



7 JUNE 2022

RISE & SHINE

LANGUAGE-ORIENTED COMPILER DESIGN

Michel Steuwer, Thomas Koehler, Bastian Köpcke, Federico Pizzuti

with contributions from **Martin Lücke, Johannes Lenfers, Rongxiao Fu, Xueying Qin**

Compiler Design With MLIR

Representation-Oriented Compiler Design

- MLIR focuses on **Representation** (i.e. Syntax)
 - *Operations* have *operands*, *attributes*, and a *type*
 - *Operations* are organised in SSA form in *Blocks*
 - *Blocks* form CFGs and are part of *Regions* (*which can be nested in Operations*)
- Consistency is ensured by implementing and calling `verify` methods on operations, blocks, and regions ...
- **Semantics** is given informally by how representations are transformed in each other

How should we design *good* IRs that:

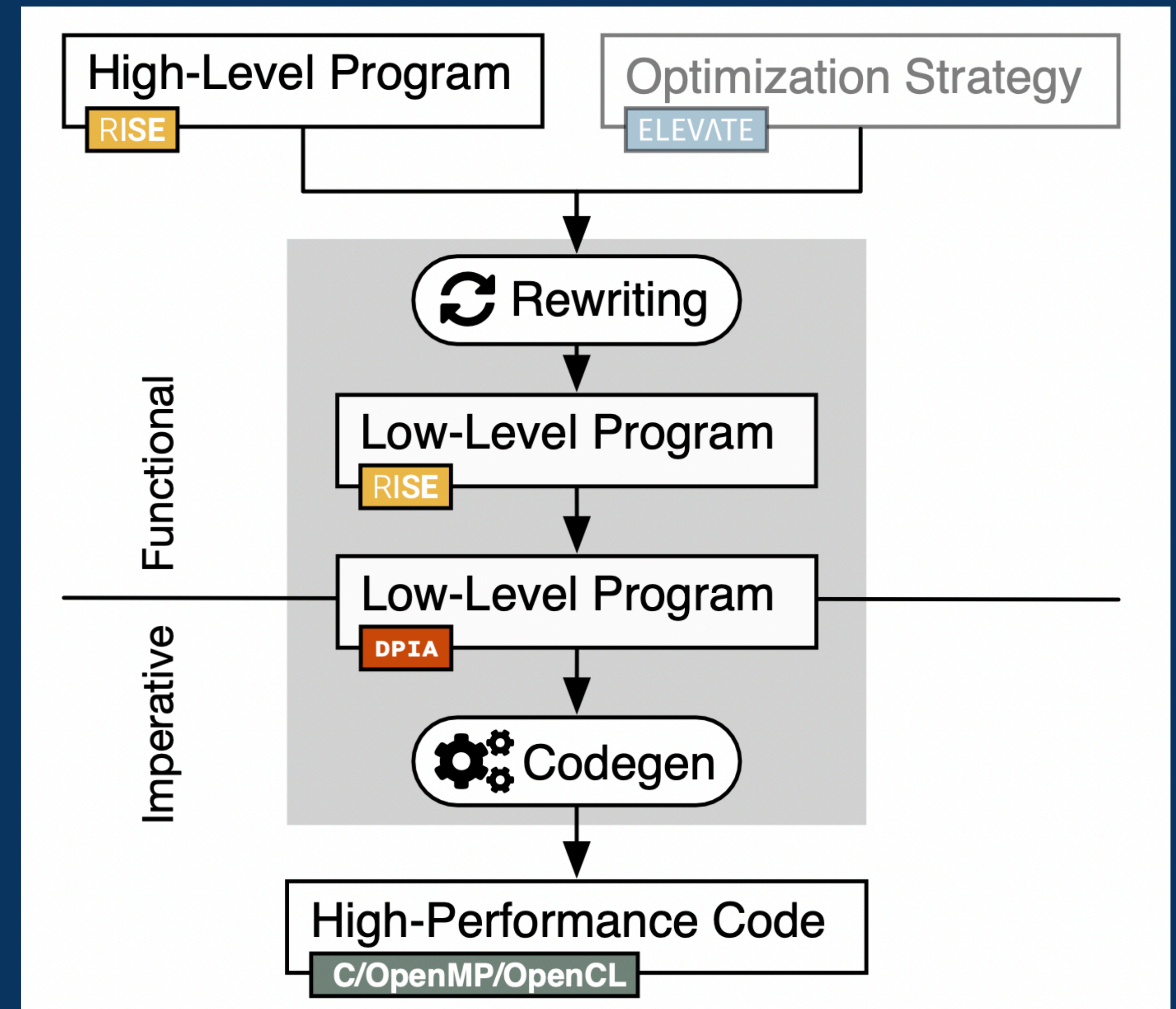
- 1) have a clear purpose and semantics;**
- 2) have formally checked invariants and assumptions;**
- 3) are easily extensible?**

