# RISE in MLIR

A functional Pattern-based Dialect

Martin Lücke | Michel Steuwer | Aaron Smith

THE UNIVERSITY *of* EDINBURGH
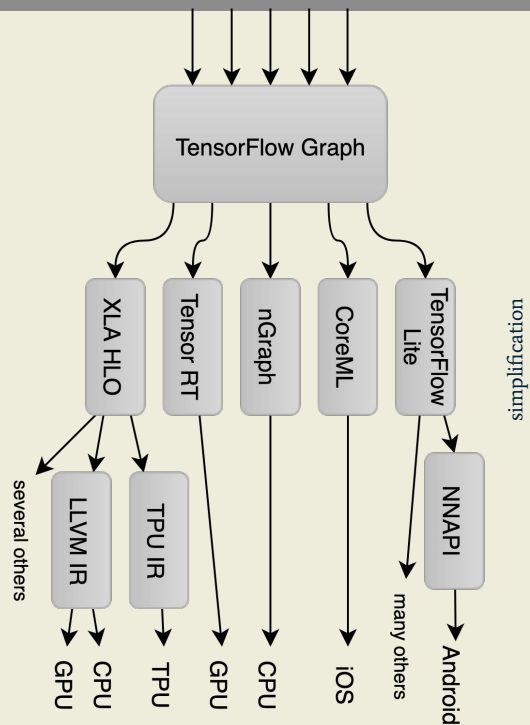
University *of* Glasgow

VIA VERITAS VITA

# Machine Learning Systems are stuck in a Rut

Paul Barham and Michael Isard [HotOS2019]

- Currently much focus on optimizing 5 year old ML benchmarks
- New ideas in ML often require new primitives that are hard to compile and optimize for the zoo of modern hardware
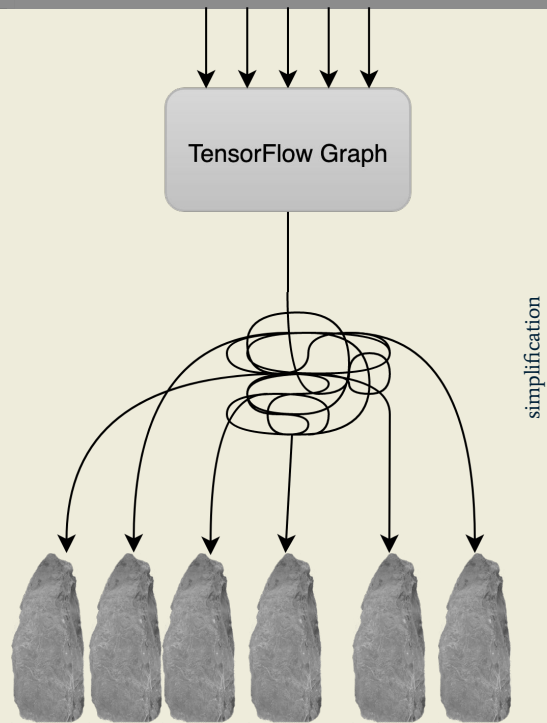
# Machine Learning Systems are stuck in a Rut

Paul Barham and Michael Isard [HotOS2019]

- Currently much focus on optimizing 5 year old ML benchmarks
- New ideas in ML often require new primitives that are hard to compile and optimize for the zoo of modern hardware

**Key Problem:**

- Reliance on high-performance but **inflexible monolithic kernels**
- The resulting lack of expressiveness, maintainability, and modularity hinders research progress
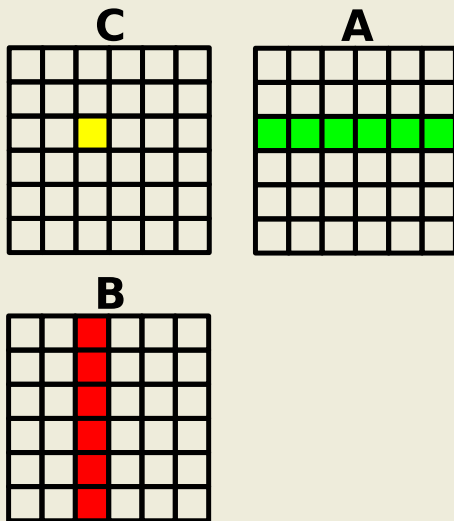
TensorFlow Graph

simplification

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

- To break up monolithic kernels we argue for:

  **representing computations using compositions of flexible and generic patterns**

- RISE is a spiritual successor to the Lift project

- Pattern-based style efficiently represents complex multidimensional computations of various domains

- Optimization choices encoded as rewrite rules are explored automatically for competitive performance

