# Achieving High-Performance the Functional Way

## Expressing High-Performance Optimisations as Rewrite Strategies

Bastian Hagedorn, Johannes Lenfers, Thomas Kœhler, Xueying Qin, Sergei Gorlatch and **Michel Steuwer**

WWU MÜNSTER · THE UNIVERSITY of EDINBURGH · University of Glasgow
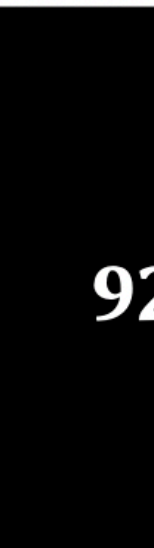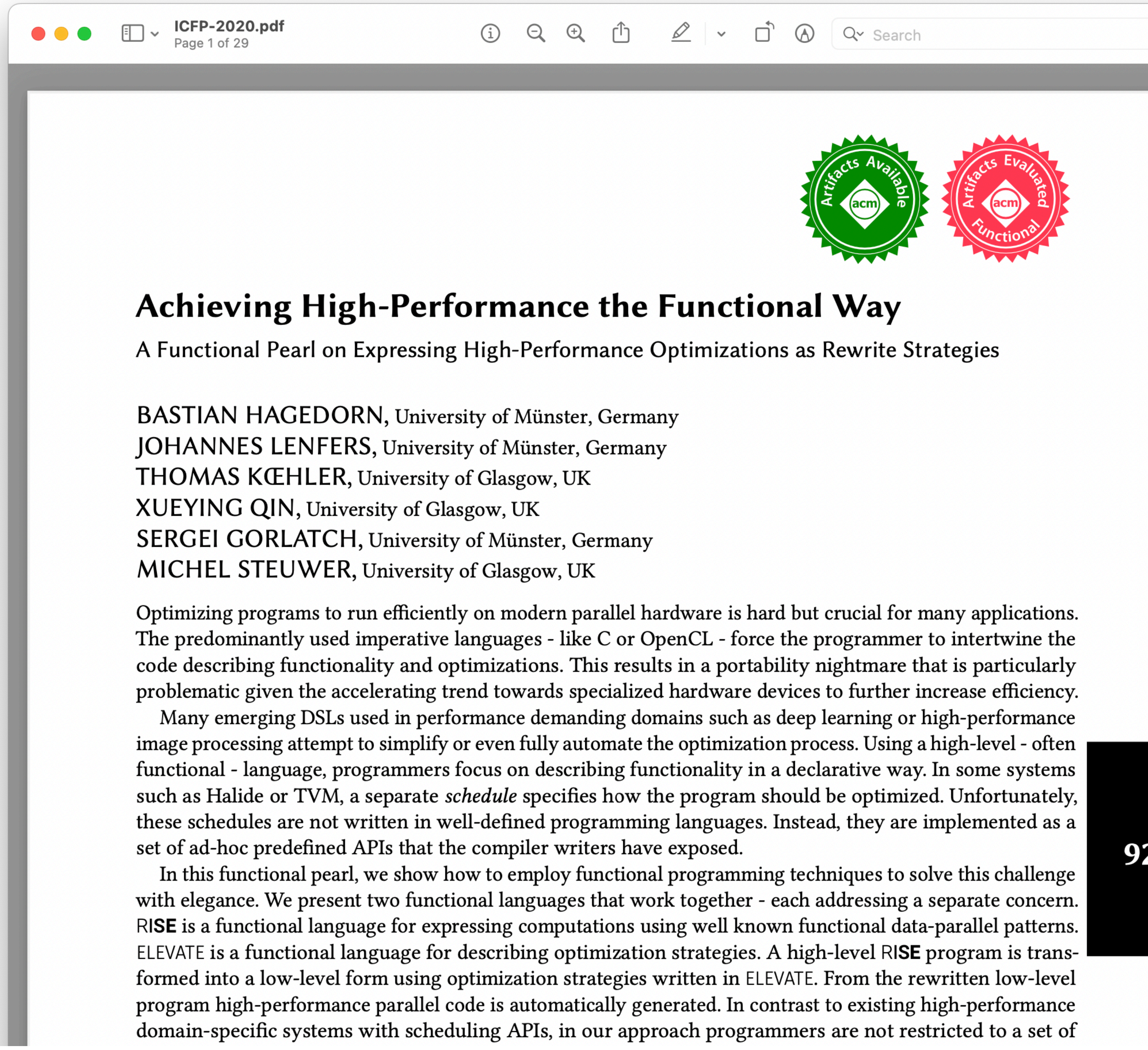
- rejected from PLDI 2020 🙁

- published at 📄 [ICFP 2020] 🙂

- selected as 1 of 4 **ACM SIGPLAN Research Highlights** from 2020 ☺️

- selected for publication as **Research Highlight** in an upcoming issue of the **Communications of the ACM** 🥰

# Achieving High-Performance the Functional Way

A Functional Pearl on Expressing High-Performance Optimizations as Rewrite Strategies

BASTIAN HAGEDORN, University of Münster, Germany
JOHANNES LENFERS, University of Münster, Germany
THOMAS KŒHLER, University of Glasgow, UK
XUEYING QIN, University of Glasgow, UK
SERGEI GORLATCH, University of Münster, Germany
MICHEL STEUWER, University of Glasgow, UK

Optimizing programs to run efficiently on modern parallel hardware is hard but crucial for many applications. The predominantly used imperative languages - like C or OpenCL - force the programmer to intertwine the code describing functionality and optimizations. This results in a portability nightmare that is particularly problematic given the accelerating trend towards specialized hardware devices to further increase efficiency.

Many emerging DSLs used in performance demanding domains such as deep learning or high-performance image processing attempt to simplify or even fully automate the optimization process. Using a high-level - often functional - language, programmers focus on describing functionality in a declarative way. In some systems such as Halide or TVM, a separate *schedule* specifies how the program should be optimized. Unfortunately, these schedules are not written in well-defined programming languages. Instead, they are implemented as a set of ad-hoc predefined APIs that the compiler writers have exposed.

In this functional pearl, we show how to employ functional programming techniques to solve this challenge with elegance. We present two functional languages that work together - each addressing a separate concern. RISE is a functional language for expressing computations using well known functional data-parallel patterns. ELEVATE is a functional language for describing optimization strategies. A high-level RISE program is transformed into a low-level form using optimization strategies written in ELEVATE. From the rewritten low-level program high-performance parallel code is automatically generated. In contrast to existing high-performance domain-specific systems with scheduling APIs, in our approach programmers are not restricted to a set of

Eelco Visser

Richard Bird

# HIGH-PERFORMANCE: *Why do we care?*

> **Elliot Turner**
> @eturner303
>
> Holy crap: It costs $245,000 to train the XLNet model (the one that's beating BERT on NLP tasks..512 TPU v3 chips * 2.5 days * $8 a TPU) – arxiv.org/abs/1906.08237

# HIGH-PERFORMANCE: *Why do we care?*

> **Elliot Turner**
> @eturner303
>
> Holy crap: It costs $245,000 to train the XLNet model (the one that's beating BERT on NLP tasks..512 TPU v3 chips * 2.5 days * $8 a TPU) - arxiv.org/abs/1906.08237

> **Elliot Turner**
> @eturner303
>
> Another way (using carbon as opposed to $$) of thinking about this experiment:  Training XLNet to convergence releases around 4.9 metric tons of $CO_2$ into the atmosphere (equivalent to driving a car around 11,000 miles)

# ...RMANCE: *Why do we care?*



**Turner**
...rner303

It costs $245,000 to train the XLNet model (the
... beating BERT on NLP tasks...512 TPU v3 chips * 2.5
...a TPU) - arxiv.org/abs/1906.08237

