# RISE in MLIR

## A functional Pattern-based Dialect

**Martin Lücke** | Michel Steuwer | Aaron Smith

THE UNIVERSITY of EDINBURGH | University of Glasgow

# Why RISE?



**Machine Learning Systems are Stuck in a Rut**

[HotOS'19]

Paul Barham
Google Brain

Michael Isard
Google Brain

Original authors
of TensorFlow

## Abstract

In this paper we argue that systems for numerical computing are stuck in a local basin of performance and programmability. Systems researchers are doing an excellent job improving the performance of 5-year-old benchmarks, but gradually making it harder to explore innovative machine learning research ideas.

We explain how the evolution of hardware accelerators favors compiler back ends that hyper-optimize large monolithic kernels, show how this reliance on high-performance but inflexible kernels reinforces the dominant style of programming model, and argue these programming abstractions lack expressiveness, maintainability, and modularity; all of which hinders research progress.

We conclude by noting promising directions in the field, and advocate steps to advance progress towards high-performance general purpose numerical computing systems on modern accelerators.
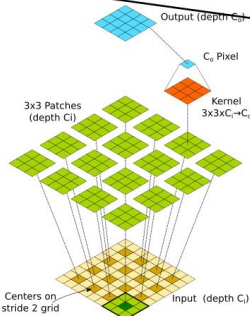
**Figure 1.** Conv2D operation with 3×3 kernel, stride=2

with 16 times fewer training parameters than the convolutional neural network (CNN) we were comparing it to, implementations in both TensorFlow[2] and PyTorch[3] were much slower and ran out of memory with much smaller models. We wanted to understand why.

## 1.1 New ideas often require new primitives

We won't discuss the full details of Capsule networks in this paper[1], but for our purposes it is sufficient to consider a simplified form of the inner loop, which is

2

# Why RISE?



## Machine Learning Systems are Stuck in a Rut

Paul Barham
Google Brain

Michael Isard
Google Brain

### Abstract

In this paper we argue that systems for numerical computing are stuck in a local basin of performance and programmability. Systems researchers are doing an excellent job improving the performance of 5-year-old benchmarks, but gradually making it harder to explore innovative machine learning research ideas.

We explain how the evolution of hardware accelerators favors compiler back ends that hyper-optimize large monolithic kernels, show how this reliance on high-performance but inflexible kernels reinforces the dominant style of programming model, and argue these programming abstractions lack expressiveness, maintainability, and modularity; all of which hinders research progress.

We conclude by noting promising directions in the field, and advocate steps to advance progress towards high-performance general purpose numerical computing systems on modern accelerators.

**Figure 1.** Conv2D operation with 3×3 kernel, stride=2

with 16 times fewer training parameters than the convolutional neural network (CNN) we were comparing it to, implementations in both TensorFlow[2] and PyTorch[3] were much slower and ran out of memory with much smaller models. We wanted to understand why.

### 1.1 New ideas often require new primitives

We won't discuss the full details of Capsule networks in this paper[1], but for our purposes it is sufficient to consider a simplified form of the inner loop, which is

3

# Why RISE?

Machine Learning Systems are Stuck in a Rut

Paul Barham
Google Brain

Michael Isard
Google Brain

**Abstract**

In this paper we argue that systems for numerical computing are stuck in a local basin of performance and programmability. Systems researchers are doing an excellent job improving the performance of 5-year-old benchmarks, but gradually making it harder to explore innovative machine learning research ideas.

We explain how the evolution of hardware accelerators favors compiler back ends that hyper-optimize large monolithic kernels, show how this reliance on high-performance but inflexible kernels reinforces the dominant style of programming model, and argue these programming abstractions lack expressiveness, maintainability, and modularity; all of which hinders research progress.

We conclude by noting promising directions in the field, and advocate steps to advance progress towards

**Figure 1.** Conv2D operation with 3×3 kernel, stride=2

with 16 times fewer training parameters than the convo-

smaller models. We wanted to understand why.

**1.1 New ideas often require new primitives**

We won't discuss the full details of Capsule networks in this paper[1], but for our purposes it is sufficient to consider a simplified form of the inner loop, which is

We should aim for more principled higher level intermediate representations

# RISE - A functional pattern-based data-parallel language

- RISE ([https://rise-lang.org/](https://rise-lang.org/)) is a spiritual successor to the Lift project

- Computations represented using compositions of flexible and generic patterns

- Rewriting system to explore optimization choices

# RISE by example: Matrix Multiplication

```
fun(A : N.K.float ⇒ fun(B : K.M.float ⇒
  A ▷ map(fun(arow ⇒
    B ▷ transpose ▷ map(fun(bcol ⇒
      zip(arow, bcol) ▷ map(*) ▷ reduce(+, 0) )) )) ))
```

# RISE by example: Matrix Multiplication

```
fun(A : N.K.float ⇒ fun(B : K.M.float ⇒
  A ▷ map(fun(arow ⇒
    B ▷ transpose ▷ map(fun(bcol ⇒
      zip(arow, bcol) ▷ map(*) ▷ reduce(+, 0) )) )) ))
```

C          A

B

dot product computation:

$$\sum arow_i * bcol_i$$

# RISE by example: Matrix Multiplication

```
fun(A : N.K.float ⇒ fun(B : K.M.float ⇒
  A ▷ map(fun(arow ⇒
    B ▷ transpose ▷ map(fun(bcol ⇒
      zip(arow, bcol) ▷ map(*) ▷ reduce(+, 0) )) )) ))
```

# Optimization choices for Matrix Multiplication

```
fun(A : N.K.float ⇒ fun(B : K.M.float ⇒
  A ▷ map(fun(arow ⇒
    B ▷ transpose ▷ map(fun(bcol ⇒
      zip(arow, bcol) ▷ map(*) ▷ reduce(+, 0) )) )) ))
```

RISE compiler

# Optimization choices for Matrix Multiplication

```
fun(A : N.K.float ⇒ fun(B : K.M.float ⇒
  A ▷ map(fun(arow ⇒
    B ▷ transpose ▷ map(fun(bcol ⇒
      zip(arow, bcol) ▷ map(*) ▷ reduce(+, 0) )) )) ))
```
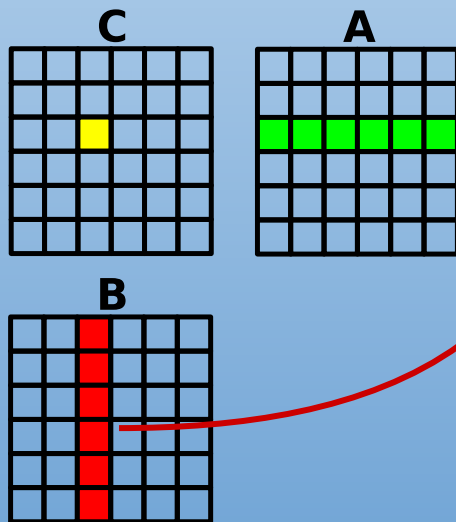
RISE compiler

**Rewrite Rules**



```
map(f, A) ⟼ join(map(map(f),
              split(n, A)))
```

# RISE Compilation

Compilation is divided into two steps:

1. **Translation from functional to imperative low level program**

2. **Produce generated code**

Compilation described formally in:

**Strategy Preserving Compilation for Parallel Functional Code** *by Robert Atkey, Michel Steuwer, Sam Lindley, and Christophe Dubach.*
https://arxiv.org/abs/1710.08332

# High Performance Results for Matrix Multiplication



Only few generated code with very good performance



Still: One can expect to find a good performing kernel quickly!



