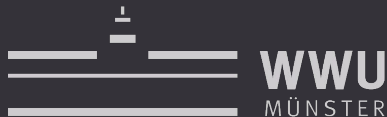


ACHIEVING HIGH-PERFORMANCE THE *Functional* WAY

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



THE UNIVERSITY
of EDINBURGH



University
of Glasgow

Bastian Hagedorn | Johannes Lenfers | Thomas Koehler | Xueying Qin | Sergei Gorlatch | Michel Steuwer

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 CO₂ into the atmosphere (equivalent to driving a car around 11,000 miles)

HIGH-PERFORMANCE

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

```
__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



HIGH-PERFO

Achieving High-Performance

```
_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



Optimized Matrix Multiplication

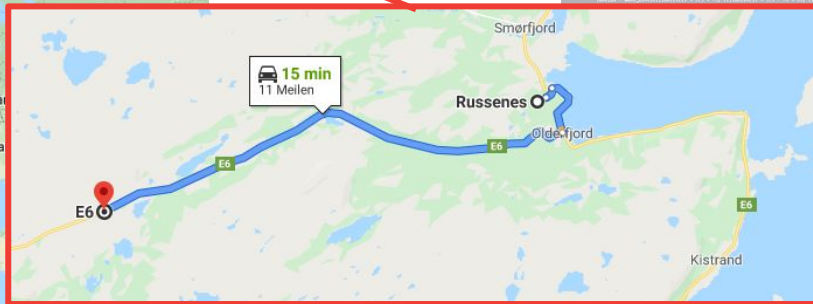
PERFO

High Performance

1000x CO2 Improvement

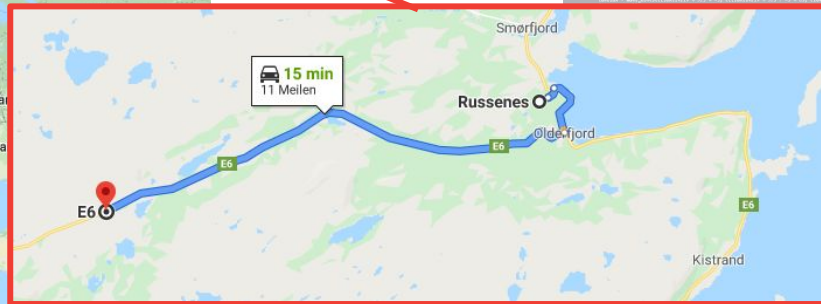
100-1000x performance

Optimized Matrix Multiplication



PERFO

High Performance



1000x CO2 Improvement

100-1000x performance

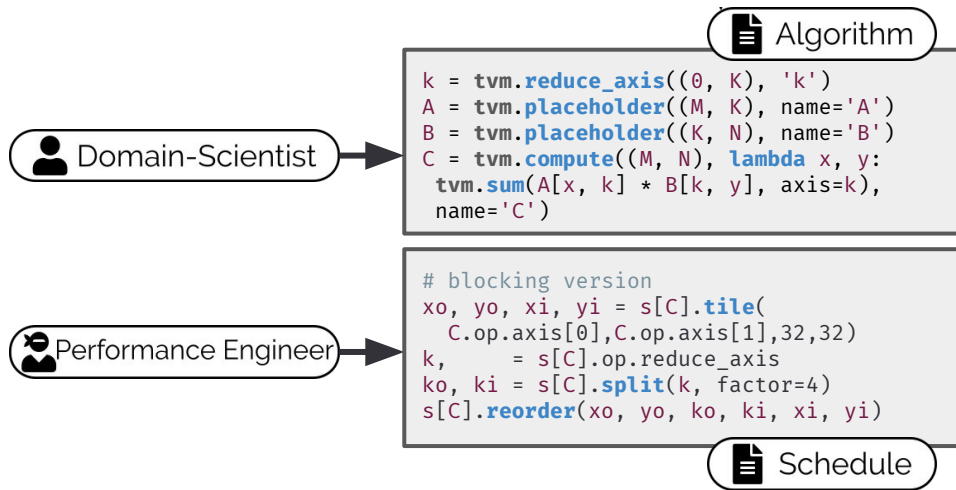
30x lines of code

time-intensive + error-prone

Optimized Matrix Multiplication

HIGH-PERFORMANCE

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



← *What to compute*

← *How to optimize*

Halide



Tiramisu-Compiler / [tiramisu](#)