**Turner**
...rner303

...way (using carbon as opposed to $$) of thinking
... experiment: Training XLNet to convergence
...round 4.9 metric tons of CO2 into the atmosphere
...t to driving a car around 11,000 miles)

# Achieving High-Performance the ~~Functional~~ Way
## Manual

```cpp
__global__ void matmul(
    float *A, float *B, float *C,
    int K, int M, int N) {

  int x = blockIdx.x * blockDim.x + threadIdx.x;
  int y = blockIdx.y * blockDim.y + threadIdx.y;
  float acc = 0.0;

  for (int k = 0; k < K; k++) {
    acc += A[y * M + k] * B[k * N + x];
  }

  C[y * N + x] = acc;
}
```

Naive Matrix Multiplication in **NVIDIA CUDA**

# Achieving High-Performanc

```cuda
__global__ void matmul(
        float *A, float *B, float *C,
        int K, int M, int N) {

    int x = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;
    float acc = 0.0;

    for (int k = 0; k < K; k++) {
        acc += A[y * M + k] * B[k * N + x];
    }

    C[y * N + x] = acc;
}
```

Naive Matrix Multiplication in

**100-1000x** performance

Optimized Matrix Multiplication

1000x CO2 Improvement

100-1000x performance

Optimized Matrix Multiplication

**1000x CO2** Improvement

**100-1000x** performance
**30x** lines of code
time-intensive + error-prone

Optimized Matrix Multiplication

# Achieving High-Performance the ~~Functional~~ Decoupled Way



```
k = tvm.reduce_axis((0, K), 'k')
A = tvm.placeholder((M, K), name='A')
B = tvm.placeholder((K, N), name='B')
C = tvm.compute((M, N), lambda x, y:
 tvm.sum(A[x, k] * B[k, y], axis=k),
 name='C')
```

Algorithm

```
# blocking version
xo, yo, xi, yi = s[C].tile(
  C.op.axis[0],C.op.axis[1],32,32)
k,    = s[C].op.reduce_axis
ko, ki = s[C].split(k, factor=4)
s[C].reorder(xo, yo, ko, ki, xi, yi)
```

Schedule

Domain-Scientist

Performance Engineer

Compiler

https://tvm.apache.org/docs/tutorials/optimize/opt_gemm.html

tvm
0.7.dev1

Search docs

Docs » Get Started Tutorials » How to optimize GEMM on CPU

View page source

**ⓘ Note**

Click here to download the full example code

## How to optimize GEMM on CPU

**Author:** Jian Weng, Ruofei Yu

HOW TO
Installation
Contribute to TVM
Deploy and Integration
Developer How-To Guide

11

# Achieving High-Performance the ~~Functional~~ Way
## Decoupled



**Algorithm**

```
k = tvm.reduce_axis((0, K), 'k')
A = tvm.placeholder((M, K), name='A')
B = tvm.placeholder((K, N), name='B')
C = tvm.compute((M, N), lambda x, y:
 tvm.sum(A[x, k] * B[k, y], axis=k),
 name='C')
```

**Domain-Scientist**

**Performance Engineer**

**Compiler**

**Schedule**

```
# "parallel schedule
s = tvm.create_schedule(C.op)
CC = s.cache_write(C, 'global')
xo, yo, xi, yi = s[C].tile(
  C.op.axis[0], C.op.axis[1], bn, bn)

s[CC].compute_at(s[C], yo)
xc, yc = s[CC].op.axis
k, = s[CC].op.reduce_axis
ko, ki = s[CC].split(k, factor=4)
s[CC].reorder(ko, xc, ki, yc)
s[CC].unroll(ki)
s[CC].vectorize(yc)
s[C].parallel(xo)
x, y, z = s[packedB].op.axis
s[packedB].vectorize(z)
s[packedB].parallel(x)
```

**200x**

TVM

Absolute Runtime (ms)

2,000
1,000
500
200
100
50

baseline    vectorization    array-packing    parallel
    blocking    loop-perm    cache-blocks