Performance close or better than hand-tuned library code

[GPGPU'16]

12

# RISE - The Caveats

- Does everyone have to write functional programs now?

- Academic work written in Scala

- Does not integrate well with existing compiler infrastructures

# MLIR - Multi-Level Intermediate Representation

- Extensible infrastructure to define compiler intermediate representations

- Dialects can capture different levels of abstraction:

    - High-level domain specific ------------ Hardware specific backend

- Existing dialects available for:

    - TensorFlow / TensorFlow Lite        - Performing polyhedral optimizations

    - Targeting GPUs                      - LLVM IR

    - ...

# MLIR - Martin Lücke Intermediate Representation

- Extensible infrastructure to define compiler intermediate representations

- Dialects can capture different levels of abstraction:

  - High-level domain specific ------------ Hardware specific backend

- Existing dialects available for:

  - TensorFlow / TensorFlow Lite       - Performing polyhedral optimizations

  - Targeting GPUs                     - LLVM IR

  - ...

15

# MLIR - Multi-Level Intermediate Representation

- Extensible infrastructure to define compiler intermediate representations

- Dialects can capture different levels of abstraction:
    - High-level domain specific  ------------ Hardware specific backend

- Existing dialects available for:
    - TensorFlow / TensorFlow Lite        - Performing polyhedral optimizations
    - Targeting GPUs                               - LLVM IR
    - ...

# RISE in MLIR = Lambda Calculus + Patterns in MLIR

- RISE in MLIR opens up opportunities to integrate with other MLIR dialects

- We do not have to write programs in RISE directly, but lower from domain-specific dialects to it

- MLIR is written in C++, using an established toolchain -> widely usable

- Natural integration with existing toolchains  `#include libMLIR`


-> to implement RISE as an MLIR dialect we implement:
- $\lambda$-calculus
- patterns

# RISE dialect by example: Matrix Multiplication

```
1 func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2   %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3   %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4   %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5     %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6       %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7       %zipped = rise.apply %zip, %arow, %bcol
8       %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9         %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10        %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11        %fst = rise.apply %fstFun, %tuple
12        %snd = rise.apply %sndFun, %tuple
13        %result = rise.embed(%fst, %snd) {
14          %res = mulf %fst, %snd : f32
15          return %res : f32
16        } : !rise.scalar<f32>
17        rise.return %result : !rise.scalar<f32>
18      }
19      %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20      %multipliedArray = rise.apply %map, %f, %zipped
21      %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22        %result = rise.embed(%a, %b) {
23          %res = addf %a, %b : f32
24          return %res : f32
25        } : !rise.scalar<f32>
26        rise.return %result : !rise.scalar<f32>
27      }
28      %init = rise.literal #rise.lit<0.0>
29      %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30      %result = rise.apply %reduce, %add, %init, %multipliedArray
31      rise.return %result : !rise.scalar<f32>
32    }
33    %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34    %result = rise.apply %mapB, %f2, %B
35    rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36  }
37  %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38  %result = rise.apply %mapA, %f1, %A
39  rise.out %out ← %result
40  return
41 }
```

# RISE dialect by example: Matrix Multiplication

```
1  func @mm(%out:memref<1024x1024xf32>, %inA:memref<1024x1024xf32>, %inB:memref<1024x1024xf32>) {
2    %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3    %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4    %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5      %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6        %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7        %zipped = rise.apply %zip, %arow, %bcol
8        %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9          %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10         %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11         %fst = rise.apply %fstFun, %tuple
12         %snd = rise.apply %sndFun, %tuple
13         %result = rise.embed(%fst, %snd) {
14           %res = mulf %fst, %snd : f32
15           return %res : f32
16         } : !rise.scalar<f32>
17         rise.return %result : !rise.scalar<f32>
18       }
19       %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20       %multipliedArray = rise.apply %map, %f, %zipped
21       %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22         %result = rise.embed(%a, %b) {
23           %res = addf %a, %b : f32
24           return %res : f32
25         } : !rise.scalar<f32>
26         rise.return %result : !rise.scalar<f32>
27       }
28       %init = rise.literal #rise.lit<0.0>
29       %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30       %result = rise.apply %reduce, %add, %init, %multipliedArray
31       rise.return %result : !rise.scalar<f32>
32     }
33     %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34     %result = rise.apply %mapB, %f2, %B
35     rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36   }
37   %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38   %result = rise.apply %mapA, %f1, %A
39   rise.out %out ← %result
40   return
41 }
```

**Data Types**

**Function Types**

# Types

## Data Types

```
2   %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
```

- Rise data types: array types, tuple types, scalar types

- Nested array types represent higher dimensional data

- Scalar wraps an arbitrary scalar MLIR type

## Function Types

```
8   %f = rise.lambda (%t : !rise.tuple<scalar<f32>, scalar<f32>>) → scalar<f32>
    // %f : !rise.fun<tuple<scalar<f32>, scalar<f32>> → scalar<f32>>
```

- Rise function types: types of lambda expressions and functional patterns

- Type system prevents mixing of function and data types

# Patterns

```
1 func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2   %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3   %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4   %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5     %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6       %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7       %zipped = rise.apply %zip, %arow, %bcol
8       %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9         %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10        %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11        %fst = rise.apply %fstFun, %tuple
12        %snd = rise.apply %sndFun, %tuple
13        %result = rise.embed(%fst, %snd) {
14          %res = mulf %fst, %snd : f32
15          return %res : f32
16        } : !rise.scalar<f32>
17        rise.return %result : !rise.scalar<f32>
18      }
19      %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20      %multipliedArray = rise.apply %map, %f, %zipped
21      %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22        %result = rise.embed(%a, %b) {
23          %res = addf %a, %b : f32
24          return %res : f32
25        } : !rise.scalar<f32>
26        rise.return %result : !rise.scalar<f32>
27      }
28      %init = rise.literal #rise.lit<0.0>
29      %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30      %result = rise.apply %reduce, %add, %init, %multipliedArray
31      rise.return %result : !rise.scalar<f32>
32    }
33    %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34    %result = rise.apply %mapB, %f2, %B
35    rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36  }
37  %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38  %result = rise.apply %mapA, %f1, %A
39  rise.out %out ← %result
40  return
41 }
```

# Patterns

```
1 func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2   %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3   %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4   %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5     %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6       %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7       %zipped = rise.apply %zip, %arow, %bcol
8       %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9         %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10        %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11        %fst = rise.apply %fstFun, %tuple
12        %snd = rise.apply %sndFun, %tuple
13        %result = rise.embed(%fst, %snd) {
14          %res = mulf %fst, %snd : f32
15          return %res : f32
16        } : !rise.scalar<f32>
17        rise.return %result : !rise.scalar<f32>
18      }
19      %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20      %multipliedArray = rise.apply %map, %f, %zipped
21      %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22        %result = rise.embed(%a, %b) {
23          %res = addf %a, %b : f32
24          return %res : f32
25        } : !rise.scalar<f32>
26        rise.return %result : !rise.scalar<f32>
27      }
28      %init = rise.literal #rise.lit<0.0>
29      %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30      %result = rise.apply %reduce, %add, %init, %multipliedArray
31      rise.return %result : !rise.scalar<f32>
32    }
33    %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34    %result = rise.apply %mapB, %f2, %B
35    rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36  }
37  %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38  %result = rise.apply %mapA, %f1, %A
39  rise.out %out ← %result
40  return
41 }
```