Shine Compiler

Language-Oriented Compiler Design

We advocate for:

- Clear separation of concerns between *optimizing* and *code generation*
- Formalisation of invariants and assumptions about IRs in *type systems*
- *Extensibility* at each level in the compiler



GEMM in RISE

High-Level GEMM

```
1 depFun((m:Nat, n:Nat, k:Nat) =>
2   fun((A: Array[m, Array[k, f32]], B: Array[k, Array[n, f32]],
3     C: Array[m, Array[n, f32]], alpha: f32, beta: f32) =>
4     zip(A)(C) |> map(fun(rowAC =>
5       zip(B |> transpose)(snd(rowAC)) |> map(fun(colBC =>
6         zip(fst(rowAC))(fst(colBC)) |>
7         map(fun((a, b) => a * b)) |> reduce(+, 0) |>
8         fun(r => (alpha * r) + (beta * snd(colBC))) )))))))
```

Optimization Strategy

ELEVATE

Rewriting

Low-Level GEMM

```
9 depFun((m:Nat, n:Nat, k:Nat) => fun(A, B, C, alpha, beta =>
10   zip(A)(C) |> mapBlock(fun(rowAC =>
11     zip(B |> transpose)(snd(rowAC)) |>
12     mapThreads(fun(colBC => zip(fst(rowAC))(fst(colBC)) |>
13       reduceSeq(Local)(fun((acc, ab) =>
14         acc + fst(ab) * snd(ab)), 0) |>
15       fun(r => (alpha * r) + (beta * snd(colBC))) )))))))
```

Translation

Imperative GEMM

```
17 depFun((m:Nat, n:Nat, k:Nat) => fun(A, B, C, alpha, beta =>
18   parForBlock(m, Array[n, f16], output, fun(rowIdx, outRow =>
19     parForThreads(n, f16, outRow, fun(colIdx, outElem =>
20       new(Local, f32, fun((accumExp, accumAcc) =>
21         accumAcc = 0.0f;
22         for(k, fun(i => accumAcc = accumExp +
23           fst(idx(i, zip(fst(idx(rowIdx, zip(A, C))),
24             fst(idx(colIdx, zip(transpose(B),
25               snd(idx(rowIdx, zip(A, C)))))))) *
26             snd(idx(i, zip(fst(idx(rowIdx, zip(A, C))),
27               fst(idx(colIdx, zip(transpose(B),
28                 snd(idx(rowIdx, zip(A, C))))))))));
29         outElem = alpha * accumExp + beta *
30           snd(idx(colIdx, zip(transpose(B),
31             snd(idx(rowIdx, zip(A, C)))))) )));
32       syncThreads()))))
```

Codegen

```
33 __global__ void gemm_kernel(float* __restrict__ output,
34   int m, int n, int k, const __half* __restrict__ A,
35   const __half* __restrict__ B,
36   const float* __restrict__ C, float alpha, float beta) {
37   for(int rowIdx=blockIdx.x;
38     blockIdx.x < m; rowIdx += blockDim.x) {
39     for(int colIdx=threadIdx.x;
40       threadIdx.x < n; rowIdx += blockDim.x) {
41       float accum = 0;
42       for (int i = 0; i < k; i++) {
43         accum = accum + A[i + rowIdx*k] * B[colIdx + i*n];
44       }
45       output[colIdx + rowIdx * n] =
46         alpha * accum + beta * C[colIdx + rowIdx*n];
47     }
48     __syncthreads(); }
```

GEMM in RISE

RISE

High-Level GEMM

```
1 depFun((m:Nat,n:Nat,k:Nat) =>
2   fun((A: Array[m,Array[k,f32]], B: Array[k,Array[n,f32]],
3     C: Array[m,Array[n,f32]], alpha: f32, beta: f32) =>
4     zip(A)(C) |> map(fun(rowAC =>
5       zip(B |> transpose)(snd(rowAC)) |> map(fun(colBC =>
6         zip(fst(rowAC))(fst(colBC)) |>
7         map(fun((a, b) => a * b)) |> reduce(+, 0) |>
8         fun(r => (alpha * r) + (beta * snd(colBC))) ))))))
```