# Achieving High-Performance the ~~Functional~~ Way
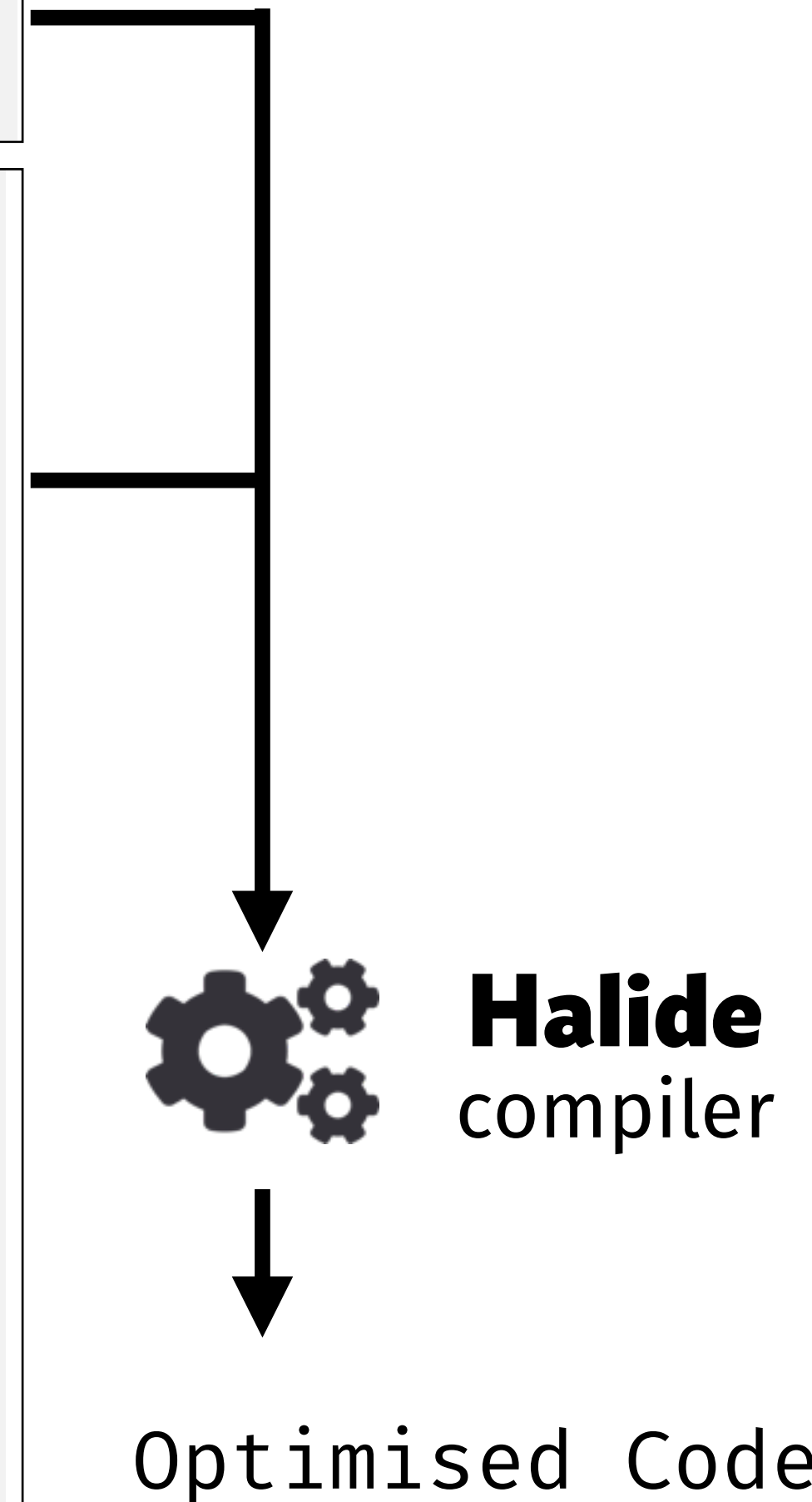## Decoupled

Program



Domain-Scientist

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```
📄 Algorithm

Performance Engineer

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile    =  3; const int y_tile    = 4;
const int y_unroll  =  8; const int r_unroll  = 1;
Var xi,yi,xio,xii,yii,xo,yo,x_pair,xiio,ty; RVar rxo,rxi;
out.bound(x, 0, size).bound(y, 0, size)
    .tile(x, y, xi, yi, x_tile * vec_size * warp_size,
          y_tile * y_unroll)
    .split(yi, ty, yi, y_unroll)
    .vectorize(xi, vec_size)
    .split(xi, xio, xii, warp_size)
    .reorder(xio, yi, xii, ty, x, y)
    .unroll(xio).unroll(yi)
    .gpu_blocks(x, y).gpu_threads(ty).gpu_lanes(xii);
prod.store_in(MemoryType::Register).compute_at(out, x)
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .unroll(xo).unroll(y).update()
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .split(r.x, rxo, rxi, warp_size)
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
     .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
     .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
     .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
     .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
     .unroll(xo).unroll(Ay);
```
📄 Schedule

## Compilers with scheduling APIs

**Halide**



Tiramisu-Compiler / **tiramisu**

Fireiron  



**Halide**
compiler

Optimised Code

Optimisation Schedule

# Problems with Scheduling APIs

## Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   = 3; const int y_tile   = 4;
const int y_unroll =  8; const int r_unroll = 1;
                                              rxo,rxi;
                                              size,
```

```
prod.store_in(MemoryType::Register).compute_at(out, x)
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .unroll(xo).unroll(y).update()
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .split(r.x, rxo, rxi, warp_size)
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
```

```
                                                   ii);
                                                t, x)
                                                Jp)

                                             es(xi)

                                                Jp)

                                             es(xi)
      .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
      .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
     .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
     .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
     .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
     .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
     .unroll(xo).unroll(Ay);
```

## Optimisation Schedule

**Halide**
compiler

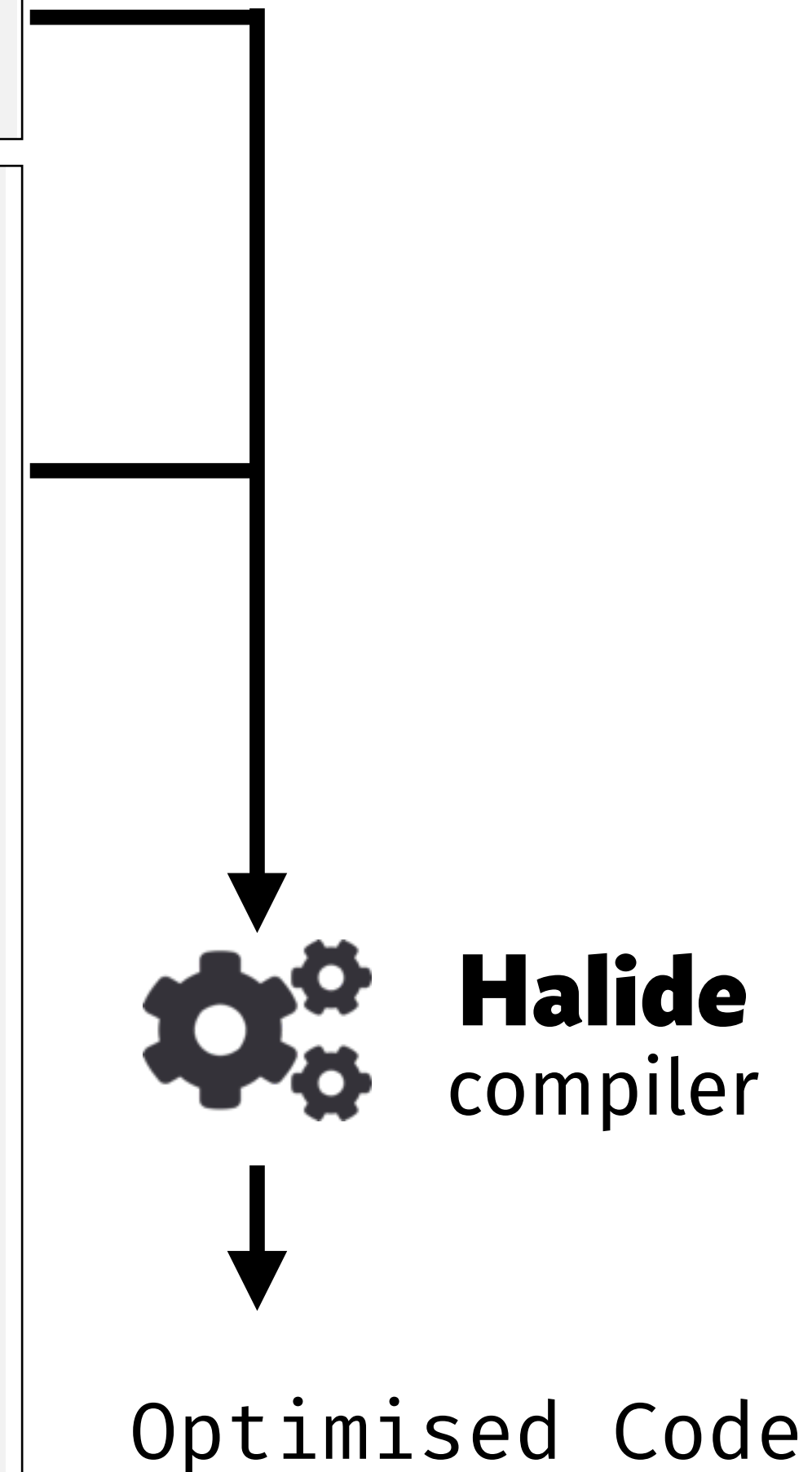Optimised Code

# Problems with Scheduling APIs

## Program

```cpp
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

No clear separation

```cpp
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   =  3; const int y_tile   = 4;
const int y_unroll =  8; const int r_unroll = 1;
```

```cpp
prod.store_in(MemoryType::Register).compute_at(out, x)
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .unroll(xo).unroll(y).update()
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .split(r.x, rxo, rxi, warp_size)
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
```

