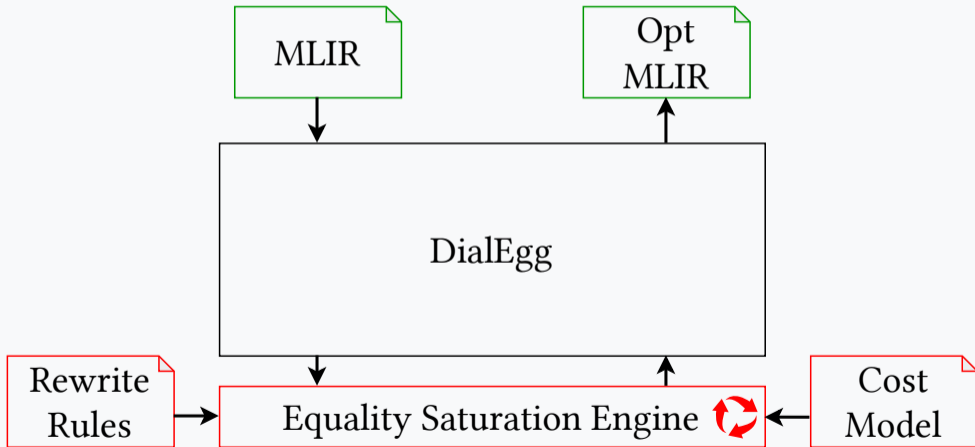
DialEgg



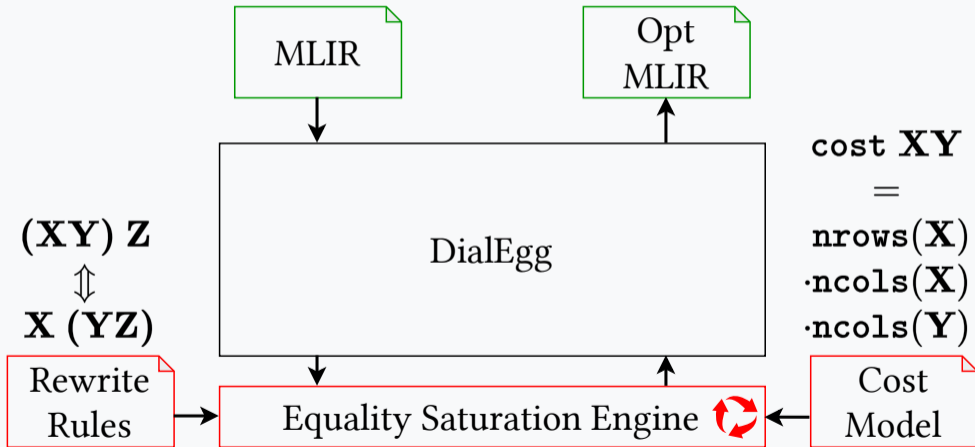
DialEgg



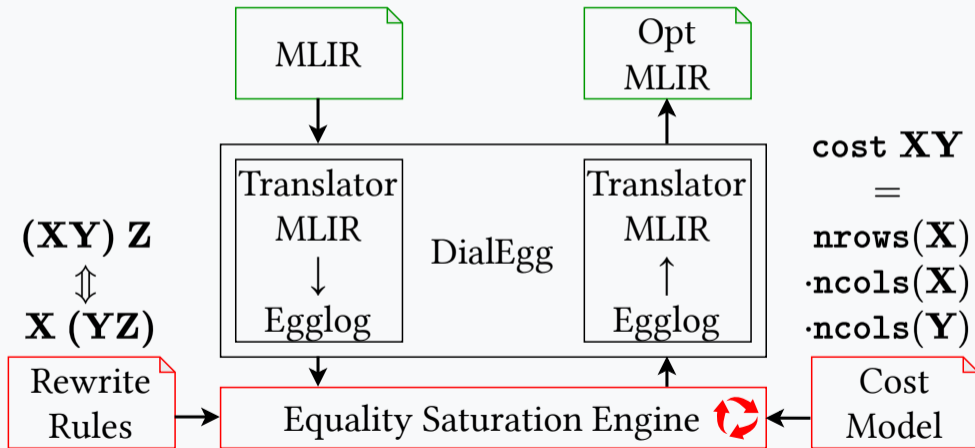
DialEgg



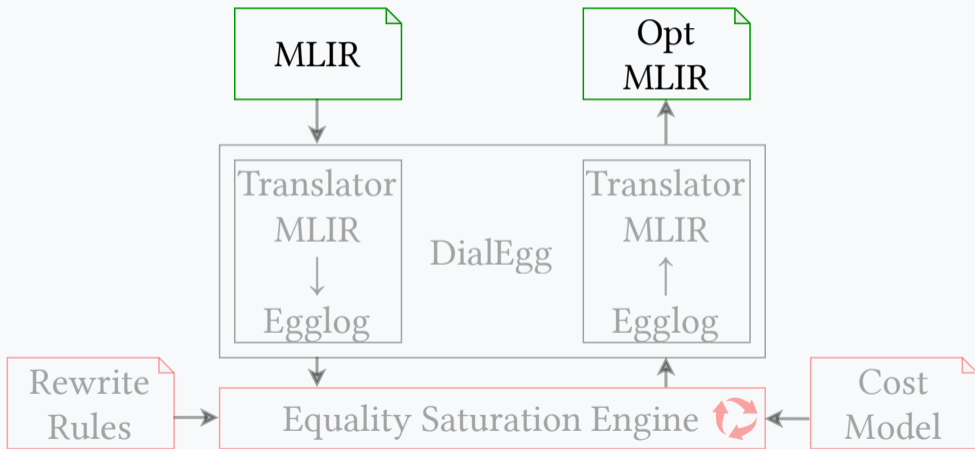
DialEgg



DialEgg



MLIR



MLIR

Operations are first-class citizens. Each operation has at least:

- A dialect and name that identifies the operation.
- A type that defines the type of the result(s).
- A list of operands.

```
1 %AmB = linalg.matmul %A %B : tensor<i64>
2 %ApB = arith.addi %A %B : tensor<i64>
3 %if   = scf.if ... { scf.yield %A }
4       else       { scf.yield %B } : tensor<i64>
```

MLIR

Operations are first-class citizens. Each operation has at least:

- A **dialect and name** that identifies the operation.
- A type that defines the type of the result(s).
- A list of operands.

```
1 %AmB = linalg.matmul %A %B : tensor<i64>
2 %ApB = arith.addi %A %B : tensor<i64>
3 %if = scf.if ... { scf.yield %A }
4     else          { scf.yield %B } : tensor<i64>
```

MLIR