Optimization Strategy

ELEVATE

Rewriting

Low-Level GEMM

```
9 depFun((m:Nat,n:Nat,k:Nat) => fun(A,B,C,alpha,beta =>
10   zip(A)(C) |> mapBlock(fun(rowAC =>
11     zip(B |> transpose)(snd(rowAC)) |>
12     mapThreads(fun(colBC => zip(fst(rowAC))(fst(colBC)) |>
13       reduceSeq(Local)(fun((acc,ab) =>
14         acc + fst(ab) * snd(ab)),0) |>
15       fun(r => (alpha * r) + (beta * snd(colBC))) ))))))
```

Translation

Imperative GEMM

DPIA

```
17 depFun((m:Nat,n:Nat,k:Nat) => fun(A,B,C,alpha,beta
18   parForBlock(m,Array[n,f16], output, fun(rowIdx,outRow =>
19     parForThreads(n,f16, outRow, fun(colIdx,outElem =>
20       new(Local,f32, fun((accumExp, accumAcc) =>
21         accumAcc = 0.0f;
22         for(k, fun(i => accumAcc = accumExp +
23           fst(idx(i, zip(fst(idx(rowIdx, zip(A,C))),
24             fst(idx(colIdx, zip(transpose(B),
25               snd(idx(rowIdx, zip(A,C)))))))) *
26             snd(idx(i, zip(fst(idx(rowIdx, zip(A,C))),
27               fst(idx(colIdx, zip(transpose(B),
28                 snd(idx(rowIdx, zip(A,C))))))))));
29         outElem = alpha * accumExp + beta *
30           snd(idx(colIdx, zip(transpose(B),
31             snd(idx(rowIdx, zip(A,C)))))) ));
32         syncThreads()))))
```

Codegen

```
33 __global__ void gemm_kernel(float* __restrict__ output,
34   int m, int n, int k, const __half* __restrict__ A,
35   const __half* __restrict__ B,
36   const float* __restrict__ C, float alpha, float beta) {
37   for(int rowIdx=blockIdx.x;
38     blockIdx.x<m; rowIdx += gridDim.x) {
39     for(int colIdx=threadIdx.x;
40       threadIdx.x<n; rowIdx += blockDim.x) {
41       float accum = 0;
42       for (int i = 0; i < k; i++) {
43         accum = accum + A[i + rowIdx*k] * B[colIdx + i*n];
44       }
45       output[colIdx + rowIdx * n] =
46         alpha * accum + beta * C[colIdx + rowIdx*n];
47     }
48     __syncthreads(); }
```

C

GEMM in RISE

RISE

High-Level GEMM

```
1 depFun((m:Nat,n:Nat,k:Nat) =>
2   fun((A: Array[m,Array[k,f32]], B: Array[k,Array[n,f32]],
3     C: Array[m,Array[n,f32]], alpha: f32, beta: f32) =>
4     zip(A)(C) |> map(fun(rowAC =>
5       zip(B |> transpose)(snd(rowAC)) |> map(fun(colBC =>
6         zip(fst(rowAC))(fst(colBC)) |>
7         map(fun((a, b) => a * b)) |> reduce(+, 0) |>
8         fun(r => (alpha * r) + (beta * snd(colBC))) ))))))
```

Optimization Strategy

ELEVATE



Rewriting

Optimization

Low-Level GEMM

```
9 depFun((m:Nat,n:Nat,k:Nat) => fun(A,B,C,alpha,beta =>
10   zip(A)(C) |> mapBlock(fun(rowAC =>
11     zip(B |> transpose)(snd(rowAC)) |>
12     mapThreads(fun(colBC => zip(fst(rowAC))(fst(colBC)) |>
13       reduceSeq(Local)(fun((acc,ab) =>
14         acc + fst(ab) * snd(ab)),0) |>
15       fun(r => (alpha * r) + (beta * snd(colBC))) ))))))
```