```cpp
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
    .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
    .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
    .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
    .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
    .unroll(xo).unroll(Ay);
```

## Optimisation Schedule

**Halide** compiler

Optimised Code

# Problems with Scheduling APIs

## Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

N**Hinders reuse**

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   = 3; const int y_tile   = 4;
const int y_unroll = 8; const int r_unroll = 1;
```
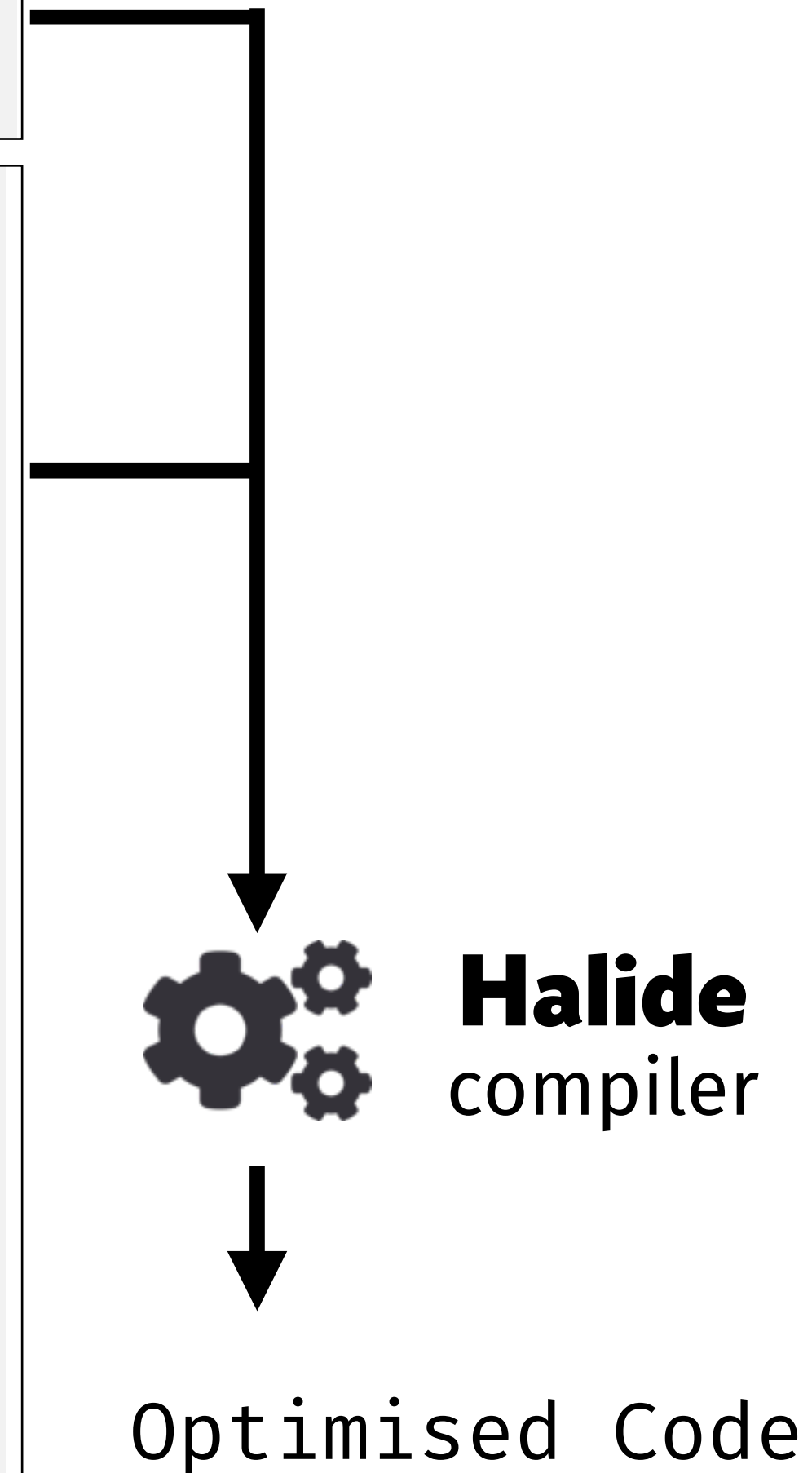
```
prod.store_in(MemoryType::Register).compute_at(out, x)
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .unroll(xo).unroll(y).update()
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .split(r.x, rxo, rxi, warp_size)
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
```

```
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
     .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
     .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
     .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
     .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
     .unroll(xo).unroll(Ay);
```

## Optimisation Schedule

**Halide**
compiler

Optimised Code

14

# Problems with Scheduling APIs

## Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

N...  **Hinders reuse**

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile    = 3; const int y_tile    = 4;
const int y_unroll  = 8; const int r_unroll  = 1;
```
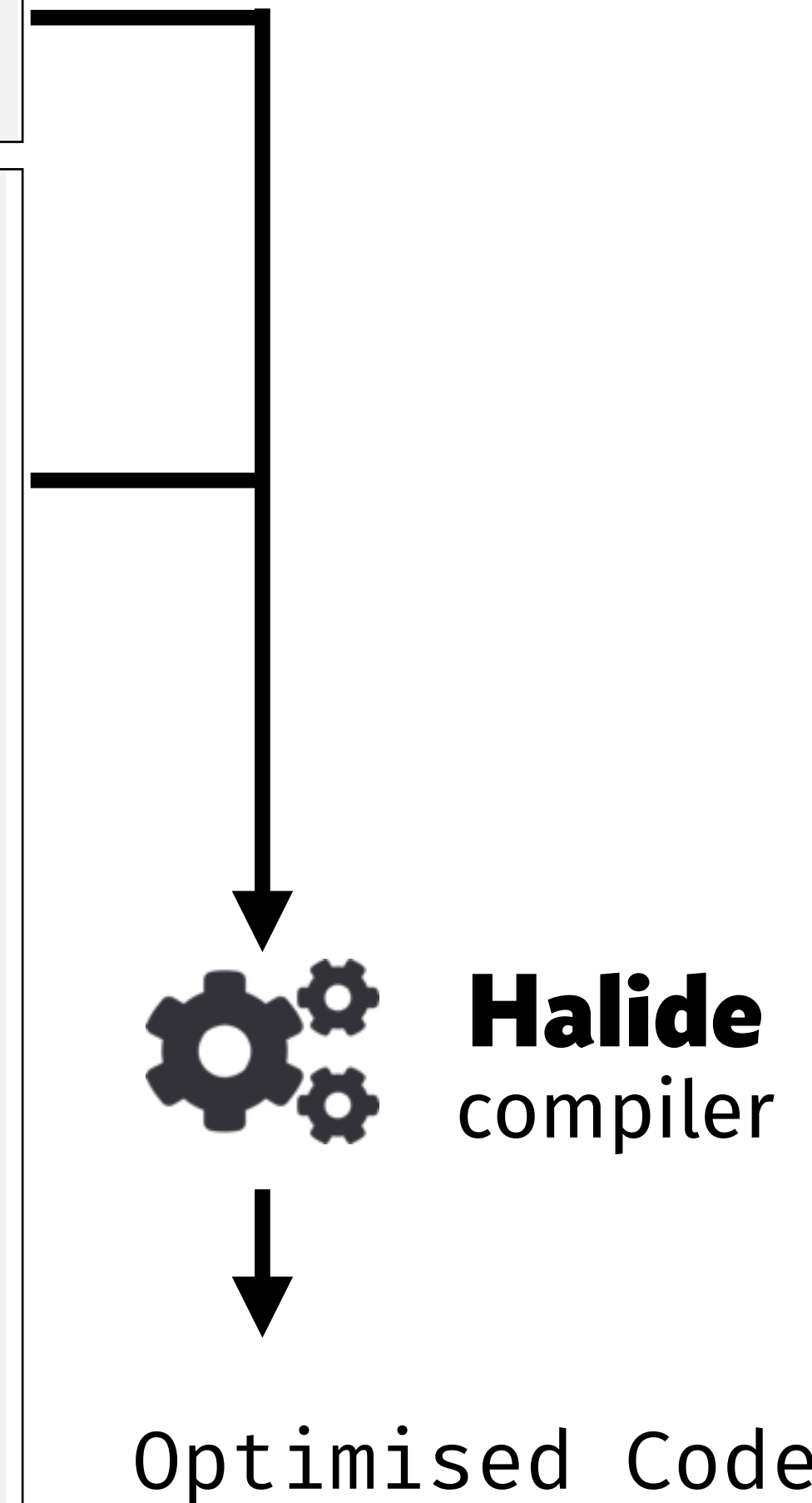
Not well defined
semantics

```
prod.store_in(MemoryType::Register).compute_at(out, x)
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .unroll(xo).unroll(y).update()
    .split(x, xo, xi, warp_size * vec_size, RoundUp)
    .split(y, ty, y, y_unroll)
    .gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
    .split(r.x, rxo, rxi, warp_size)
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
```