# Patterns: zip

```
6    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
```

# Patterns: zip

```
6   %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>

    // type: () → !rise.fun<array<1024, scalar<f32>> →
                        fun<array<1024, scalar<f32>> →
                            array<1024, tuple<scalar<f32>, scalar<f32>>>>

    // type: () →  ▢▢▢   →   ▣▣▣   →   <▢><▢><▢>
```

- Each RISE pattern is implemented as an MLIR operation

- Operations customized with attributes specifying information for the type

- Patterns encoded as operations have a RISE function type

# Function application

```
4

5

6         %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7
```

# Function application

```
4  %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) →
                                    array<1024, scalar<f32>> {
5    %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) →
                                    array<1024, scalar<f32>> {
6        %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7        %zipped = rise.apply %zip, %arow, %bcol
```

- `rise.apply` models function application:

  expects an SSA value with a RISE `function type` (%zip)

  and arguments to the function (%arow, %bcol)

# Function application

```
4  %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) →
                                    array<1024, scalar<f32>> {
5     %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) →
                                       array<1024, scalar<f32>> {
6          %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7          %zipped = rise.apply %zip, %arow, %bcol

 //type: (%zip.type, array<1024, scalar<f32>>, array<1024, scalar<f32>>) →
                                   array<1024,tuple<scalar<f32>,scalar<f32>>>
```

- `rise.apply` models function application:

  expects an SSA value with a RISE function type (%zip)

  and arguments to the function (%arow, %bcol)

# Function abstraction, aka λ-expressions

```
1 func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2   %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3   %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4   %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5     %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6       %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7       %zipped = rise.apply %zip, %arow, %bcol
8       %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9         %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10        %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11        %fst = rise.apply %fstFun, %tuple
12        %snd = rise.apply %sndFun, %tuple
13        %result = rise.embed(%fst, %snd) {
14          %res = mulf %fst, %snd : f32
15          return %res : f32
16        } : !rise.scalar<f32>
17        rise.return %result : !rise.scalar<f32>
18      }
19      %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20      %multipliedArray = rise.apply %map, %f, %zipped
21      %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22        %result = rise.embed(%a, %b) {
23          %res = addf %a, %b : f32
24          return %res : f32
25        } : !rise.scalar<f32>
26        rise.return %result : !rise.scalar<f32>
27      }
28      %init = rise.literal #rise.lit<0.0>
29      %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30      %result = rise.apply %reduce, %add, %init, %multipliedArray
31      rise.return %result : !rise.scalar<f32>
32    }
33    %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34    %result = rise.apply %mapB, %f2, %B
35    rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36  }
37  %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38  %result = rise.apply %mapA, %f1, %A
39  rise.out %out ← %result
40  return
41 }
```

# Function abstraction, aka λ-expressions

```
1  func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2    %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3    %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4    %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5      %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6        %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7        %zipped = rise.apply %zip, %arow, %bcol
8        %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9          %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10         %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11         %fst = rise.apply %fstFun, %tuple
12         %snd = rise.apply %sndFun, %tuple
13         %result = rise.embed(%fst, %snd) {
14           %res = mulf %fst, %snd : f32
15           return %res : f32
16         } : !rise.scalar<f32>
17         rise.return %result : !rise.scalar<f32>
18       }
19       %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20       %multipliedArray = rise.apply %map, %f, %zipped
21       %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22         %result = rise.embed(%a, %b) {
23           %res = addf %a, %b : f32
24           return %res : f32
25         } : !rise.scalar<f32>
26         rise.return %result : !rise.scalar<f32>
27       }
28       %init = rise.literal #rise.lit<0.0>
29       %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30       %result = rise.apply %reduce, %add, %init, %multipliedArray
31       rise.return %result : !rise.scalar<f32>
32     }
33     %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34     %result = rise.apply %mapB, %f2, %B
35     rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36   }
37   %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38   %result = rise.apply %mapA, %f1, %A
39   rise.out %out ← %result
40   return
41 }
```

# Function abstraction, aka λ-expressions

```
21  %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32> {

22    %result = rise.embed(%a, %b) {

23      %res = addf %a, %b : f32
24      return %res
25    } : !rise.scalar<f32>
26    rise.return %result : !rise.scalar<f32>
27  }
```

- **rise.lambda** has an MLIR region of exactly one block

- Arbitrary number of arguments and a result

- **rise.lambda** associates the region with a RISE function type

30

# Function abstraction, aka λ-expressions

```
21  %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32> {

22    %result = rise.embed(%a, %b) {

23      %res = addf %a, %b : f32
24      return %res
25    } : !rise.scalar<f32>
26    rise.return %result : !rise.scalar<f32>
27  } // type: !rise.fun<scalar<f32> → fun<scalar<f32> → scalar<f32>>>
```

- rise.lambda  has an MLIR region of exactly one block

- Arbitrary number of arguments and a result

- rise.lambda  associates the region with a RISE function type

31

# Embedding of other dialects

```
21  %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32> {

22    %result = rise.embed(%a, %b) {

23      %res = addf %a, %b : f32
24      return %res
25    } : !rise.scalar<f32>
26    rise.return %result : !rise.scalar<f32>
27  }
```

- rise.embed unwraps the operands of type !rise.scalar<t> and exposes the values of the wrapped types t inside the region

- The region may contain operations from arbitrary MLIR dialects

- rise.embed returns the result value of !rise.scalar type

# Embedding of other dialects

```
21  %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32> {
       // %a, %b : !rise.scalar<f32>
22     %result = rise.embed(%a, %b) {
         // %a, %b : f32
23       %res = addf %a, %b : f32
24       return %res
25     } : !rise.scalar<f32>
26     rise.return %result : !rise.scalar<f32>
27  }
```

- `rise.embed` unwraps the operands of type `!rise.scalar<t>` and
  exposes the values of the wrapped types `t` inside the region

- The region may contain operations from arbitrary MLIR dialects

- `rise.embed` returns the result value of `!rise.scalar` type

# RISE dialect by example: Matrix Multiplication

```
1 func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2    %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3    %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4    %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5       %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6          %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7          %zipped = rise.apply %zip, %arow, %bcol
8          %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9             %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10            %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11            %fst = rise.apply %fstFun, %tuple
12            %snd = rise.apply %sndFun, %tuple
13            %result = rise.embed(%fst, %snd) {
14               %res = mulf %fst, %snd : f32
15               return %res : f32
16            } : !rise.scalar<f32>
17            rise.return %result : !rise.scalar<f32>
18         }
19         %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20         %multipliedArray = rise.apply %map, %f, %zipped
21         %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22            %result = rise.embed(%a, %b) {
23               %res = addf %a, %b : f32
24               return %res : f32
25            } : !rise.scalar<f32>
26            rise.return %result : !rise.scalar<f32>
27         }
28         %init = rise.literal #rise.lit<0.0>
29         %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30         %result = rise.apply %reduce, %add, %init, %multipliedArray
31         rise.return %result : !rise.scalar<f32>
32      }
33      %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34      %result = rise.apply %mapB, %f2, %B
35      rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36   }
37   %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38   %result = rise.apply %mapA, %f1, %A
39   rise.out %out ← %result
40   return
41 }
```