**We believe this style will allows us to express optimizations on a higher level**

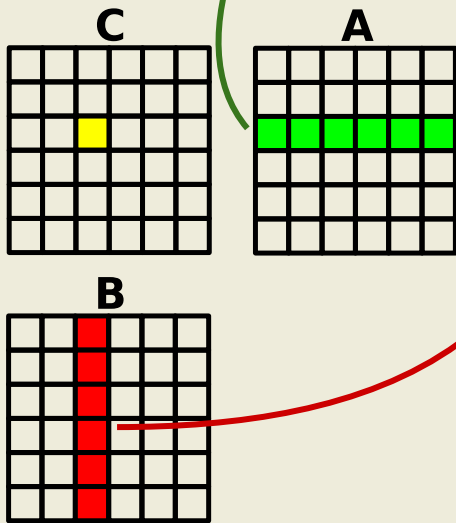**than what is currently available in the existing ML frameworks**

# 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) )) )) ))
```
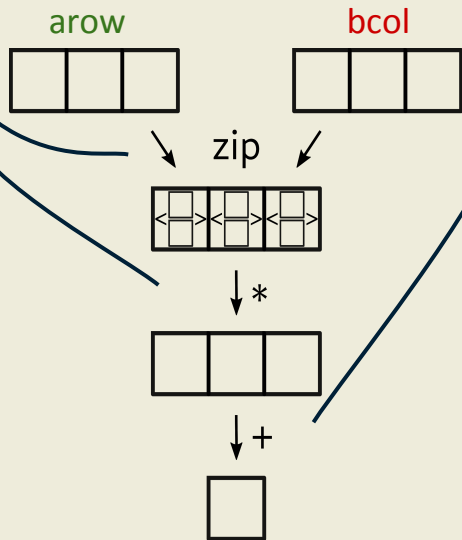


**C** **A**

**B**

dotproduct 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) )) )) ))
```
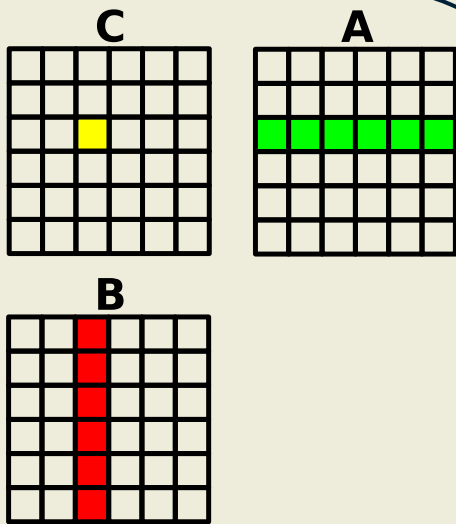
# RISE - The Caveats

Does everyone have to write functional programs now?

Academic work written in Scala

Does not integrate well with existing ML 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                                     - Wrapping 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                      - Wrapping LLVM IR

# 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                        - Wrapping LLVM IR

# RISE in `MLIR` - **Lambda Calculus 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 $\rightarrow$ usable for industry

- Natural integration with machine learning toolchains

   $\rightarrow$  to implement RISE we implement $\lambda$-**calculus** as an `MLIR` dialect

# RISE dialect by example: Matrix multiplication

```
1   func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2            (!rise.data<array<3, array<3, int>>>) {
3
4     %f1 = rise.lambda (%arow) :
5              !rise.fun<data<array<3, int>> → data<array<3, int>>> {
6         %f2 = rise.lambda (%bcol) :
7              !rise.fun<data<array<3, int>> → data<array<3, int>>> {
8            %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9            %zipped = rise.apply %zip, %arow, %bcol
10
11           %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12              %fst = rise.fst #rise.int #rise.int
13              %snd = rise.snd #rise.int #rise.int
14              %t1  = rise.apply %fst, %t
15              %t2  = rise.apply %snd, %t
16              %mul = rise.mult #rise.int
17              %res = rise.apply %mul, %t1, %t2
18              rise.return %res
19           }
20           %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21           %mapped = rise.apply %map, %f, %zipped
22
23           %add = rise.add #rise.int
24           %init = rise.literal #rise.lit<int<0>>
25           %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26           %res = rise.apply %reduce, %add, %init, %mapped
27           rise.return %res
28        }
29
30        %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31        %res = rise.apply %map, %f2, %B
32        rise.return %res
33     }
34
35     %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36     %res = rise.apply %map, %f1, %A
37     return %res
38  }
```