```
        .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
        .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
      .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
      .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
      .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
      .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
      .unroll(xo).unroll(Ay);
```

## Optimisation Schedule

**Halide**
compiler

Optimised Code

# Problems with Scheduling APIs

Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

N...

**Hinders reuse**

Not well d...

**Hinders understanding**

...ntics

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   =  3; const int y_tile   = 4;
const int y_unroll =  8; const int r_unroll = 1;
```
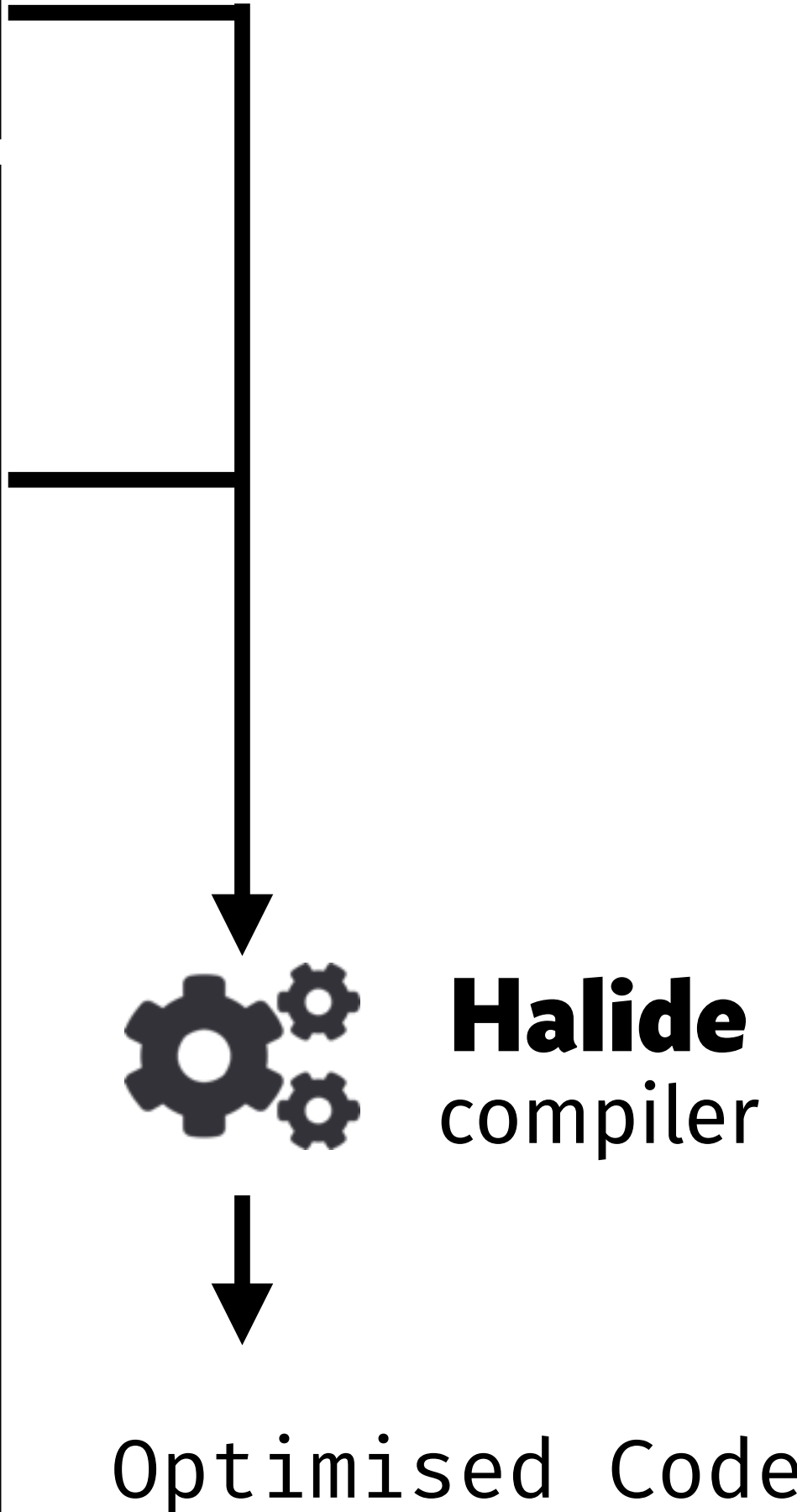
```
store_in(MemoryType::Register).compute_at(out, x)              rxo,rxi;
.split(x, xo, xi, warp_size * vec_size, RoundUp)              size,
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.unroll(xo).unroll(y).update()
.split(x, xo, xi, warp_size * vec_size, RoundUp)              ii);
.split(y, ty, y, y_unroll)                                    t, x)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)           Jp)
.split(r.x, rxo, rxi, warp_size)                              es(xi)
.unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)       Jp)
.unroll(xo).unroll(y);                                        es(xi)
```

```
        .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
        .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
        .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
        .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
        .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
        .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
        .unroll(xo).unroll(Ay);
```

Optimisation Schedule

**Halide**
compiler

Optimised Code

14

# Problems with Scheduling APIs

Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

**Hinders reuse**

Not ... ...ation

**Hinders understanding**

Not well ...

...ntics

Only fixed built-in
optimisations

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   =  3; const int y_tile   = 4;
const int y_unroll =  8; const int r_unroll = 1;
```

```
.store_in(MemoryType::Register).compute_at(out, x)
.split(x, xo, xi, warp_size * vec_size, RoundUp)
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.unroll(xo).unroll(y).update()
.split(x, xo, xi, warp_size * vec_size, RoundUp)
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.split(r.x, rxo, rxi, warp_size)
.unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
.unroll(xo).unroll(y);
```

```
          .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
          .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
       .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
       .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
       .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
       .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
       .unroll(xo).unroll(Ay);
```

Optimisation Schedule

**Halide**
compiler

Optimised Code

# Problems with Scheduling APIs

Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

N̶ Hinders reuse ̶a̶t̶i̶o̶n

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   =  3; const int y_tile   = 4;
const int y_unroll =  8; const int r_unroll = 1;
```

Not well d̶ Hinders understanding

```
store_in(MemoryType::Register).compute_at(out, x)
.split(x, xo, xi, warp_size * vec_size, RoundUp)
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.unroll(xo).unroll(y).update()
.split(x, xo, xi, warp_size * vec_size, RoundUp)
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.split(r.x, rxo, rxi, warp_size)
.unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
.unroll(xo).unroll(y);
```

Only fix̶ No extensibility ̶a̶t̶i̶o̶n̶s̶

Optimisation Schedule

```
.unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
.unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];
B.in().compute_at(prod, ty).split(Bx, xo, xi, warp_size)
    .gpu_lanes(xi).unroll(xo).unroll(By);
A.in().compute_at(prod, rxo).vectorize(Ax, vec_size)
    .split(Ax,xo,xi,warp_size).gpu_lanes(xi).unroll(xo)
    .split(Ay,yo,yi,y_tile).gpu_threads(yi).unroll(yo);
A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
    .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
    .unroll(xo).unroll(Ay);
```

**Halide**
compiler

Optimised Code

# Problems with Scheduling APIs

Program

```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
 out(x, y)  = prod(x, y);
```

N                    ation

**Hinders reuse**

Not well d-f-

**Hinders understanding**

```
store_in(MemoryType::Register).compute_at(out, x)
.split(x, xo, xi, warp_size * vec_size, RoundUp)
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.unroll(xo).unroll(y).update()
.split(x, xo, xi, warp_size * vec_size, RoundUp)
.split(y, ty, y, y_unroll)
.gpu_threads(ty).unroll(xi, vec_size).gpu_lanes(xi)
.split(r.x, rxo, rxi, warp_size)
.unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
.unroll(xo).unroll(y);
```