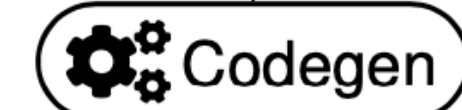
Translation

Translation

DPIA

Imperative GEMM

```
17 depFun((m:Nat,n:Nat,k:Nat) => fun(A,B,C,alpha,beta
18   parForBlock(m,Array[n,f16], output, fun(rowIdx,outRow =>
19     parForThreads(n,f16, outRow, fun(colIdx,outElem =>
20       new(Local,f32, fun((accumExp, accumAcc) =>
21         accumAcc = 0.0f;
22         for(k, fun(i => accumAcc = accumExp +
23           fst(idx(i, zip(fst(idx(rowIdx, zip(A,C))),
24             fst(idx(colIdx, zip(transpose(B),
25               snd(idx(rowIdx, zip(A,C)))))))) *
26             snd(idx(i, zip(fst(idx(rowIdx, zip(A,C))),
27               fst(idx(colIdx, zip(transpose(B),
28                 snd(idx(rowIdx, zip(A,C))))))))));
29         outElem = alpha * accumExp + beta *
30           snd(idx(colIdx, zip(transpose(B),
31             snd(idx(rowIdx, zip(A,C)))))) ))));
32       syncThreads()))))
```



Codegen

Translation

C

```
33 __global__ void gemm_kernel(float* __restrict__ output,
34   int m, int n, int k, const __half* __restrict__ A,
35   const __half* __restrict__ B,
36   const float* __restrict__ C, float alpha, float beta) {
37   for(int rowIdx=blockIdx.x;
38     blockIdx.x<m; rowIdx += gridDim.x) {
39     for(int colIdx=threadIdx.x;
40       threadIdx.x<n; rowIdx += blockDim.x) {
41       float accum = 0;
42       for (int i = 0; i < k; i++) {
43         accum = accum + A[i + rowIdx*k] * B[colIdx + i*n];
44       }
45       output[colIdx + rowIdx * n] =
46         alpha * accum + beta * C[colIdx + rowIdx*n];
47     }
48     __syncthreads(); }
```

RISE: a Purely Functional Language for Optimizing via Rewriting

$E := x \mid 0.0f \mid$ *variables and literals*
 $\text{fun}(x \Rightarrow E) \mid$ *function abstraction*
 $E \mid > E \mid E(E) \mid$ *function application*
 $\text{depFun}(x: K \Rightarrow E) \mid$ *dependent fun. abstraction*
 $E(N) \mid E(DT) \mid$ *dependent fun. application*
 $\text{map} \mid \text{reduce} \mid \text{zip} \mid \dots$ *primitives*

$T := t \mid DT \mid$ *type variables & data types*
 $T \rightarrow T \mid (x: K) \rightarrow T$ *function types*

$DT := f32 \mid \dots \mid$ *scalar types*
 $\text{Array}[N, DT] \mid \text{Tuple}[DT, DT]$ *array & tuple types*

$N := 0 \mid 1 \mid \dots \mid N + N \mid N * N \mid \dots$ *natural numbers*

$K := \text{Nat} \mid \text{DataType} \mid \text{AddrSp}$ *Kinds*

Type system enforces that no functions can be stored in memory

RISE: a Purely Functional Language for Optimizing via Rewriting

map: $\{n: \text{Nat}\} \rightarrow \{s: \text{DataType}\} \rightarrow \{t: \text{DataType}\} \rightarrow$
 $(s \rightarrow t) \rightarrow \text{Array}[n, s] \rightarrow \text{Array}[n, t]$

reduce: $\{n: \text{Nat}\} \rightarrow \{t: \text{DataType}\} \rightarrow$
 $(t \rightarrow t \rightarrow t) \rightarrow t \rightarrow \text{Array}[n, t] \rightarrow t$

zip: $\{n: \text{Nat}\} \rightarrow \{s: \text{DataType}\} \rightarrow \{t: \text{DataType}\} \rightarrow$
 $\text{Array}[n, s] \rightarrow \text{Array}[n, t] \rightarrow \text{Array}[n, \text{Tuple}[s, t]]$

mapSeq: $\{n: \text{Nat}\} \rightarrow \{s: \text{DataType}\} \rightarrow \{t: \text{DataType}\} \rightarrow$
 $(s \rightarrow t) \rightarrow \text{Array}[n, s] \rightarrow \text{Array}[n, t]$

mapPar: $\{n: \text{Nat}\} \rightarrow \{s: \text{DataType}\} \rightarrow \{t: \text{DataType}\} \rightarrow$
 $(s \rightarrow t) \rightarrow \text{Array}[n, s] \rightarrow \text{Array}[n, t]$

reduceSeq: $\{n: \text{Nat}\} \rightarrow \{s: \text{DataType}\} \rightarrow \{t: \text{DataType}\} \rightarrow$
 $(t \rightarrow s \rightarrow t) \rightarrow t \rightarrow \text{Array}[n, s] \rightarrow t$