Operations are first-class citizens. Each operation has at least:

- A dialect and name that identifies the operation.
- A **type** that defines the type of the result(s).
- A list of operands.

```
1 %AmB = linalg.matmul %A %B : tensor<i64>
2 %ApB = arith.addi %A %B : tensor<i64>
3 %if   = scf.if ... { scf.yield %A }
4       else      { scf.yield %B } : tensor<i64>
```

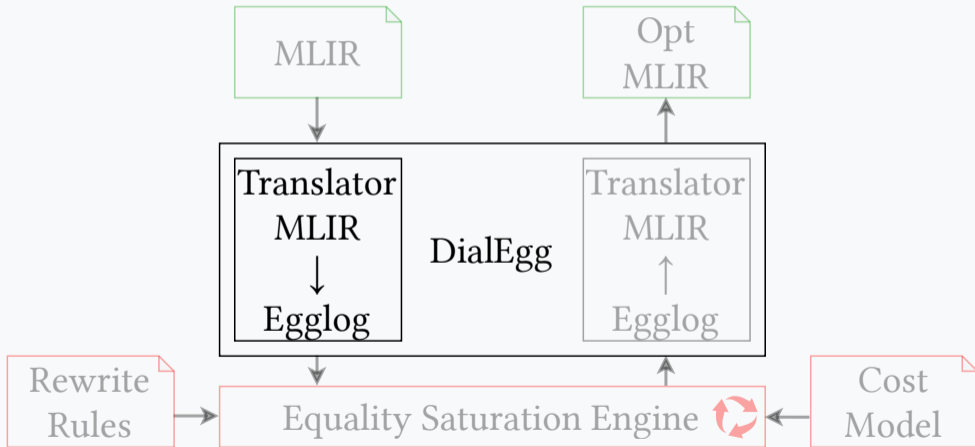

MLIR

Operations are first-class citizens. Each operation has at least:

- A dialect and name that identifies the operation.
- A type that defines the type of the result(s).
- A list of **operands**.

```
1 %AmB = linalg.matmul %A %B : tensor<i64>
2 %ApB = arith.addi %A %B : tensor<i64>
3 %if = scf.if ... { scf.yield %A }
4     else { scf.yield %B } : tensor<i64>
```

MLIR → Egglog



MLIR → Egglog: Types

DialEgg provides most built-in MLIR types in Egglog.

```
1 (datatype Type
2   (F32)
3   (I64)
4   (Tensor IntVec Type)
5   (OpaqueType String String)
6   ...
7 )
```