# RISE dialect by example: Matrix Multiplication

```
1 func @mm(%out:memref<1024×1024xf32>, %inA:memref<1024×1024xf32>, %inB:memref<1024×1024xf32>) {
2   %A = rise.in %inA : !rise.array<1024, array<1024, scalar<f32>>>
3   %B = rise.in %inB : !rise.array<1024, array<1024, scalar<f32>>>
4   %f1 = rise.lambda (%arow : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
5     %f2 = rise.lambda (%bcol : !rise.array<1024, scalar<f32>>) → !rise.array<1024, scalar<f32>> {
6       %zip = rise.zip #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
7       %zipped = rise.apply %zip, %arow, %bcol
8       %f = rise.lambda (%tuple : !rise.tuple<scalar<f32>, scalar<f32>>) → !rise.scalar<f32> {
9         %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
10        %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
11        %fst = rise.apply %fstFun, %tuple
12        %snd = rise.apply %sndFun, %tuple
13        %result = rise.embed(%fst, %snd) {
14          %res = mulf %fst, %snd : f32
15          return %res : f32
16        } : !rise.scalar<f32>
17        rise.return %result : !rise.scalar<f32>
18      }
19      %map = rise.mapSeq #rise.nat<1024> #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
20      %multipliedArray = rise.apply %map, %f, %zipped
21      %add = rise.lambda (%a : !rise.scalar<f32>, %b : !rise.scalar<f32>) → !rise.scalar<f32>> {
22        %result = rise.embed(%a, %b) {
23          %res = addf %a, %b : f32
24          return %res : f32
25        } : !rise.scalar<f32>
26        rise.return %result : !rise.scalar<f32>
27      }
28      %init = rise.literal #rise.lit<0.0>
29      %reduce = rise.reduceSeq #rise.nat<1024> #rise.scalar<f32> #rise.scalar<f32>
30      %result = rise.apply %reduce, %add, %init, %multipliedArray
31      rise.return %result : !rise.scalar<f32>
32    }
33    %mapB = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
34    %result = rise.apply %mapB, %f2, %B
35    rise.return %result : !rise.array<1024, array<1024, scalar<f32>>>
36  }
37  %mapA = rise.mapSeq #rise.nat<1024> #rise.array<1024, scalar<f32>> #rise.array<1024, scalar<f32>>
38  %result = rise.apply %mapA, %f1, %A
39  rise.out %out ← %result
40  return
41 }
```
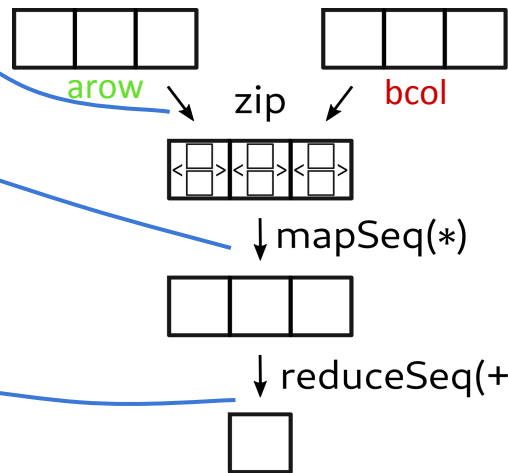
mapSeq (
  mapSeq (

arow      zip      bcol

mapSeq(*)

reduceSeq(+)

  )
)

# RISE dialect - modelling λ-calculus + patterns

**Core λ-calculus:**  `rise.lambda`, `rise.apply`, *`rise.let`*,
                `!rise.array, !rise.tuple, !rise.fun`

**Patterns:** `rise.zip`, `rise.mapSeq`, `rise.reduceSeq`, *`rise.mapPar, rise.reducePar,`*
        `rise.literal rise.tuple`, `rise.fst`, `rise.snd`,
        *`rise.split`, `rise.join`, `rise.transpose`*, **...**

**Interoperability:** `rise.embed, rise.in, rise.out`

**-> direct correspondence of MLIR dialect with original functional representation**

`italic → not implemented yet`

# Lowering of the RISE dialect

**Lowering of patterns in not context free, e.g. zip influences the code generation of the next pattern**

**Lowering is divided into two steps:**

1. **Lowering functional to imperative**

2. **Lowering to final target code**

**Lowering described formally in:**
Strategy Preserving Compilation for Parallel Functional Code
*by Robert Atkey, Michel Steuwer, Sam Lindley, and Christophe Dubach*
https://arxiv.org/abs/1710.08332

# Lowering of the RISE dialect

Lowering of patterns in not context free, e.g. zip influences the code generation of the next pattern
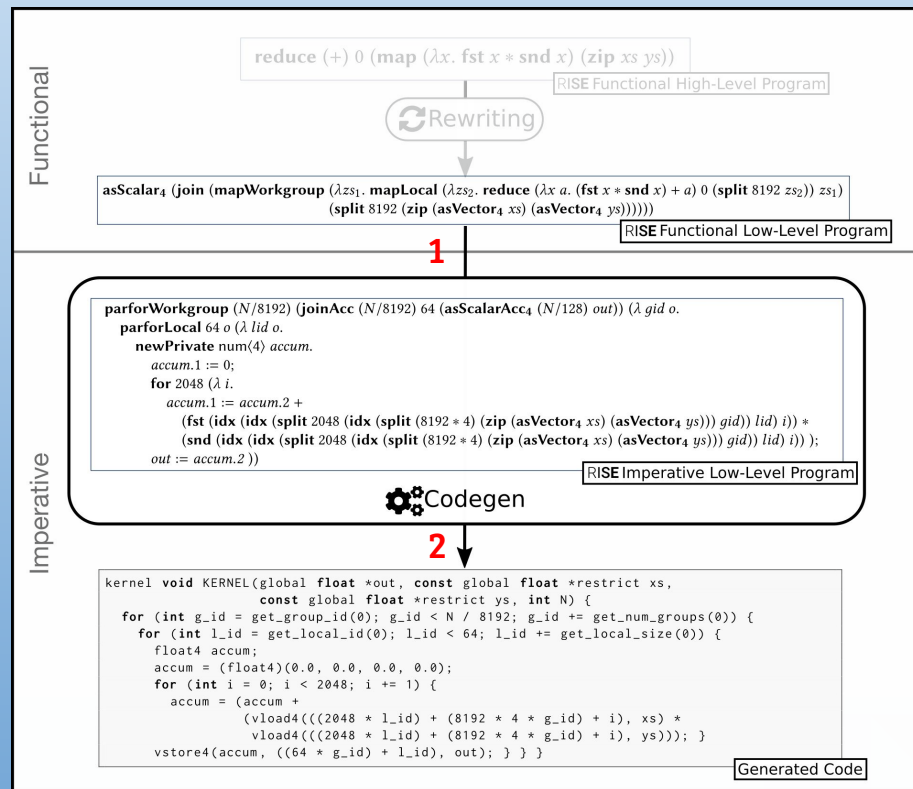
Lowering is divided into two steps:

1. **Lowering functional to imperative**

2. Lowering to final target code

Lowering described formally in:
Strategy Preserving Compilation for Parallel Functional Code
by Robert Atkey, Michel Steuwer, Sam Lindley, and Christophe Dubach
https://arxiv.org/abs/1710.08332

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                 %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
            #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
            #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```
`RISE`

lowering

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```
`Loop + RISE Intermediate`

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
           #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering →

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
           : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

- Start translation from `rise.out`

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
            #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
            #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering →

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                        %b:memref<1024xf32>) {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
          #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
            #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

<span>RISE</span>

lowering
reduceSeq

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

<span>Loop + RISE Intermediate</span>

$$out :=_{\delta_2} \textbf{reduceSeq}(N, \delta_1, \delta_2, F, I, E) \quad = \quad C([I])_{\delta_2}(\lambda\ init.$$
$$C([E])_{N.\delta_1}(\lambda\ x.$$
$$out :=_{\delta_1} init;$$
$$\textbf{for } N\ (\lambda i.\ out :=_{\delta_2} F(x[i], out)); ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>) {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
          #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
**reduceSeq**

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \mathbf{reduceSeq}(N, \delta_1, \delta_2, F, I, E) \quad = \quad C(\!(I)\!)_{\delta_2}(\lambda\ init.$$
$$C(\!(E)\!)_{N.\delta_1}(\lambda\ x.$$
$$out :=_{\delta_1} init;$$
$$\boxed{\mathbf{for}\ N}(\lambda i.\ out :=_{\delta_2} F(x[i], out)); ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                      %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
          #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
          #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
literal

→

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
            : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \textbf{reduceSeq}(N, \delta_1, \delta_2, F, \boxed{I,} E) \quad = \quad \boxed{C(\!(I)\!)}_{\delta_2} (\lambda\ init.$$
$$C(\!(E)\!)_{N.\delta_1} (\lambda\ x.$$
$$\boxed{out :=_{\delta_1} init;}$$
$$\textbf{for } N\ (\lambda i.\ out :=_{\delta_2} F(x[i], out)); ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>) {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
            #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
zip
→

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
            : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \mathbf{reduceSeq}(N, \delta_1, \delta_2, F, I, \boxed{E}) \quad = \quad C(\![I]\!)_{\delta_2} (\lambda\ init.$$
$$\boxed{C(\![E]\!)_{N.\delta_1}(\lambda\ x.}$$
$$out :=_{\delta_1} init;$$
$$\mathbf{for}\ N\ (\lambda i.\ out :=_{\delta_2} F(\boxed{x[i]}, out)); )$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>) {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
         #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```
RISE

lowering
zip

→

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
         : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```
Loop + RISE Intermediate

$$out :=_{\delta_2} \textbf{reduceSeq}(N, \delta_1, \delta_2, F, I, \boxed{E}) \quad = \quad C(\![I]\!)_{\delta_2}(\lambda\ init.$$
$$\boxed{C(\![E]\!)_{N.\delta_1}(\lambda\ x.}$$
$$out :=_{\delta_1} init;$$
$$\textbf{for } N\ (\lambda i.\ out :=_{\delta_2} F(\boxed{x[i]}, out)); ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
          #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
          #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

$$\xrightarrow[\text{zip}]{\text{lowering}}$$

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
           : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \mathbf{reduceSeq}(N, \delta_1, \delta_2, F, I, \boxed{E}) \quad = \quad C(\!(I)\!)_{\delta_2}(\lambda\ init. \\ \boxed{C(\!(E)\!)_{N.\delta_1}(\lambda\ x.} \\ out :=_{\delta_1} init; \\ \mathbf{for}\ N\ (\lambda i.\ out :=_{\delta_2}\ F(\boxed{x[i]}, out));\ ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>) {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
           #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
lambda{ ... }

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
         : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \textbf{reduceSeq}(N, \delta_1, \delta_2, \boxed{F}, I, E) \quad = \quad C(\!|I|\!)_{\delta_2} (\lambda\ init.$$
$$C(\!|E|\!)_{N.\delta_1} (\lambda\ x.$$
$$out :=_{\delta_1} init;$$
$$\textbf{for } N \boxed{(\lambda i.\ out :=_{\delta_2}\ F(x[i], out))}; ))$$

- **Start translation from** `rise.out`
- **Traverse program back to front and lower patterns as specified by their arguments**

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                    %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
           #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
lambda{ ... }

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \mathbf{reduceSeq}(N, \delta_1, \delta_2, \boxed{F}, I, E) \quad = \quad C(\![I]\!)_{\delta_2}(\lambda\ init.$$
$$C(\![E]\!)_{N.\delta_1}(\lambda\ x.$$
$$out :=_{\delta_1} init;$$
$$\mathbf{for}\ N\ \boxed{(\lambda i.\ out :=_{\delta_2}\ F(x[i], out))}; ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                     %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
           #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
lambda{ ... }

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \textbf{reduceSeq}(N, \delta_1, \delta_2, \boxed{F}, I, E) \quad = \quad \begin{aligned} &C(\!(I)\!)_{\delta_2}(\lambda\ init. \\ &C(\!(E)\!)_{N.\delta_1}(\lambda\ x. \\ &\quad out :=_{\delta_1} init; \\ &\quad \textbf{for } N\ \boxed{(\lambda i.\ out :=_{\delta_2}\ F(x[i], out))}; )) \end{aligned}$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>)  {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
          #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
          #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering
embed{ … }

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \textbf{reduceSeq}(N, \delta_1, \delta_2, \boxed{F}, I, E) \quad = \quad \begin{aligned} & C(\!|I|\!)_{\delta_2}(\lambda\ init. \\ & C(\!|E|\!)_{N.\delta_1}(\lambda\ x. \\ & \quad out :=_{\delta_1} init; \\ & \quad \textbf{for}\ N\ \boxed{(\lambda i.\ out :=_{\delta_2}\ F(x[i], out))}; )) \end{aligned}$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

IR trimmed for brevity

# Lowering to rise.codegen: dot

```
1  func @dot(%out:memref<1xf32>, %a:memref<1024xf32>,
                                  %b:memref<1024xf32>) {
2    %lhs = rise.in %a : !rise.array<1024, scalar<f32>>
3    %rhs = rise.in %b : !rise.array<1024, scalar<f32>>
4
5    %zip = rise.zip #rise.nat<1024> #rise.scalar<f32>
           #rise.scalar<f32>
6    %zipped = rise.apply %zip, %lhs, %rhs
7    %f = rise.lambda (%tuple, %acc) → !rise.scalar<f32> {
8      %fstFun = rise.fst #rise.scalar<f32> #rise.scalar<f32>
9      %sndFun = rise.snd #rise.scalar<f32> #rise.scalar<f32>
10     %fst = rise.apply %fstFun, %tuple
11     %snd = rise.apply %sndFun, %tuple
12     %result = rise.embed(%fst, %snd, %acc) {
13       %product = mulf %fst, %snd : f32
14       %res = addf %product, %acc : f32
15       return res : f32
16     } : !rise.scalar<f32>
17     rise.return %result : !rise.scalar<f32>
18   }
19   %init = rise.literal #rise.lit<0.0>
20   %reduce = rise.reduceSeq #rise.nat<1024>
           #rise.tuple<scalar<f32>, scalar<f32>> #rise.scalar<f32>
21   %result = rise.apply %reduce, %f, %init, %zipped
22   rise.out %out ← %result
23   return
24 }
```