Only fixe-th

**No extensibility**

```
// schedule for Nvidida GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile   =  3; const int y_tile   = 4;
const int y_unroll =  8; const int r_unroll = 1;
                                                rxo,rxi;
                                                size,
                                                ii);
                                                t, x)
                                                Up)
                                                es(xi)
                                                Up)
                                                es(xi)
    .unroll(rxi, r_unroll).reorder(xi, xo, y, rxi, ty, rxo)
    .unroll(xo).unroll(y);
Var Bx = B.in().args()[0], By = B.in().args()[1];
Var Ax = A.in().args()[0], Ay = A.in().args()[1];

A.in().in().compute_at(prod, rxi).vectorize(Ax, vec_size)
    .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
    .unroll(xo).unroll(Ay);
```

Optimisation Schedule

**Halide**
compiler

Code

**We should aim for more principled ways to describe and apply optimisations**

# The Need for a Principled Way to Separate, Describe and Apply Optimizations

Our goals:

1. ## Separate concerns
   Computations should be expressed at a high abstraction level only.
   They should not be changed to express optimizations;

2. ## Facilitate reuse
   Optimization strategies should be defined clearly separated from the computational program facilitating reusability of computational programs and strategies;

3. ## Enable composability
   Computations *and* strategies should be written as compositions of user-defined building blocks (possibly domain-specific ones); both languages should facilitate the creation of higher-level abstractions;

4. ## Allow reasoning
   *Computational patterns, but also especially strategies, should have a precise, well-defined semantics allowing reasoning about them;*

5. ## Be explicit
   *Implicit default behavior should be avoided to empower users to be in control.*

Our goals:

1. *Separate concerns*

Computations should be expressed at a high abstraction level only.

Fundamentally we argue that a more principled high-performance code generation approach should be holistic by considering *computation* and *optimization strategies* **equally important.**

**As a consequence, a strategy language should be built with the same standards as a language describing computation.**

Computational patterns, but also especially strategies, should have a precise, well-defined semantics allowing reasoning about them;

5. *Be explicit*

Implicit default behavior should be avoided to empower users to be in control.

# Achieving High-Performance the **Functional** Way

```
// Matrix Matrix Multiplication in RISE
val dot = fun(as, fun(bs,
  zip(as)(bs) |> map(fun(ab, mult(fst(ab))(snd(ab)))) |> reduce(add)(0) ) )
val mm  = fun(a : M.K.float, fun(b : K.N.float,
  a |> map(fun(arow,                    // iterating over M
    transpose(b) |> map(fun(bcol,       // iterating over N
      dot(arow)(bcol) )))) ) )          // iterating over K
```

```
val loopPerm = (
  tile(32,32)         '@' outermost(mapNest(2))        ';;'
  fissionReduceMap '@' outermost(appliedReduce) ';;'
  split(4)            '@' innermost(appliedReduce) ';;'
  reorder(Seq(1,2,5,3,6,4))                            ';;'
  vectorize(32)       '@' innermost(isApp(isApp(isMap))))
(loopPerm ';' lowerToC)(mm)
```



based on Lift
[ICFP 2015] by Steuwer et. al.

based on Stratego
[ICFP 1998] by Visser et. al.

# ELEVATE **A Language for Describing Optimisation Strategies**

- A **Strategy** encodes a program transformation as a function:

```
type Strategy[P] = P ⇒ RewriteResult[P]
```

- A **RewriteResult** encodes its success or failure:

```
RewriteResult[P] = Success[P](p: P)
                 | Failure[P](s: Strategy[P])
```

# Rewrite Rules in ELEVATE

- *Rewrite rules* are basic strategies

$$\textbf{map}(\,f\,)\,\textbf{<<}\,\textbf{map}(\,g\,)\,\rightsquigarrow\,\textbf{map}(\,f\,\textbf{<<}\,g\,)$$

```
def mapFusion: Strategy[Rise] =
  (p: Rise) ⟹ p match {
    case app(app(map, f),
             app(app(map, g), xs)) =
      Success( map(fun(x ⟹ f(g(x))), xs) )
    case _ = Failure(mapFusion)
}
```

mapFusion(  ) =

# Combinators in ELEVATE

- Building more complex strategies from simpler once

- Sequential Composition (;)

```
def seq[P]: Strategy[P] ⇒ Strategy[P] ⇒ Strategy[P] =
        fs ⇒ ss ⇒ p ⇒ fs(p).flatMapSuccess(ss)
```

- Left Choice (<+)

```
def lChoice[P]: Strategy[P] ⇒ Strategy[P] ⇒ Strategy[P] =
        fs ⇒ ss ⇒ p ⇒ fs(p).flatMapFailure(_ ⇒ ss(p))
```

- Try

```
def try[P]: Strategy[P] ⇒ Strategy[P] =
        s ⇒ p ⇒ (s <+ id)(p)
```

- Repeat

```
def repeat[P]: Strategy[P] ⇒ Strategy[P] =
        s ⇒ p ⇒ try(s ; repeat(s))(p)
```

# Traversals in ELEVATE

- Describing Precise Locations

# Traversals in ELEVATE

- Describing Precise Locations

```
def body: Strategy[Rise] ⇒ Strategy[Rise] =
 s ⇒ p ⇒ p match {
   case fun(x,b) ⇒ s(b).mapSuccess(nb ⇒
fun(x,nb))
   case _ ⇒ Failure( body(s) )
}
```

body(mapFusion) (  ) = ?

threemaps = fun(xs, map(f)(map(g)(map(h)(xs))))

# Traversals in ELEVATE

- Describing Precise Locations

```
def body: Strategy[Rise] ⇒ Strategy[Rise] =
 s ⇒ p ⇒ p match {
  case fun(x,b) ⇒ s(b).mapSuccess(nb ⇒
fun(x,nb))
  case _ ⇒ Failure( body(s) )
}
```

body(argument(mapFusion)) ( ) = ?

```
def argument: Strategy[Rise] ⇒ Strategy[Rise] =
 s ⇒ p ⇒ p match {
  case app(f,a) ⇒ s(a).mapSuccess(na ⇒
app(f,na))
  case _ ⇒ Failure( argument(s) )
}
```



threemaps = fun(xs, map(f)(map(g)(map(h)(xs))))

# Complex Traversals + Normalization in ELEVATE

- With three basic generic traversals

```
type Traversal[P] = Strategy[P] => Strategy[P]
def all[P]: Traversal[P];     def one[P]: Traversal[P];     def some[P]: Traversal[P]
```

- we define more complex traversals:

```
def      topDown[P]: Traversal[P] = s => p => (s <+ one(topDown(s)))(p)
def     bottomUp[P]: Traversal[P] = s => p => (one(bottomUp(s)) <+ s)(p)
def  allTopDown[P]: Traversal[P] = s => p => (s ';' all(allTopDown(s)))(p)
def allBottomUp[P]: Traversal[P] = s => p => (all(allBottomUp(s)) ';' s)(p)
def      tryAll[P]: Traversal[P] = s => p => (all(tryAll(try(s))) ';' try(s))(p)
```

- With these traversals we define normal forms, e.g. $\beta\eta$-normal-form:

```
def normalize[P]: Strategy[P] => Strategy[P] = s => p => repeat(topDown(s))(p)

def BENF = normalize(betaReduction <+ etaReduction)
```

# Complex optimisations defined as strategies



```
def tile: Int → Int → Strategy =
  (dim) ⟹ (n) ⟹ dim match {
    case 1 = function(splitJoin(n))

    case 2 = fmap(function(splitJoin(n))) ;
             function(splitJoin(n)) ; interchange(2)
    case i = fmap(tile(dim-1, n)) ;
             function(splitJoin(n)) ; interchange(n)
  }
```

Tiling defined as composition of rewrites not a built-in!