# Types of the RISE **dialect**

```
1    func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2                     (!rise.data<array<3, array<3, int>>>) {
3
4      %f1 = rise.lambda (%arow) :
5                     !rise.fun<data<array<3, int>> → data<array<3, int>>> {
6          %f2 = rise.lambda (%bcol) :
7                       !rise.fun<data<array<3, int>> → data<array<3, int>>> {
8              %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9              %zipped = rise.apply %zip, %arow, %bcol
10
11             %f = rise.lambda (%t) :   !rise.fun<data<tuple<int, int>> → data<int>> {
12                 %fst = rise.fst #rise.int #rise.int
13                 %snd = rise.snd #rise.int #rise.int
14                 %t1  = rise.apply %fst, %t
15                 %t2  = rise.apply %snd, %t
16                 %mul = rise.mult #rise.int
17                 %res = rise.apply %mul, %t1, %t2
18                 rise.return %res
19             }
20             %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21             %mapped = rise.apply %map, %f, %zipped
22
23             %add = rise.add #rise.int
24             %init = rise.literal #rise.lit<int<0>>
25             %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26             %res = rise.apply %reduce, %add, %init, %mapped
27             rise.return %res
28         }
29
30         %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31         %res = rise.apply %map, %f2, %B
32         rise.return %res
33     }
34
35     %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36     %res = rise.apply %map, %f1, %A
37     return %res
38 }
```

**DataTypes**

**FunctionTypes**

# Types of the RISE **dialect**

## DataTypes

```
1    func @mm(%A: !rise.data<array<4, array<4, int>>>,
2             %B: !rise.data<array<4, array<4, int>>>) →
3               (!rise.data<array<4, array<4, int>>>) {
```

- Rise DataTypes: array types, tuple types, scalar types
- Nested array types represent higher dimensional data

## FunctionTypes

```
11   %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>>
```

- Rise function types: types of lambda expressions
- Type system prevents mixing of function and data types

# Patterns in the RISE **dialect**

```
1   func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2              (!rise.data<array<3, array<3, int>>>) {
3
4     %f1 = rise.lambda (%arow) :
5                !rise.fun<data<array<3, int>> → data<array<3, int>>> {
6         %f2 = rise.lambda (%bcol) :
7                !rise.fun<data<array<3, int>> → data<array<3, int>>> {
8             %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9             %zipped = rise.apply %zip, %arow, %bcol
10
11            %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12                %fst = rise.fst #rise.int #rise.int
13                %snd = rise.snd #rise.int #rise.int
14                %t1  = rise.apply %fst, %t
15                %t2  = rise.apply %snd, %t
16                %mul = rise.mult #rise.int
17                %res = rise.apply %mul, %t1, %t2
18                rise.return %res
19            }
20            %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21            %mapped = rise.apply %map, %f, %zipped
22
23            %add = rise.add #rise.int
24            %init = rise.literal #rise.lit<int<0>>
25            %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26            %res = rise.apply %reduce, %add, %init, %mapped
27            rise.return %res
28        }
29
30        %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31        %res = rise.apply %map, %f2, %B
32        rise.return %res
33    }
34
35    %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36    %res = rise.apply %map, %f1, %A
37    return %res
38 }
```

# Patterns in the RISE **dialect**

```
1   func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2              (!rise.data<array<3, array<3, int>>>) {
3
4     %f1 = rise.lambda (%arow) :
5                !rise.fun<data<array<3, int>> → data<array<3, int>> {
6         %f2 = rise.lambda (%bcol) :
7                !rise.fun<data<array<3, int>> → data<array<3, int>> {
8             %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9             %zipped = rise.apply %zip, %arow, %bcol
10
11            %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12                %fst = rise.fst #rise.int #rise.int
13                %snd = rise.snd #rise.int #rise.int
14                %t1  = rise.apply %fst, %t
15                %t2  = rise.apply %snd, %t
16                %mul = rise.mult #rise.int
17                %res = rise.apply %mul, %t1, %t2
18                rise.return %res
19            }
20            %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21            %mapped = rise.apply %map, %f, %zipped
22
23            %add = rise.add #rise.int
24            %init = rise.literal #rise.lit<int<0>>
25            %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26            %res = rise.apply %reduce, %add, %init, %mapped
27            rise.return %res
28        }
29
30        %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31        %res = rise.apply %map, %f2, %B
32        rise.return %res
33    }
34
35    %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36    %res = rise.apply %map, %f1, %A
37    return %res
38 }
```

# Patterns in the RISE dialect: zip

```
8    %zip = rise.zip #rise.nat<3> #rise.int #rise.int
```

# Patterns in the RISE **dialect:** `zip`

```
8   %zip = rise.zip #rise.nat<3> #rise.int #rise.int

    // type: () → !rise.fun<data<array<3, int>> →
                          fun<data<array<3, int>> →
                              data<array<3, tuple<int, int>>>>>

    // type: () →
```