Optimizing via Rewriting

High-Level GEMM

```
1 depFun((m:Nat, n:Nat, k:Nat) =>
2   fun((A: Array[m, Array[k, f32]], B: Array[k, Array[n, f32]],
3     C: Array[m, Array[n, f32]], alpha: f32, beta: f32) =>
4   zip(A)(C) |> map(fun(rowAC =>
5     zip(B |> transpose)(snd(rowAC)) |> map(fun(colBC =>
6       zip(fst(rowAC))(fst(colBC)) |>
7       map(fun((a, b) => a * b)) |> reduce(+, 0) |>
8       fun(r => (alpha * r) + (beta * snd(colBC))) ))))))))
```

Optimization Strategy

ELEVATE

↻ Rewriting

Low-Level GEMM

```
9 depFun((m:Nat, n:Nat, k:Nat) => fun(A, B, C, alpha, beta =>
10  zip(A)(C) |> mapBlock(fun(rowAC =>
11    zip(B |> transpose)(snd(rowAC)) |>
12    mapThreads(fun(colBC => zip(fst(rowAC))(fst(colBC)) |>
13      reduceSeq(Local)(fun((acc, ab) =>
14        acc + fst(ab) * snd(ab)), 0) |>
15      fun(r => (alpha * r) + (beta * snd(colBC))) ))))))))
```

Discussed 4 weeks ago

Correctness Proof of Rewrite Rule

```
1 mapSplit : (n: ℕ) → {m: ℕ} → {s t: Set} → (f: s → t) → (xs: Vec s (m * n)) →
2     map (map f) (split n {m} xs) ≡ split n {m} (map f xs)
3 simplification : (n: ℕ) → {m: ℕ} → {t: Set} → (xs: Vec t (m*n)) → (join ∘ split n {m}) xs ≡ xs
4 {- Split-join rule proof -}
5 splitJoin : {m: ℕ} → {s: Set} → {t: Set} → (n: ℕ) → (f: s → t) → (xs: Vec s (m * n)) →
6     (join ∘ map (map f) ∘ split n {m}) xs ≡ map f xs
7 splitJoin {m} n f xs =
8     begin
9     (join ∘ map (map f) ∘ split n {m}) xs
10 ≡⟨⟩
11   join (map (map f) (split n {m} xs))
12 ≡⟨ cong join (mapSplit n {m} f xs) ⟩
13   join (split n {m} (map f xs))
14 ≡⟨ simplification n {m} (map f xs) ⟩
15   map f xs
16   ■
```

Listing 3. Proof of correctness of the `splitJoin` rewrite rule in Agda

DPIA: Combining Functional and Imperative

$P := x \mid 0.0f \mid$	<i>variables and literals</i>
$\text{fun}(x \Rightarrow P) \mid$	<i>function abstraction</i>
$P \mid > P \mid P(P) \mid$	<i>function application</i>
$\text{depFun}(x: K' \Rightarrow E) \mid$	<i>dependent fun. abstraction</i>
$P(N) \mid P(DT) \mid$	<i>dependent fun. application</i>
$\text{mapPar} \mid \text{reduceSeq} \mid \text{zip} \mid \dots$	<i>functional primitives</i>
$P = P \mid ; \mid \text{new} \mid \text{parFor} \mid \dots$	<i>imperative primitives</i>
$T' := t \mid T' \rightarrow T' \mid (x: K') \rightarrow T'$	<i>type var. & function types</i>
$T' \times T' \mid$	<i>phrase pair type</i>
$\text{Exp}[DT, RW] \mid$	<i>expression type</i>
$\text{Acc}[DT] \mid$	<i>acceptor type</i>
Comm	<i>command type</i>
$RW := \text{Rd} \mid \text{Wr}$	<i>read-write annotations</i>
$K' := K \mid RW$	<i>kinds</i>

Type system separates functional and imperative parts

DPIA: Combining Functional and Imperative

```
mapPar(n: Nat, s: DataType, t: DataType,  
        f: Exp[s, Rd] -> Exp[t, Wr],  
        in: Exp[Array[n, s], Rd]): Exp[Array[n, t], Wr]
```

```
reduceSeq(n: Nat, s: DataType, t: DataType,  
           f: Exp[t, Rd] -> Exp[s, Rd] -> Exp[t, Wr],  
           init: Exp[t, Wr],  
           in: Exp[Array[n, s], Rd]): Exp[t, Rd]
```

```
assign(t: DataType, lhs: Acc[t], rhs: Exp[t, Rd]): Comm
```

```
seq(c1: Comm, c2: Comm): Comm
```