# Case Study: Implementing TVM's Scheduling API

- We attempt to express the same optimizations described in the TVM tutorial:

# Optimizing Matrix Matrix Multiplication with ELEVATE **Strategies**

## RISE

```
1  // matrix multiplication in RISE
2  val dot = fun(as, fun(bs, zip(as)(bs) |>
3    map(fun(ab, mult(fst(ab))(snd(ab)))) |>
4    reduce(add)(0) ) )
5  val mm = fun(a, fun(b, a |>
6    map( fun(arow, transpose(b) |>
7      map( fun(bcol,
8        dot(arow)(bcol) )))) ))
```

```
1  // baseline strategy in ELEVATE
2  val baseline = ( DFNF ';'
3    fuseReduceMap '@' topDown )
4  (baseline ';' lowerToC)(mm)
```

## tvm

```
1   # Naive matrix multiplication algorithm
2   k = tvm.reduce_axis((0, K), 'k')
3   A = tvm.placeholder((M, K), name='A')
4   B = tvm.placeholder((K, N), name='B')
5   C = tvm.compute((M, N),lambda x, y:
6       tvm.sum(A[x, k] * B[k, y],
7       axis=k),name='C')
8
9
10
11
12  # TVM default schedule
13  s = tvm.create_schedule(C.op)
```

## ELEVATE

*Baseline* Strategy

# Optimizing Matrix Matrix Multiplication with ELEVATE Strategies

Clear separation of concerns

## RISE

```
1  // matrix multiplication in RISE
2  val dot = fun(as, fun(bs, zip(as)(bs) |>
3    map(fun(ab, mult(fst(ab))(snd(ab)))) |>
4    reduce(add)(0) ) )
5  val mm = fun(a, fun(b, a |>
6    map( fun(arow, transpose(b) |>
7      map( fun(bcol,
8        dot(arow)(bcol) )))) ))
```

```
1  // baseline strategy in ELEVATE
2  val baseline = ( DFNF ';'
3    fuseReduceMap '@' topDown )
4  (baseline ';' lowerToC)(mm)
```



```
1   # Naive matrix multiplication algorithm
2   k = tvm.reduce_axis((0, K), 'k')
3   A = tvm.placeholder((M, K), name='A')
4   B = tvm.placeholder((K, N), name='B')
5   C = tvm.compute((M, N),lambda x, y:
6       tvm.sum(A[x, k] * B[k, y],
7       axis=k),name='C')
8
9
10
11
12  # TVM default schedule
13  s = tvm.create_schedule(C.op)
```

## ELEVATE

Be explicit

Implicit behavior

Enable composability

*Baseline* Strategy

28

# Optimizing Matrix Matrix Multiplication with ELEVATE **Strategies**

ELEVATE



```
1  val loopPerm = (
2   tile(32,32)         '@' outermost(mapNest(2))      ';;'
3   fissionReduceMap '@' outermost(appliedReduce) ';;'
4   split(4)            '@' innermost(appliedReduce) ';;'
5   reorder(Seq(1,2,5,3,6,4))                          ';;'
6   vectorize(32)       '@' innermost(isApp(isApp(isMap))))
7  (loopPerm ';' lowerToC)(mm)
```

```
1  xo, yo, xi, yi = s[C].tile(
2    C.op.axis[0],C.op.axis[1],32,32)
3  k,              = s[C].op.reduce_axis
4  ko, ki          = s[C].split(k, factor=4)
5  s[C].reorder(xo, yo, ko, xi, ki, yi)
6  s[C].vectorize(yi)
```

*Loop Permutation with blocking* Strategy

# Optimizing Matrix Matrix Multiplication with ELEVATE **Strategies**

ELEVATE

User-defined vs. build in

Facilitate reuse

```
1  val loopPerm = (
2    tile(32,32)        '@' outermost(mapNest(2))           ';;'
3    fissionReduceMap '@' outermost(appliedReduce) ';;'
4    split(4)           '@' innermost(appliedReduce) ';;'
5    reorder(Seq(1,2,5,3,6,4))                              ';;'
6    vectorize(32)      '@' innermost(isApp(isApp(isMap))))
7  (loopPerm ';' lowerToC)(mm)
```

```
1  xo, yo, xi, yi = s[C].tile(
2    C.op.axis[0],C.op.axis[1],32,32)
3  k,               = s[C].op.reduce_axis
4  ko, ki           = s[C].split(k, factor=4)
5  s[C].reorder(xo, yo, ko, xi, ki, yi)
6  s[C].vectorize(yi)
```

No clear separation
of concerns

*Loop Permutation with blocking* Strategy

# Optimizing Matrix Matrix Multiplication with ELEVATE Strategies

## ELEVATE

```
1   val appliedMap = isApp(isApp(isMap))
2   val isTransposedB = isApp(isTranspose)
3
4   val packB = storeInMemory(isTransposedB,
5    permuteB ';;'
6    vectorize(32) '@' innermost(appliedMap) ';;'
7    parallel       '@' outermost(isMap)
8   ) '@' inLambda
9
10  val arrayPacking = packB ';;' loopPerm
11  (arrayPacking ';' lowerToC )(mm)
```
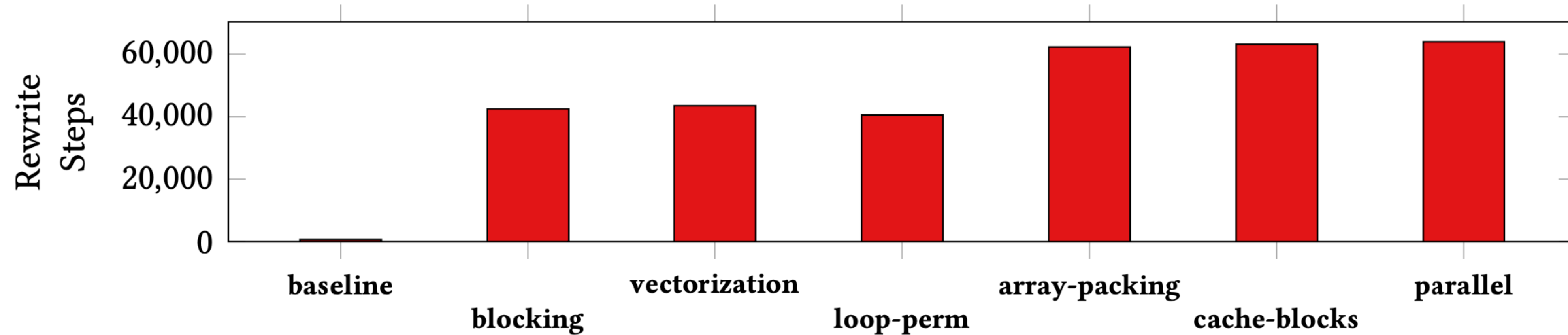
## tvm

```
1   # Modified algorithm
2   bn = 32
3   k = tvm.reduce_axis((0, K), 'k')
4   A = tvm.placeholder((M, K), name='A')
5   B = tvm.placeholder((K, N), name='B')
6   pB = tvm.compute((N / bn, K, bn),
7     lambda x, y, z: B[y, x * bn + z], name='pB')
8   C = tvm.compute((M,N), lambda x,y:
9     tvm.sum(A[x,k] * pB[y//bn,k,
10    tvm.indexmod(y,bn)], axis=k),name='C')
11  # Array packing schedule
12  s = tvm.create_schedule(C.op)
13  xo, yo, xi, yi = s[C].tile(
14    C.op.axis[0], C.op.axis[1], bn, bn)
15  k,              = s[C].op.reduce_axis
16  ko, ki = s[C].split(k, factor=4)
17  s[C].reorder(xo, yo, ko, xi, ki, yi)
18  s[C].vectorize(yi)
19  x, y, z          = s[pB].op.axis
20  s[pB].vectorize(z)
21  s[pB].parallel(x)
```