- One `MLIR` operation per RISE pattern

- Operations customized with attributes specifying information for the type

- Patterns encoded as operations have a RISE function type

# Function application in the RISE dialect

```
4
5
6
7
8        %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9
```

# Function application in the RISE dialect

```
4  %f1 = rise.lambda (%arow) : !rise.fun<data<array<3, int>> →
5                                        data<array<3, int>> {
6    %f2 = rise.lambda (%bcol) : !rise.fun<data<array<3, int>> →
7                                        data<array<3, int>> {
8        %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9        %zipped = rise.apply %zip, %arow, %bcol
```

- `rise.apply` models function applications

- SSA value with RISE `function type` and arguments to the function are passed to `rise.apply`

- Partial function application naturally supported (i.e. not specifying all function arguments at once)

# Function application in the RISE **dialect**

```
4  %f1 = rise.lambda (%arow) : !rise.fun<data<array<3, int>> →
5                                        data<array<3, int>>> {
6      %f2 = rise.lambda (%bcol) : !rise.fun<data<array<3, int>> →
7                                          data<array<3, int>>> {
8          %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9          %zipped = rise.apply %zip, %arow, %bcol

 //type: (λ, data<array<3, int>>, data<array<3, int>>) →
                                    data<array<3,tuple<int,int>>>
```

- `rise.apply` models function applications

- SSA value with RISE `function type` and arguments to the function are passed to `rise.apply`

- Partial function application naturally supported (i.e. not specifying all function arguments at once)

# Function abstraction in the RISE **dialect**

```
1   func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2            (!rise.data<array<3, array<3, int>>>) {
3
4     %f1 = rise.lambda (%arow) :
5              !rise.fun<data<array<3, int>> → data<array<3, int>>> {
6        %f2 = rise.lambda (%bcol) :
7               !rise.fun<data<array<3, int>> → data<array<3, int>>> {
8           %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9           %zipped = rise.apply %zip, %arow, %bcol
10
11          %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12             %fst = rise.fst #rise.int #rise.int
13             %snd = rise.snd #rise.int #rise.int
14             %t1  = rise.apply %fst, %t
15             %t2  = rise.apply %snd, %t
16             %mul = rise.mult #rise.int
17             %res = rise.apply %mul, %t1, %t2
18             rise.return %res
19          }
20          %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21          %mapped = rise.apply %map, %f, %zipped
22
23          %add = rise.add #rise.int
24          %init = rise.literal #rise.lit<int<0>>
25          %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26          %res = rise.apply %reduce, %add, %init, %mapped
27          rise.return %res
28       }
29
30       %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31       %res = rise.apply %map, %f2, %B
32       rise.return %res
33    }
34
35    %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36    %res = rise.apply %map, %f1, %A
37    return %res
38 }
```

# Function abstraction in the RISE **dialect**

```
1   func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2            (!rise.data<array<3, array<3, int>>>) {
3
4     %f1 = rise.lambda (%arow) :
5              !rise.fun<data<array<3, int>> → data<array<3, int>> {
6         %f2 = rise.lambda (%bcol) :
7                  !rise.fun<data<array<3, int>> → data<array<3, int>> {
8             %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9             %zipped = rise.apply %zip, %arow, %bcol
10
11            %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12                %fst = rise.fst #rise.int #rise.int
13                %snd = rise.snd #rise.int #rise.int
14                %t1  = rise.apply %fst, %t
15                %t2  = rise.apply %snd, %t
16                %mul = rise.mult #rise.int
17                %res = rise.apply %mul, %t1, %t2
18                rise.return %res
19            }
20            %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21            %mapped = rise.apply %map, %f, %zipped
22
23            %add = rise.add #rise.int
24            %init = rise.literal #rise.lit<int<0>>
25            %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26            %res = rise.apply %reduce, %add, %init, %mapped
27            rise.return %res
28        }
29
30        %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31        %res = rise.apply %map, %f2, %B
32        rise.return %res
33    }
34
35    %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36    %res = rise.apply %map, %f1, %A
37    return %res
38 }
```