RISE

lowering →

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
         : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

$$out :=_{\delta_2} \mathbf{reduceSeq}(N, \delta_1, \delta_2, F, I, E) \quad = \quad C(\![I]\!)_{\delta_2}(\lambda\ init.$$
$$C(\![E]\!)_{N.\delta_1}(\lambda\ x.$$
$$out :=_{\delta_1} init;$$
$$\mathbf{for}\ N\ (\lambda i.\ out :=_{\delta_2} F(x[i], out)); ))$$

- Start translation from `rise.out`
- Traverse program back to front and lower patterns as specified by their arguments

52

IR trimmed for brevity

# Lowering of the RISE dialect

Lowering of patterns in not context free, e.g. zip influences the code generation of the next pattern
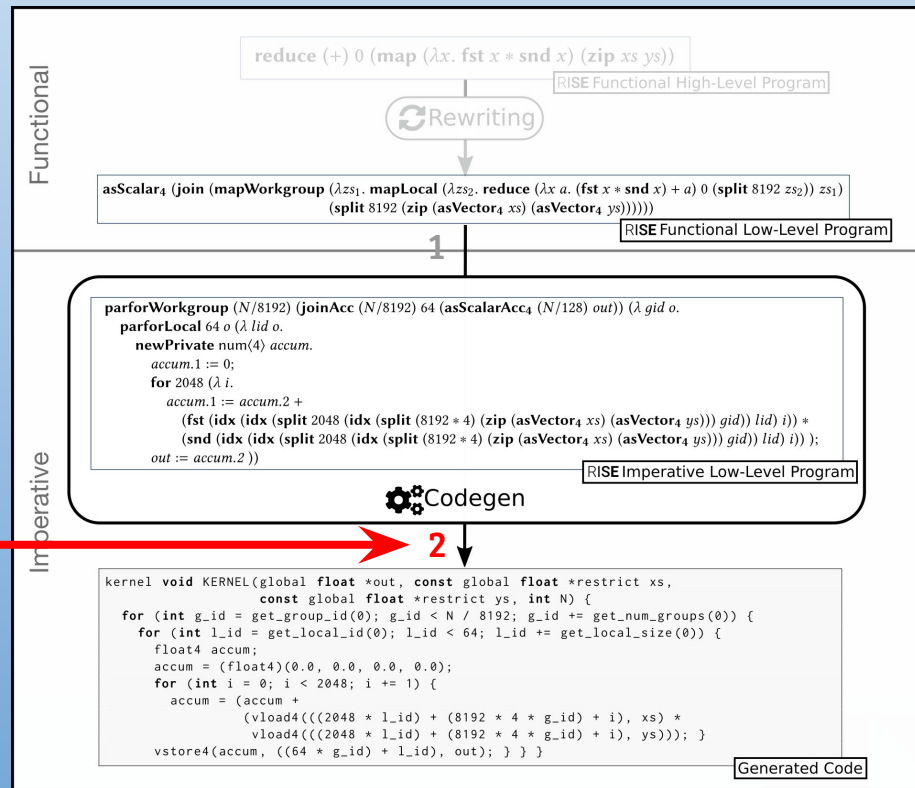
Lowering is divided into two steps:

1. Lowering functional to intermediate imperative

**2. Lowering to final target code**

Lowering described formally in:
Strategy Preserving Compilation for Parallel Functional Code
by Robert Atkey, Michel Steuwer, Sam Lindley, and Christophe Dubach
https://arxiv.org/abs/1710.08332

# Lowering rise.codegen to loop: Matrix Multiplication

```
func @mm(%out, %A, %B) {
   // map(map(dot))
}
```

lowering

```
1  func @mm(%out, %A, %B) {
      %cst0 = constant 0.000000e+00 : f32
2    %c0 = constant 0 : index
3    %c1024 = constant 1024 : index
4    %c1 = constant 1 : index
5    loop.for %i = %c0 to %c1024 step %c1 {
6      %0 = rise.codegen.idx(%A, %i)
7      %1 = rise.codegen.idx(%out, %i)
8      loop.for %j = %c0 to %c1024 step %c1 {
9        %2 = rise.codegen.idx(%B, %j)
10       %3 = rise.codegen.idx(%1, %j)
11       %4 = rise.codegen.zip(%0, %2)
12       %5 = rise.codegen.idx(%3, %c0)
13       rise.codegen.assign(%cst0, %5)
14       loop.for %k = %c0 to %c1024 step %c1 {
15         %6 = rise.codegen.idx(%4, %k)
16         %7 = rise.codegen.idx(%3, %c0)
17         %8 = rise.codegen.fst(%6)
18         %9 = rise.codegen.snd(%6)
19         %10 = mulf %8, %9 : f32
20         %11 = addf %10, %7 : f32
21         rise.codegen.assign(%11, %7)
22       }
23     }
24   }
25   return
26 }
```

Loop + RISE Intermediate

```
1  func @dot(%out:memref<1xf32>,%a:memref<1024xf32>,%b:memref<1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    %0 = rise.codegen.zip(%a, %b)
          : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
7    %1 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
8    rise.codegen.assign(%cst0, %1) : (f32, f32) → ()
9    loop.for %i = %c0 to %c1024 step %c1 {
10     %2 = rise.codegen.idx(%0, %i) : (memref<1024xf32>, index) → f32
11     %3 = rise.codegen.idx(%out, %c0) : (memref<1xf32>, index) → f32
12     %4 = rise.codegen.fst(%2) : (f32) → f32
13     %5 = rise.codegen.snd(%2) : (f32) → f32
14     %6 = mulf %4, %5 : f32
15     %7 = addf %6, %3 : f32
16     rise.codegen.assign(%7, %3) : (f32, f32) → ()
17   }
18   return
19 }
```

Loop + RISE Intermediate

- Matrix multiply builds on dot -> similar inner structure
- Only differences:
  - different **input** to `rise.codegen.zip`
  - accumulating **indexed** with j and i