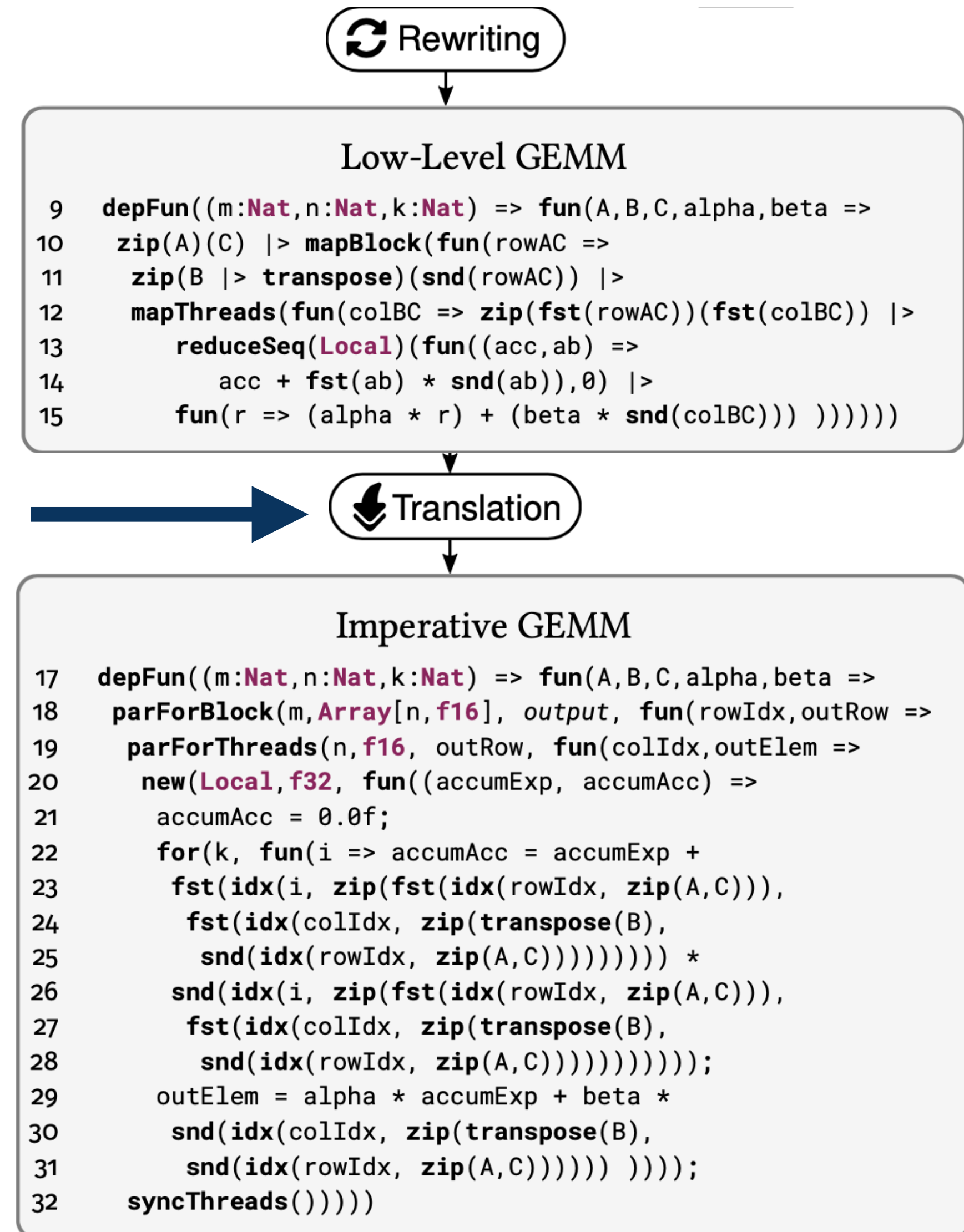
```
new(t: DataType, body: (Exp[t, Rd] x Acc[t]) -> Comm): Comm
```

```
for(n: Nat, body: Exp[Idx[n], Rd] -> Comm): Comm
```

```
parFor(n: Nat, t: DataType, out: Acc[Array[n, t]],  
        body: Exp[Idx[n], Rd] -> Acc[t] -> Comm): Comm
```