MLIR → Egglog: Types

DialEgg provides most built-in MLIR types in Egglog.

```
1 (datatype Type
2   (F32)
3   (I64)
4   (Tensor IntVec Type)
5   (OpaqueType String String)
6   ...
7 )
```

Tensor Type in Egglog

```
1 (Tensor (vec-of 2 3) (I64)) ; tensor<2x3xi64>
```

MLIR → Egglog: Operations

Operations are translated piece by piece to Egglog.

```
1 (datatype Op
2   (Value i64 Type) ; (id, type)
3   (arith_const i64 Type) ; simplified constant
4   (arith_addi Op Op Type)
5   (linalg_matmul Op Op Type)
6   (scf_if Op Region Region Type)
7   ...
8 )
```

MLIR → Egglog: Operations

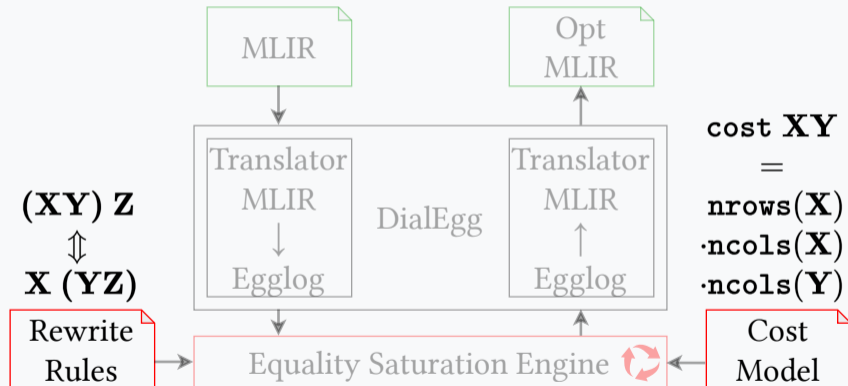
Operations are translated piece by piece to Egglog.

```
1 (datatype Op
2   (Value i64 Type) ; (id, type)
3   (arith_const i64 Type) ; simplified constant
4   (arith_addi Op Op Type)
5   (linalg_matmul Op Op Type)
6   (scf_if Op Region Region Type)
7   ...
8 )
```

MatMul Operation in Egglog

```
1 (linalg_matmul A B (Tensor (vec-of 20 10) (I64)))
```

MM Associativity Rewrite Rule in Egglog



MM Associativity Rewrite Rule in Egglog

$$(XY)Z \Rightarrow X(YZ)$$

where $\mathbf{X}: a \times b$, $\mathbf{Y}: b \times c$, and $\mathbf{Z}: c \times d$

```
1 (rule
2   ((= ?LHS (linalg_matmul (linalg_matmul ?X ?Y ?XY_t) ?Z ?XYZ_t))
3     (= ?b (nrows type-of ?Y))
4     (= ?d (ncols type-of ?Z))
5     (= ?XYZ_t (Tensor ?_ ?t))))
6
7   ((let YZ_t (Tensor (vec-of ?b ?d) ?t)
8     (union ?LHS
9       (linalg_matmul ?X (linalg_matmul ?Y ?Z ?YZ_t) ?XYZ_t))))
10 )
```


MM Associativity Rewrite Rule in Egglog

$$(XY)Z \Rightarrow X(YZ)$$

where $\mathbf{X}: a \times b$, $\mathbf{Y}: b \times c$, and $\mathbf{Z}: c \times d$

```

1 (rule
2   ((= ?LHS (linalg_matmul (linalg_matmul ?X ?Y ?XY_t) ?Z ?XYZ_t))
3     (= ?b (nrows type-of ?Y))
4     (= ?d (ncols type-of ?Z))
5     (= ?XYZ_t (Tensor ?_ ?t)))
6
7   ((let YZ_t (Tensor (vec-of ?b ?d) ?t)
8     (union ?LHS
9       (linalg_matmul ?X (linalg_matmul ?Y ?Z ?YZ_t) ?XYZ_t))))
10 )

```

MM Associativity Rewrite Rule in Egglog

$$(XY)Z \Rightarrow X(YZ)$$

where $\mathbf{X}: a \times b$, $\mathbf{Y}: b \times c$, and $\mathbf{Z}: c \times d$

```
1 (rule
2   ((= ?LHS (linalg_matmul (linalg_matmul ?X ?Y ?XY_t) ?Z ?XYZ_t))
3     (= ?b (nrows type-of ?Y))
4     (= ?d (ncols type-of ?Z))
5     (= ?XYZ_t (Tensor ?_ ?t))))
6
7   ((let YZ_t (Tensor (vec-of ?b ?d) ?t)
8     (union ?LHS
9       (linalg_matmul ?X (linalg_matmul ?Y ?Z ?YZ_t) ?XYZ_t))))
10 )
```