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
     %cst0 = constant 0.000000e+00 : f32
2    %c0 = constant 0 : index
3    %c1024 = constant 1024 : index
4    %c1 = constant 1 : index
5    loop.for %i = %c0 to %c1024 step %c1 {
6      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
7      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      loop.for %j = %c0 to %c1024 step %c1 {
9        %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
10       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
11       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
12       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
13       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
14       loop.for %k = %c0 to %c1024 step %c1 {
15         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
16         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
17         %8 = rise.codegen.fst(%6) : (f32) → f32
18         %9 = rise.codegen.snd(%6) : (f32) → f32
19         %10 = mulf %8, %9 : f32
20         %11 = addf %10, %7 : f32
21         rise.codegen.assign(%11, %7) : (f32, f32) → ()
22       }
23     }
24   }
25   return
26 }
```

`Loop + RISE Intermediate`

IR trimmed for brevity

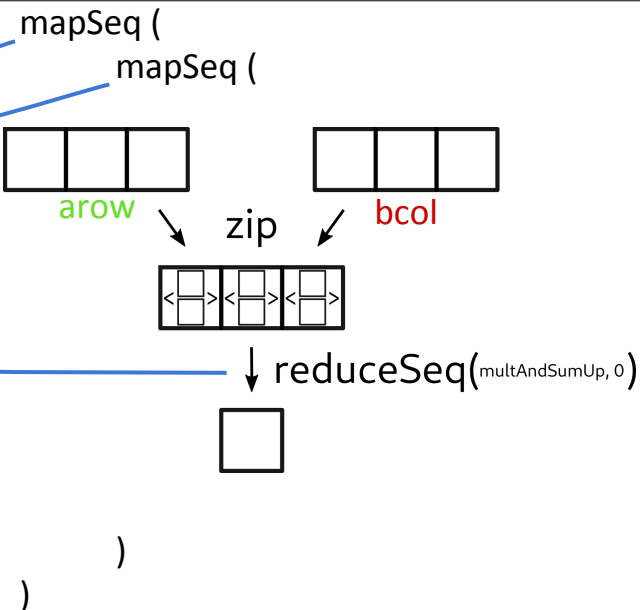# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
     %cst0 = constant 0.000000e+00 : f32
2    %c0 = constant 0 : index
3    %c1024 = constant 1024 : index
4    %c1 = constant 1 : index
5    loop.for %i = %c0 to %c1024 step %c1 {
6      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
7      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      loop.for %j = %c0 to %c1024 step %c1 {
9        %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
10       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
11       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
12       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
13       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
14       loop.for %k = %c0 to %c1024 step %c1 {
15         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
16         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
17         %8 = rise.codegen.fst(%6) : (f32) → f32
18         %9 = rise.codegen.snd(%6) : (f32) → f32
19         %10 = mulf %8, %9 : f32
20         %11 = addf %10, %7 : f32
21         rise.codegen.assign(%11, %7) : (f32, f32) → ()
22       }
23     }
24   }
25   return
26 }
```

Loop + RISE Intermediate

mapSeq (

    mapSeq (

arow    zip    bcol

reduceSeq(multAndSumUp, 0)

    )

)

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

? = ?

- Start codegen from `rise.codegen.assign`

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

? = ? + ?

- Start codegen from `rise.codegen.assign`
- Traverse program back to front following references

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

`Loop + RISE Intermediate`

`? = ? * ? + ?`

- Start codegen from `rise.codegen.assign`
- Traverse program back to front following references

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

**Path:**

fst

? = ? * ? + ?

- **Start codegen from** `rise.codegen.assign`
- **Traverse program back to front following references**
- **build up path**

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

Path:

```
idx(%k)
fst
```

```
? = ?[?] * ? + ?
```

- Start codegen from `rise.codegen.assign`
- Traverse program back to front following references
- build up path

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

Path:

```
idx(%k)
fst
```

codegen(lhs)

? = ?[?] * ? + ?

- Start codegen from `rise.codegen.assign`
- Traverse program back to front following references
- build up path

62

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

**Path:**

```
idx(%i)
idx(%k)
```

? = ?[?,?] * ? + ?

- **Start codegen from `rise.codegen.assign`**
- **Traverse program back to front following references**
- **build up path**

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

`Loop + RISE Intermediate`

**Path:**

```
idx(%i)
idx(%k)
```

```
? = A[?,?] * ? + ?
```

- Start codegen from `rise.codegen.assign`
- Traverse program back to front following references
- build up path

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

**Path:**

```
idx(%i)
idx(%k)
```

? = A[k,i] * ? + ?

- **Start codegen from `rise.codegen.assign`**
- **Traverse program back to front following references**
- **build up path**
- **reverse path**

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7      %0 = rise.codegen.idx(%A, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
8      %1 = rise.codegen.idx(%out, %i) : (memref<1024×1024xf32>, index) → memref<1024xf32>
9      loop.for %j = %c0 to %c1024 step %c1 {
10       %2 = rise.codegen.idx(%B, %j) : (memref<1024×1024xf32>, index) → memref<1024xf32>
11       %3 = rise.codegen.idx(%1, %j) : (memref<1024xf32>, index) → memref<1024xf32>
12       %4 = rise.codegen.zip(%0, %2) : (memref<1024xf32>, memref<1024xf32>) → memref<1024xf32>
13       %5 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
14       rise.codegen.assign(%cst0, %5) : (f32, f32) → ()
15       loop.for %k = %c0 to %c1024 step %c1 {
16         %6 = rise.codegen.idx(%4, %k) : (memref<1024xf32>, index) → f32
17         %7 = rise.codegen.idx(%3, %c0) : (memref<1024xf32>, index) → f32
18         %8 = rise.codegen.fst(%6) : (f32) → f32
19         %9 = rise.codegen.snd(%6) : (f32) → f32
20         %10 = mulf %8, %9 : f32
21         %11 = addf %10, %7 : f32
22         rise.codegen.assign(%11, %7) : (f32, f32) → ()
23       }
24     }
25   }
26   return
27 }
```

Loop + RISE Intermediate

generate()   %new8 = load %A[%k, %i]

**Path:**

idx(%i)
idx(%k)

? = A[k,i] * ? + ?

- **Start codegen from** `rise.codegen.assign`
- **Traverse program back to front following references**
- **build up path**
- **reverse path and generate new operations**

# Lowering rise.codegen to loop: Matrix Multiplication

```
1  func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>, %B: memref<1024×1024xf32>) {
2    %cst0 = constant 0.000000e+00 : f32
3    %c0 = constant 0 : index
4    %c1024 = constant 1024 : index
5    %c1 = constant 1 : index
6    loop.for %i = %c0 to %c1024 step %c1 {
7
8
9      loop.for %j = %c0 to %c1024 step %c1 {
10
11
12
13
14        store %cst0, %out[%j, %i] : memref<1024×1024xf32>
15        loop.for %k = %c0 to %c1024 step %c1 {
16          %0 = load %A[%k, %i] : memref<1024×1024xf32>
17          %1 = load %B[%k, %j] : memref<1024×1024xf32>
18          %2 = load %out[%j, %i] : memref<1024×1024xf32>
19
20          %3 = mulf %0, %1 : f32
21          %4 = addf %3, %2 : f32
22          store %4, %out[%j, %i] : memref<1024×1024xf32>
23        }
24      }
25    }
26    return
27 }
```

Loop

# Lowering RISE to affine