*Array Packing* Strategy

# Optimizing Matrix Matrix Multiplication with ELEVATE Strategies

Clear separation of concerns     vs     No clear separation of concerns

## ELEVATE

```
1   val appliedMap = isApp(isApp(isMap))
2   val isTransposedB = isApp(isTranspose)
3
4   val packB = storeInMemory(isTransposedB,
5    permuteB ';;'
6    vectorize(32) '@' innermost(appliedMap) ';;'
7    parallel        '@' outermost(isMap)
8   ) '@' inLambda
9
10  val arrayPacking = packB ';;' loopPerm
11  (arrayPacking ';' lowerToC )(mm)
```

Facilitate reuse

```
1   # Modified algorithm
2   bn = 32
3   k = tvm.reduce_axis((0, K), 'k')
4   A = tvm.placeholder((M, K), name='A')
5   B = tvm.placeholder((K, N), name='B')
6   pB = tvm.compute((N / bn, K, bn),
7     lambda x, y, z: B[y, x * bn + z], name='pB')
8   C = tvm.compute((M,N), lambda x,y:
9     tvm.sum(A[x,k] * pB[y//bn,k,
10    tvm.indexmod(y,bn)], axis=k),name='C')
11  # Array packing schedule
12  s = tvm.create_schedule(C.op)
13  xo, yo, xi, yi = s[C].tile(
14    C.op.axis[0], C.op.axis[1], bn, bn)
15  k,           = s[C].op.reduce_axis
16  ko, ki = s[C].split(k, factor=4)
17  s[C].reorder(xo, yo, ko, xi, ki, yi)
18  s[C].vectorize(yi)
19  x, y, z      = s[pB].op.axis
20  s[pB].vectorize(z)
21  s[pB].parallel(x)
```

*Array Packing* Strategy

# Optimizing Matrix Matrix Multiplication with ELEVATE **Strategies**
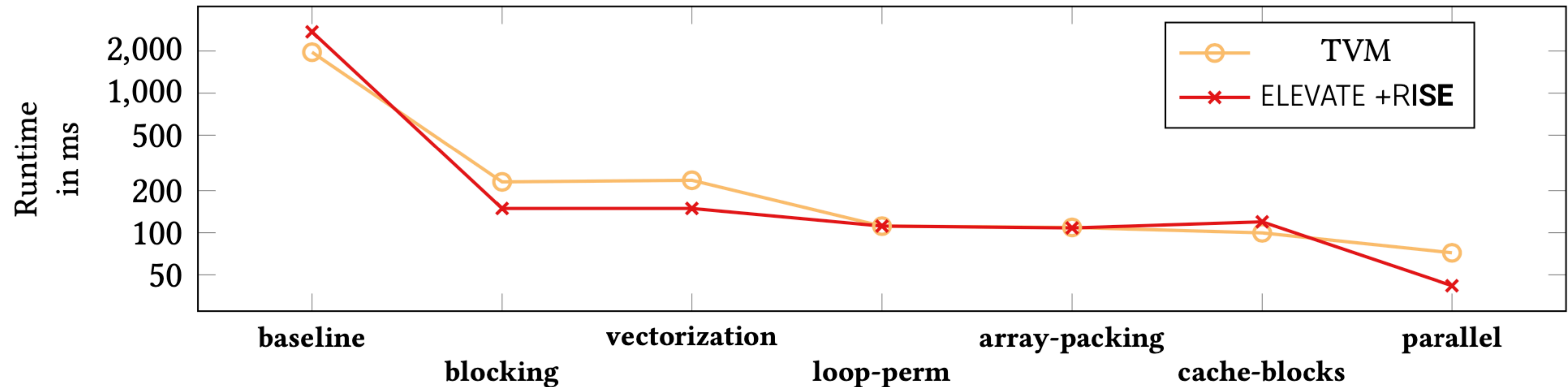
Number of successful rewrite steps



Rewriting took less than 2 seconds with our unoptimised implementation

**Rewrite based approach scales to complex optimizations**

# Optimizing Matrix Matrix Multiplication with ELEVATE **Strategies**

Performance of generated code



**Competitive performance compared to TVM compiler**

# Types for ELEVATE?

**?** Can we build a type system for ELEVATE
to statically reject bad compositions
of rewrites?

Ongoing work using row-polymorphic
types for this.

Preliminary result in an arXiv paper:
https://arxiv.org/abs/2103.13390

---

## Row-Polymorphic Types for Strategic Rewriting

Rongxiao Fu
University of Glasgow
United Kingdom
rongxiao.fu@glasgow.ac.uk

Xueying Qin
The University of Edinburgh
United Kingdom
xueying.qin@ed.ac.uk

Ornela Dardha
University of Glasgow
United Kingdom
ornela.dardha@glasgow.ac.uk

Michel Steuwer
The University of Edinburgh
United Kingdom
michel.steuwer@ed.ac.uk

**Abstract**

We present a type system for *strategy languages* that express program transformations as compositions of rewrite rules. Our row-polymorphic type system assists compiler engineers to write correct strategies by statically rejecting non meaningful compositions of rewrites that otherwise would fail during rewriting at runtime. Furthermore, our type system enables reasoning about how rewriting transforms the shape of the computational program. We present a formalization of our language at its type system and demonstrate its practical use for expressing compiler optimization strategies.

Our type system builds the foundation for many interesting future applications, including verifying the correctness of program transformations and synthesizing program transformations from specifications encoded as types.

Hagedorn et al. [15] describe how the ELEVATE st
language is used to encode and control the application
ditional compiler optimizations such as loop-tiling a
ing performance comparable to the traditionally de
TVM compiler [5] for deep learning. This picks up the
of increased importance of efficiency in many appli
domains of today and the future. For example, the
through success of deep learning has only been p
thanks to carefully optimized software making efficie
of modern parallel hardware. In the TVM compile
mization decisions are encoded in a so-called *schedule*
performance engineers select from a fixed set of ex
compiler transformations to optimize their deep le
application. In ELEVATE, strategic rewriting gives de
ers even greater flexibility as they are free to encode
program transformations – possibly domain- or har
specific – as strategies and precisely control their a
tion.

However, developing strategies that encode mean
program transformations is not easy. One reason is th
rent strategy languages provide little to no support for

## 1 Introduction

Rewrite systems find applications in many domains ranging from logic [26] and theorem provers [17] to program trans-

# Sketch-Guided Equality Saturation

## Automation vs. Manual control

A *sketch* describes
a desired program shape

```
containsMap(m / 32,            |  for m / 32:
 containsMap(n / 32,           |   for n / 32:
  containsReduceSeq(k / 4,     |    for k / 4:
   containsReduceSeq(4,        |     for 4:
    containsMap(32,            |      for 32:
     containsMap(32,           |       for 32:
      containsAddMul))))))     |        .. + .. × ..
```

💡 **Idea**:
 Describe rewrite *goal* rather than rewrite sequence:



Break intractable equality saturation search
into multi tractable one, by human *guidance*.

All optimizations from this paper
are found in  < 7 seconds *automatically*

Talk by **Thomas Kœhler** earlier this week at the E-Graph workshop. Paper: https://arxiv.org/abs/2111.13040

# Achieving High-Performance the Functional Way

## Expressing High-Performance Optimisations as Rewrite Strategies

Bastian Hagedorn, Johannes Lenfers, Thomas Kœhler, Xueying Qin, Sergei Gorlatch and Michel Steuwer



WWU MÜNSTER  THE UNIVERSITY of EDINBURGH  University of Glasgow

https://rise-lang.org/

michel.steuwer@ed.ac.uk

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

https://michel.steuwer.info/

https://github.com/elevate-lang/elevate