MM Associativity Cost Model in Egglog

cost **XY**

where **X**: $a \times b$ and **Y**: $b \times c$

```
1 (rule ((linalg_matmul ?X ?Y ?XY_t)
2   (= ?a (nrows type-of ?X))
3   (= ?b (ncols type-of ?X))
4   (= ?c (ncols type-of ?Y)))
5
6   ((cost (linalg_matmul ?X ?Y ?XY_t) (* (* ?a ?b) ?c)))
7 )
```

MM Associativity Cost Model in Egglog

cost **XY**

where **X**: $a \times b$ and **Y**: $b \times c$

```
1 (rule ((linalg_matmul ?X ?Y ?XY_t)
2   (= ?a (nrows type-of ?X))
3   (= ?b (ncols type-of ?X))
4   (= ?c (ncols type-of ?Y)))
5
6   ((cost (linalg_matmul ?X ?Y ?XY_t) (* (* ?a ?b) ?c)))
7 )
```

Optimized 3MM

```
1 func.func @3mm(%A,%B,%C,%D) {  
2   %AB =linalg.matmul %A %B  
3   %ABC =linalg.matmul %AB %C  
4   %ABCD=linalg.matmul %ABC %D  
5   func.return %ABCD  
6 }
```

cost $((\mathbf{AB}) \mathbf{C}) \mathbf{D} = 10,000$ mults.

Optimized 3MM

```
1 func.func @3mm(%A,%B,%C,%D) {  
2   %AB =linalg.matmul %A %B  
3   %ABC =linalg.matmul %AB %C  
4   %ABCD=linalg.matmul %ABC %D  
5   func.return %ABCD  
6 }
```

cost **((AB) C) D** = 10,000 mults.

```
1 func.func @3mm(%A,%B,%C,%D) {  
2   %CD =linalg.matmul %C %D  
3   %BCD =linalg.matmul %B %CD  
4   %ABCD=linalg.matmul %A %BCD  
5   func.return %ABCD  
6 }
```

cost **A (B (CD))** = 5,000 mults.

Case Study: Polynomial Evaluation

Reduce from $O(n^2)$ to $O(n)$ multiplications using Horner's method.

$$\begin{aligned} P(x) &= a_0 + a_1x + a_2x^2 + \cdots + a_{n-1}x^{n-1} + a_nx^n \\ &= a_0 + x(a_1 + x(a_2 + \cdots + x(a_{n-1} + xa_n))) \end{aligned}$$

Case Study: Polynomial Evaluation

Reduce from $O(n^2)$ to $O(n)$ multiplications using Horner's method.

$$\begin{aligned} P(x) &= a_0 + a_1x + a_2x^2 + \dots + a_{n-1}x^{n-1} + a_nx^n \\ &= a_0 + x(a_1 + x(a_2 + \dots + x(a_{n-1} + xa_n))) \end{aligned}$$

- Exponentiation: $x^0 \Rightarrow 1$ and $x^n \Leftrightarrow x \cdot x^{n-1}$

Case Study: Polynomial Evaluation

Reduce from $O(n^2)$ to $O(n)$ multiplications using Horner's method.

$$\begin{aligned} P(x) &= a_0 + a_1x + a_2x^2 + \cdots + a_{n-1}x^{n-1} + a_nx^n \\ &= a_0 + x(a_1 + x(a_2 + \cdots + x(a_{n-1} + xa_n))) \end{aligned}$$

- Exponentiation: $x^0 \Rightarrow 1$ and $x^n \Leftrightarrow x \cdot x^{n-1}$
- Commutativity: $x + y \Leftrightarrow y + x$ and $x \cdot y \Leftrightarrow y \cdot x$
- Associativity: $(x + y) + z \Leftrightarrow x + (y + z)$ and $(x \cdot y) \cdot z \Leftrightarrow x \cdot (y \cdot z)$
- Distributivity: $x \cdot (y + z) \Leftrightarrow x \cdot y + x \cdot z$
- Identity: $x \cdot 1 \Rightarrow x$

Case Study: Polynomial Evaluation

Reduce from $O(n^2)$ to $O(n)$ multiplications using Horner's method.

$$\begin{aligned} P(x) &= a_0 + a_1x + a_2x^2 + \cdots + a_{n-1}x^{n-1} + a_nx^n \\ &= a_0 + x(a_1 + x(a_2 + \cdots + x(a_{n-1} + xa_n))) \end{aligned}$$