```
1    func @mm(%out:memref<1024×1024xf32>,%A:memref<1024×1024xf32>,
2             %B:memref<1024×1024xf32>) {
3      %cst0 = constant 0.000000e+00 : f32
4      %c0 = constant 0 : index
5
6
7      affine.for %i = 0 to 1024 {
8        affine.for %j = 0 to 1024 {
9          affine.store %cst0, %out[%j, %i] : memref<1024×1024xf32>
10         affine.for %k = 0 to 1024 {
11           %0 = affine.load %arg1[%k, %i] : memref<1024×1024xf32>
12           %1 = affine.load %arg2[%k, %j] : memref<1024×1024xf32>
13           %2 = affine.load %arg0[%j, %i] : memref<1024×1024xf32>
14           %3 = mulf %0, %1 : f32
15           %4 = addf %3, %2 : f32
16           affine.store %4, %arg0[%j, %i] : memref<1024×1024xf32>
17         }
18       }
19     }
20     return
21   }
```

Affine

- Similar to the loop lowering presented before

- Lowering target specified using Attributes
    - rise.reduceSeq {to = "affine"}

- We can always use affine loops with our patterns

- Enables usage of polyhedral optimizations

# Lowering RISE to affine vs. loop

```
1
2   func @mm(%out:memref<1024×1024xf32>,%A:memref<1024×1024xf32>,
3       %B:memref<1024×1024xf32>) {
4     %cst0 = constant 0.000000e+00 : f32
5     %c0 = constant 0 : index
6
7
8     affine.for %i = 0 to 1024 {
9       affine.for %j = 0 to 1024 {
10        affine.store %cst0, %out[%j, %i] : memref<1024×1024xf32>
11        affine.for %k = 0 to 1024 {
12          %0 = affine.load %arg1[%k, %i] : memref<1024×1024xf32>
13          %1 = affine.load %arg2[%k, %j] : memref<1024×1024xf32>
14          %2 = affine.load %arg0[%j, %i] : memref<1024×1024xf32>
15          %3 = mulf %0, %1 : f32
16          %4 = addf %3, %2 : f32
17          affine.store %4, %arg0[%j, %i] : memref<1024×1024xf32>
18        }
19      }
20    }
21    return
   }
```

Affine
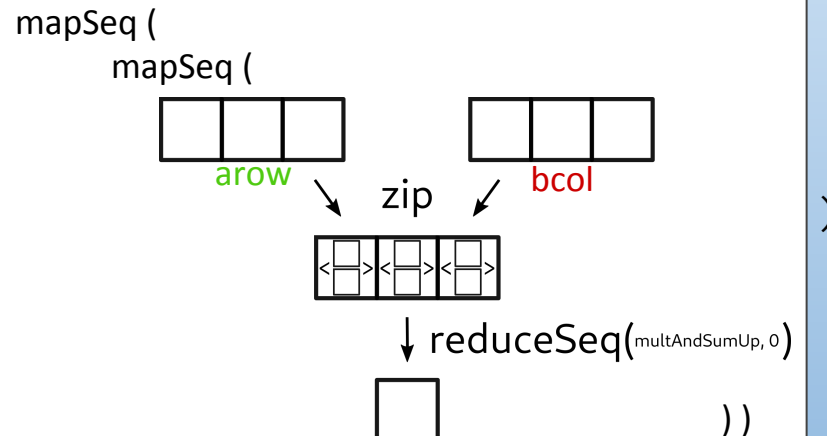
```
1   func @mm(%out: memref<1024×1024xf32>, %A: memref<1024×1024xf32>,
2       %B: memref<1024×1024xf32>) {
3     %cst0 = constant 0.000000e+00 : f32
4     %c0 = constant 0 : index
5     %c1024 = constant 1024 : index
6     %c1 = constant 1 : index
7     loop.for %i = %c0 to %c1024 step %c1 {
8       loop.for %j = %c0 to %c1024 step %c1 {
9         store %cst0, %out[%j, %i] : memref<1024×1024xf32>
10        loop.for %k = %c0 to %c1024 step %c1 {
11          %0 = load %A[%k, %i] : memref<1024×1024xf32>
12          %1 = load %B[%k, %j] : memref<1024×1024xf32>
13          %2 = load %out[%j, %i] : memref<1024×1024xf32>
14          %3 = mulf %0, %1 : f32
15          %4 = addf %3, %2 : f32
16          store %4, %out[%j, %i] : memref<1024×1024xf32>
17        }
18      }
19    }
20    return
21 }
```

Loop

# Lowering RISE to library calls

# Lowering RISE to library calls



match_for(

matchSuccess

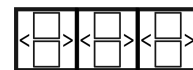mapSeq (
  mapSeq (

  arow    zip    bcol

  reduceSeq(multAndSumUp, 0)
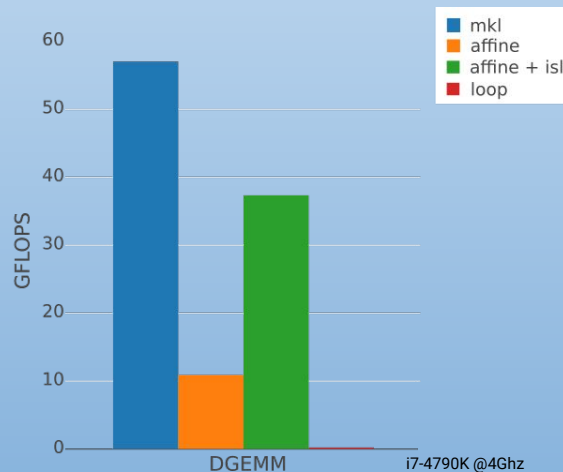
))

```
1   func @mm(%out:memref<1024×1024xf32>,%A:memref<1024×1024xf32>,%B:memref<1024×1024xf32>) {
2     %CblasRowMajor = constant 101 : i32
3     %CblasNoTrans = constant 111 : i32
4     %M = constant 1024 : i32
5     %N = constant 1024 : i32
6     %K = constant 1024 : i32
7     %LDA = constant 1024 : i32
8     %LDB = constant 1024 : i32
9     %LDC = constant 1024 : i32
10    %alpha = constant 1.0 : f32
11    %beta = constant 1.0 : f32
12    call @cblas_sgemm_wrapper(%CblasRowMajor, %CblasNoTrans, %CblasNoTrans, %M, %N, %K, %alpha, %A, %LDA, %B, %LDB, %beta, %out, %LDC) :
            (i32, i32, i32, i32, i32, i32, f32, memref<1024×1024xf32>, i32, memref<1024×1024xf32>, i32, f32, memref<1024×1024xf32>, i32) → ()
13    return
14  }
```

# State of our implementation

- Lambda calculus foundations implemented

- Basic set of patterns implemented

- Basic lowering system implemented

- Lowering to Library Calls quite ad-hoc

- Preliminary experiments shows
  performance as expected

# RISE in MLIR: Summary

**Today**

- Lambda-Calculus + Composable Patterns in MLIR
- Ability to embed arbitrary code from other dialects
- Different lowering approaches

**Future directions:**

- Explore other lowering approaches (e.g. GPU, OpenMP dialects)
- Introduce symbolic sizes via dependent types
- Rewriting enables composable and reusable high-level transformations
    - New language for controlling rewriting (ELEVATE - https://elevate-lang.org)

# RISE in MLIR

## A functional Pattern-based Dialect

# We are Open Source!

https://rise-lang.org/mlir

https://github.com/rise-lang/mlir

**Martin Lücke** | Michel Steuwer | Aaron Smith

THE UNIVERSITY
*of* EDINBURGH    |    University *of* Glasgow

VIA VERITAS VITA