# Function abstraction in the RISE dialect

```
11  %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12       %fst = rise.fst #rise.int #rise.int
13       %snd = rise.snd #rise.int #rise.int
14       %t1  = rise.apply %fst, %t
15       %t2  = rise.apply %snd, %t
16       %mul = rise.mult #rise.int
17       %res = rise.apply %mul, %t1, %t2
18       rise.return %res
19  }
```
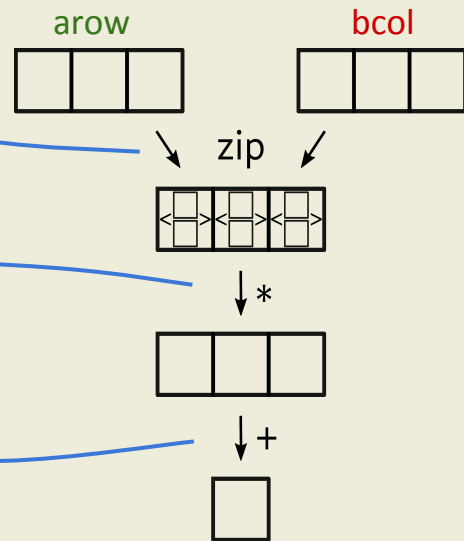
- rise.lambda wraps an MLIR region of exactly one block

- Arbitrary number of arguments and one result

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

# RISE dialect by example: Matrix multiply

```
1   func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2            (!rise.data<array<3, array<3, int>>>) {
3
4     %f1 = rise.lambda (%arow) :
5              !rise.fun<data<array<3, int>> → data<array<3, int>>> {
6         %f2 = rise.lambda (%bcol) :
7              !rise.fun<data<array<3, int>> → data<array<3, int>>> {
8            %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9            %zipped = rise.apply %zip, %arow, %bcol
10
11           %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12               %fst = rise.fst #rise.int #rise.int
13               %snd = rise.snd #rise.int #rise.int
14               %t1  = rise.apply %fst, %t
15               %t2  = rise.apply %snd, %t
16               %mul = rise.mult #rise.int
17               %res = rise.apply %mul, %t1, %t2
18               rise.return %res
19           }
20           %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21           %mapped = rise.apply %map, %f, %zipped
22
23           %add = rise.add #rise.int
24           %init = rise.literal #rise.lit<int<0>>
25           %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26           %res = rise.apply %reduce, %add, %init, %mapped
27           rise.return %res
28         }
29
30         %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31         %res = rise.apply %map, %f2, %B
32         rise.return %res
33     }
34
35     %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36     %res = rise.apply %map, %f1, %A
37     return %res
38  }
```

# RISE dialect by example: Matrix multiply

```
1    func @mm(%A: !rise.data<array<3, array<3, int>>>, %B: !rise.data<array<3, array<3, int>>>) →
2              (!rise.data<array<3, array<3, int>>>) {
3
4      %f1 = rise.lambda (%arow) :
5                  !rise.fun<data<array<3, int>> → data<array<3, int>>> {
6          %f2 = rise.lambda (%bcol) :
7                  !rise.fun<data<array<3, int>> → data<array<3, int>>> {
8              %zip = rise.zip #rise.nat<3> #rise.int #rise.int
9              %zipped = rise.apply %zip, %arow, %bcol
10
11             %f = rise.lambda (%t) : !rise.fun<data<tuple<int, int>> → data<int>> {
12                 %fst = rise.fst #rise.int #rise.int
13                 %snd = rise.snd #rise.int #rise.int
14                 %t1  = rise.apply %fst, %t
15                 %t2  = rise.apply %snd, %t
16                 %mul = rise.mult #rise.int
17                 %res = rise.apply %mul, %t1, %t2
18                 rise.return %res
19             }
20             %map = rise.map #rise.nat<3> #rise.tuple<int, int> #rise.int
21             %mapped = rise.apply %map, %f, %zipped
22
23             %add = rise.add #rise.int
24             %init = rise.literal #rise.lit<int<0>>
25             %reduce = rise.reduce #rise.nat<3> #rise.int #rise.int
26             %res = rise.apply %reduce, %add, %init, %mapped
27             rise.return %res
28         }
29
30         %map = rise.map #rise.nat<3> #rise.array<3, int>  #rise.array<3, int>
31         %res = rise.apply %map, %f2, %B
32         rise.return %res
33     }
34
35     %map = rise.map #rise.nat<3> #rise.array<3,int> #rise.array<3,int>
36     %res = rise.apply %map, %f1, %A
37     return %res
38  }
```

# RISE in MLIR next steps

- Explore interaction and integration with other MLIR dialects:

    - Lowering RISE to LLVM IR or other dialects (e.g. polyhedral)

    - Lowering from the TensorFlow dialect to RISE

    - Raise RISE to TensorFlow or other domain-specific dialects?

- Express optimisations as rewrite rules over the RISE dialect

# 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

Fork me on GitHub

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

VIA VERITAS VITA