- Exponentiation: $x^n \Leftrightarrow x \cdot x^{n-1}$

```
1 (rule
2   ((= ?lhs (math_powf ?x (arith_const ?n ?t) ?t))
3     (>= ?n 1))
4
5   ((union ?lhs
6     (arith_mulf ?x (math_powf ?x (arith_const (- ?n 1) ?t)) ?t))
7   )
8 )
```

DialEgg



azizzayed.com

- More case studies are available in the paper:
 - Image conversion: rewrites with computation.
 - Vector normalization: rewrites matching on MLIR attributes.
- DialEgg is open-source.
- I look forward to see how you use DialEgg.

MM Associativity Helpers

```
1 (rule ((= ?A (Value ?id ?t))) ((set (type-of ?A) ?t)))
2 (rule ((= ?A (linalg_matmul ?x ?y ?t))) ((set (type-of ?A) ?t)))
3
4 (function nrows (Type) i64)
5 (function ncols (Type) i64)
6 (rule ((= ?t (RankedTensor ?dims ?tp)))
7       ((set (nrows ?t) (vec-get ?dims 0))
8            (set (ncols ?t) (vec-get ?dims 1))))
9
10 (rule ((linalg_matmul ?x ?y ?t)
11       (= a (nrows (type-of ?x)))
12       (= b (ncols (type-of ?x)))
13       (= c (ncols (type-of ?y))))
14
15       ((cost (linalg_matmul ?x ?y ?t) (* (* a b) c))))
```

Poly Rewrite Rules

```

1 (rewrite (arith_addf ?x ?y ?a ?t) (arith_addf ?y ?x ?a ?t))
2 (rewrite (arith_mulf ?x ?y ?a ?t) (arith_mulf ?y ?x ?a ?t))
3 (rewrite ; (x + y) + z = x + (y + z)
4   (arith_addf (arith_addf ?x ?y ?a ?t) ?z ?a ?t)
5   (arith_addf ?x (arith_addf ?y ?z ?a ?t) ?a ?t))
6 (rewrite ; (x * y) * z = x * (y * z)
7   (arith_mulf (arith_mulf ?x ?y ?a ?t) ?z ?a ?t)
8   (arith_mulf ?x (arith_mulf ?y ?z ?a ?t) ?a ?t))
9
10 (rewrite (arith_mulf ?x ; x * 1 = x
11   (arith_constant (NamedAttr "value" (FloatAttr 1.0 ?t)) ?t)
12   ?a ?t) ?x)
13 (rewrite (math_powf ?x ; x^0 = 1
14   (arith_constant (NamedAttr "value" (FloatAttr 0.0 ?t)) ?t) ?a
15   ?t)
16   (arith_constant (NamedAttr "value" (FloatAttr 1.0 ?t)) ?t))

```

Poly Rewrite Rules

```

1 (rewrite ; mx + nx = x(m + n)
2   (arith_addf (arith_mulf ?m ?x ?a ?t) (arith_mulf ?n ?x ?a ?t)
3     ?a ?t)
4   (arith_mulf ?x (arith_addf ?m ?n ?a ?t) ?a ?t))
5 (rule ((= ?lhs (math_powf ?x ; x^n = x * x^(n - 1)
6   (arith_constant (NamedAttr "value" (FloatAttr ?n ? t)) ?t)
7   ?a ?t)) (>= ?n 1.0))
8 ((union ?lhs
9   (arith_mulf ?x (math_powf ?x
10    (arith_constant
11      (NamedAttr "value" (FloatAttr (- ?n 1.0) ?t))
12      ?t)
13    ?a ?t) ?a ?t)
14 ))
15 )

```

Compile Time and Scalability

Bench.	#Rules	#Ops	MLIR		Saturat.	Canon.
			↓↑ Egglog	Egglog		
Img Conv	1	29	0.4ms	14.6ms	<0.1ms	0.2ms
Vec Norm	1	44	0.5ms	21.6ms	<0.1ms	0.2ms
Poly	8	26	0.5ms	18.9ms	2ms	0.2ms
3MM	5	8	0.3ms	8.7ms	1ms	0.1ms
10MM	5	22	0.5ms	14.4ms	4ms	0.1ms
20MM	5	42	1.0ms	41.3ms	23ms	0.2ms
40MM	5	82	1.8ms	296.2ms	235ms	0.3ms
80MM	5	162	7.3ms	4939.3ms	3732ms	0.6ms

Qualitative Analysis