Translating From RISE to DPIA

- Translation via two mutual recursive functions
 - **Acceptor translation:**
translate an *expression* E into a command by writing the translated result into *acceptor* A :
$$accT(E, A) \approx A = E$$
 - **Continuation translation:**
translate an *expression* E into a command by passing the translated result to a *continuation* C :
$$conT(E, C) \approx C(E)$$



Translating From RISE to DPIA




```
def accT(E: Phrase[Exp[t,Wr]],
         A: Phrase[Acc[t]]): Phrase[Comm] = expr match {
  case mapSeq(n, s, t, f, in) =>
    cont(in, fun(inT =>
      for(n, fun(i =>
        accT( f(inT[i]), A[i] ) ))))
  case mapPar(n, s, t, f, in) =>
    cont(in, fun(inT =>
      parFor(n, t, A, fun((i, a) =>
        accT( f(inT[i]), a ) ))))
  case ...
}
```

Translating From RISE to DPIA

```
def conT(E: Phrase[Exp[t,Rd]],
         C: Phrase[Expr[t,Rd]] -> Phrase[Comm]): Phrase[Comm] = expr match {
  case reduceSeq(n, t, s, f, init, in) =>
    conT(in, fun(inT =>
      new(t, fun((accumAcc, accumExp) =>
        accT(init, accumAcc) ;
        for(n, fun(i =>
          accT( f(accumExpr, inT[i]), accumAcc) )) ;
        conT(accumExp, C) )) ))
  case ...
}
```

Systematically Extending Shine With Support for Tensor Cores

Bottom-up approach:

1. Add new *low-level imperative primitives* corresponding to the CUDA Tensor Core API and implement  Codegen for these primitives.
2. Add *low-level functional primitives* and implement  Translation to their imperative counterparts
3. Add *rewrite* rules to enable exploiting Tensor Cores via  Rewriting

1. Low-Level Imperative Primitives and



```
template<typename FragmKind, int m, int n, int k,  
        typename T, typename Layout=void> class fragment;  
  
void mma_sync(  
    fragment<...> &D,  
    const fragment<...> &A,  
    const fragment<...> &B,  
    const fragment<...> &C);  
void load_matrix_sync(fragment<...> &A,  
    const T* tile, unsigned l_dim, layout_t layout);  
void store_matrix_sync(T* tile,  
    const fragment<...> &A,  
    unsigned l_dim, layout_t layout);  
void fill_fragment(  
    fragment<...> &A, const T& value);
```

```
Fragment[m: Nat, n: Nat, k: Nat, t: DataType, f: FragmKind]  
  
def mmaFragment(m:Nat, n:Nat, k:Nat, s:DataType, t:DataType,  
    A: Exp[Fragment[m,k,n,s,AMatrix], Rd],  
    B: Exp[Fragment[k,n,m,s,BMatrix], Rd],  
    C: Exp[Fragment[m,n,k,t,Accum], Rd],  
    D: Acc[Fragment[m,n,k,t,Accum]]): Comm  
def loadFragment(f:FragmKind, m:Nat, n:Nat, k:Nat, t:DataType,  
    tile: Exp[Array[m,Array[n,t]], Rd], A: Acc[Fragment[m,n,k,t,f]]): Comm  
def storeFragment(m:Nat, n:Nat, k:Nat, t:DataType,  
    A: Exp[Fragment[m,n,k,t,Accum],Rd], tile: Acc[Array[m,Array[n,t]]]): Comm  
def fillFragment(f:FragmKind, m:Nat, n:Nat, k:Nat, t:DataType,  
    A: Acc[Fragment[m,n,k,t,f]], value: Exp[t, Rd]): Comm
```

- Direct representation of CUDA API as imperative primitives in RISE
- Fragment types needed to be added to RISE
- Code generation is straightforward

2. Low-Level Functional Primitives and



functional primitives

```
tensorMatMulAdd: {m: Nat} -> {n: Nat} -> {k: Nat} ->
  {s: DataType} -> {t: DataType} ->
  Fragment[m,k,n,s, AMatrix] ->
  Fragment[k,m,n,s, BMatrix] ->
  Fragment[m,n,k,t, Accum] -> Fragment[m,n,k,t, Accum]
asFragment: {m: Nat} -> {n: Nat} -> {k: Nat} ->
  {t: DataType} -> {f: FragmKind} ->
  Array[m, Array[n, t]] -> Fragment[m,n,k,t, f]
asMatrix: {m: Nat} -> {n: Nat} -> {k: Nat} -> {t: DataType} ->
  Fragment[m,n,k,t, Accum] -> Array[m, Array[n, t]]
generateFragment: {m: Nat} -> {n: Nat} -> {k: Nat} ->
  {t: DataType} -> {f: FragmKind} ->
  t -> Fragment[m,n,k,t, f]
```

imperative primitives

```
Fragment[m: Nat, n: Nat, k: Nat, t: DataType, f: FragmKind]
def mmaFragment(m:Nat, n:Nat, k:Nat, s:DataType, t:DataType,
  A: Exp[Fragment[m,k,n,s,AMatrix], Rd],
  B: Exp[Fragment[k,n,m,s,BMatrix], Rd],
  C: Exp[Fragment[m,n,k,t,Accum], Rd],
  D: Acc[Fragment[m,n,k,t,Accum]]): Comm
def loadFragment(f:FragmKind, m:Nat, n:Nat, k:Nat, t:DataType,
  tile: Exp[Array[m,Array[n,t]], Rd], A: Acc[Fragment[m,n,k,t,f]]): Comm
def storeFragment(m:Nat, n:Nat, k:Nat, t:DataType,
  A: Exp[Fragment[m,n,k,t,Accum],Rd], tile: Acc[Array[m,Array[n,t]]]): Comm
def fillFragment(f:FragmKind, m:Nat, n:Nat, k:Nat, t:DataType,
  A: Acc[Fragment[m,n,k,t,f]], value: Exp[t, Rd]): Comm
```

- One *low-level functional primitive* per *imperative primitive*
- Functional primitives have return values, rather than returning nothing (i.e. `void` / `Comm`)
- loading / storing a fragment corresponds to turning a matrix into a fragment (and reverse)

2. Low-Level Functional Primitives and



- Translation by a case for each low-level functional primitive

```
def accT(expr: Phrase[Exp[d,Wr]],
         output: Phrase[Acc[t]]): Phrase[Comm] = expr match {
case tensorMatMulAdd(m,n,k,dt,dtAcc,aMatrix,bMatrix,cMatrix)
=> conT(aMatrix, fun(aMatrix => conT(bMatrix,
fun(bMatrix => conT(cMatrix, fun(cMatrix =>
mmaFragment(m, n, k, dt,
dtAcc, aMatrix, bMatrix, cMatrix, A))))))
case asFragment(m, n, k, dt, f, tile)
=> conT(tile, fun(tile: =>
loadFragment(f, m, n, k, dt, tile, A)))
case asMatrix(m, n, k, dt, frag)
=> conT(frag, fun(frag: =>
storeFragment(m, n, k, dt, frag, A)))
case generateFragment(m, n, k, dt, f, fill)
=> conT(fill, fun(fill =>
fillFragment(f, m, n, k, dt, fill, A)))
... }
```

3. Add Rewrite Rules To Enable Rewriting

- Rewrite rules enable automatic exploitation of Tensor Cores
- Examples shows automatic use of Tensor Cores for high-level matrix multiplication code
- Rewrite rules can be applied *automatically* [GPGPU'16, ICFP'15], *manually* [ICFP'20], or *guided* [arXiv:2111.13040].

```
aTile: Array[16, Array[16, f16]] |> map(fun(aRow =>
bTile: Array[16, Array[16, f16]] |> map(fun(bCol =>
zip(aRow, bCol) |>
reduceSeq(fun(ac, ab =>
add(ac, mul(fst(ab), snd(ab)))))(0.0))))))
```



```
tensorMatMulAdd
(aTile: Array[16, Array[16, f16]] |> asFragment |> toMem(Local))
(bTile: Array[16, Array[16, f16]] |> transpose
|> asFragment |> toMem(Local))
(generateFragment(0.0) |> toMem(Local))
|> toMem(Local) |> asMatrix
```

RISE & Shine: Language-Oriented Compiler Design

- **Shine** demonstrates an extensible compiler design allowing targeting specialised hardware

- Progressive compilation is a good idea:

High-level functional primitives via  Rewriting to

low-level functional primitives via  Translation to

low-level imperative primitives via  Codegen to

low-level imperative code.

<https://rise-lang.org/>