Fireiron NVIDIA

HIGH-PERFORMANCE

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

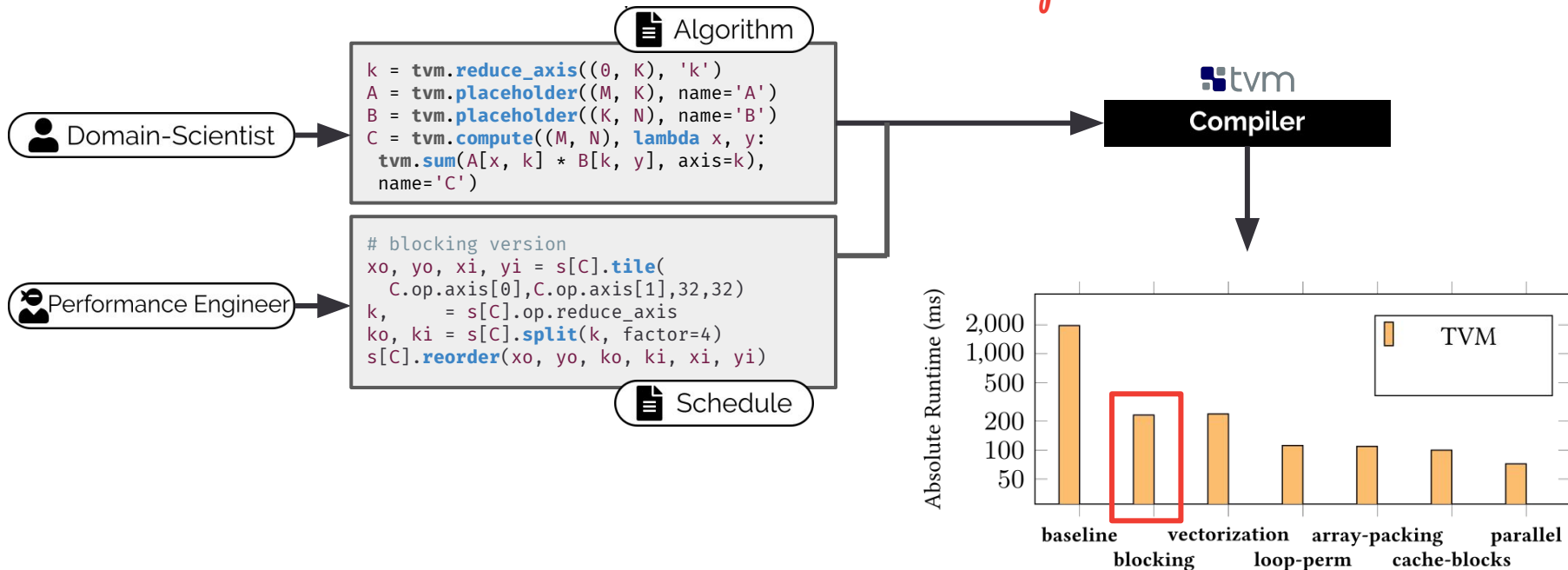
Algorithm

The screenshot shows a web browser displaying the TVM documentation page for "How to optimize GEMM on CPU". The browser's address bar shows the URL `https://tvm.apache.org/docs/tutorials/optimize/opt_gemm.html`. The page features a sidebar with navigation links under the TVM logo (0.7.dev1). The main content area includes a "Note" box with a link to download example code, the article title "How to optimize GEMM on CPU" by Jian Weng and Ruofei Yu, a TL;DR summary, and an introduction paragraph. The sidebar contains the following links:

- HOW TO
 - Installation
 - Contribute to TVM
 - Deploy and Integration
 - Developer How-To Guide
- TUTORIALS
 - Get Started Tutorials
 - Compile Deep Learning Models
 - Tensor Expression and Schedules
- Optimize Tensor Operators
 - How to optimize convolution on

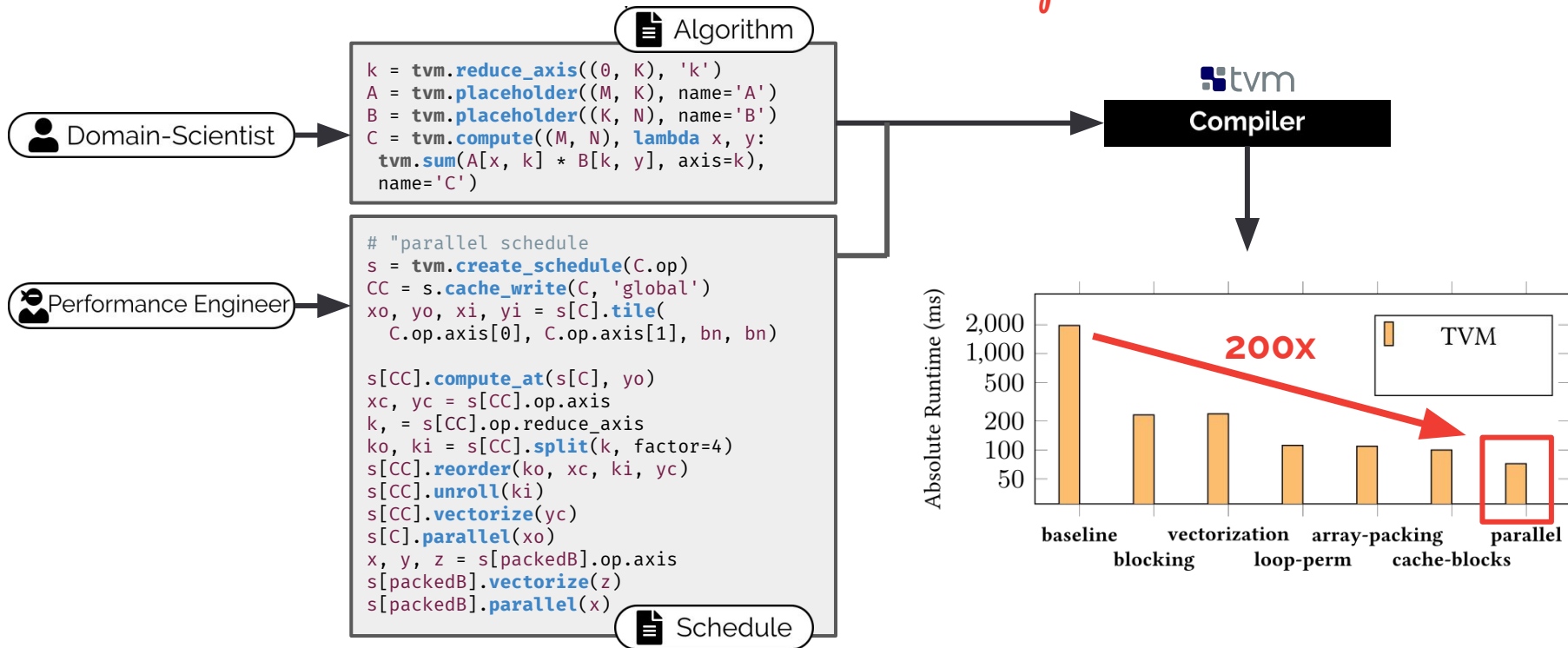
HIGH-PERFORMANCE

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



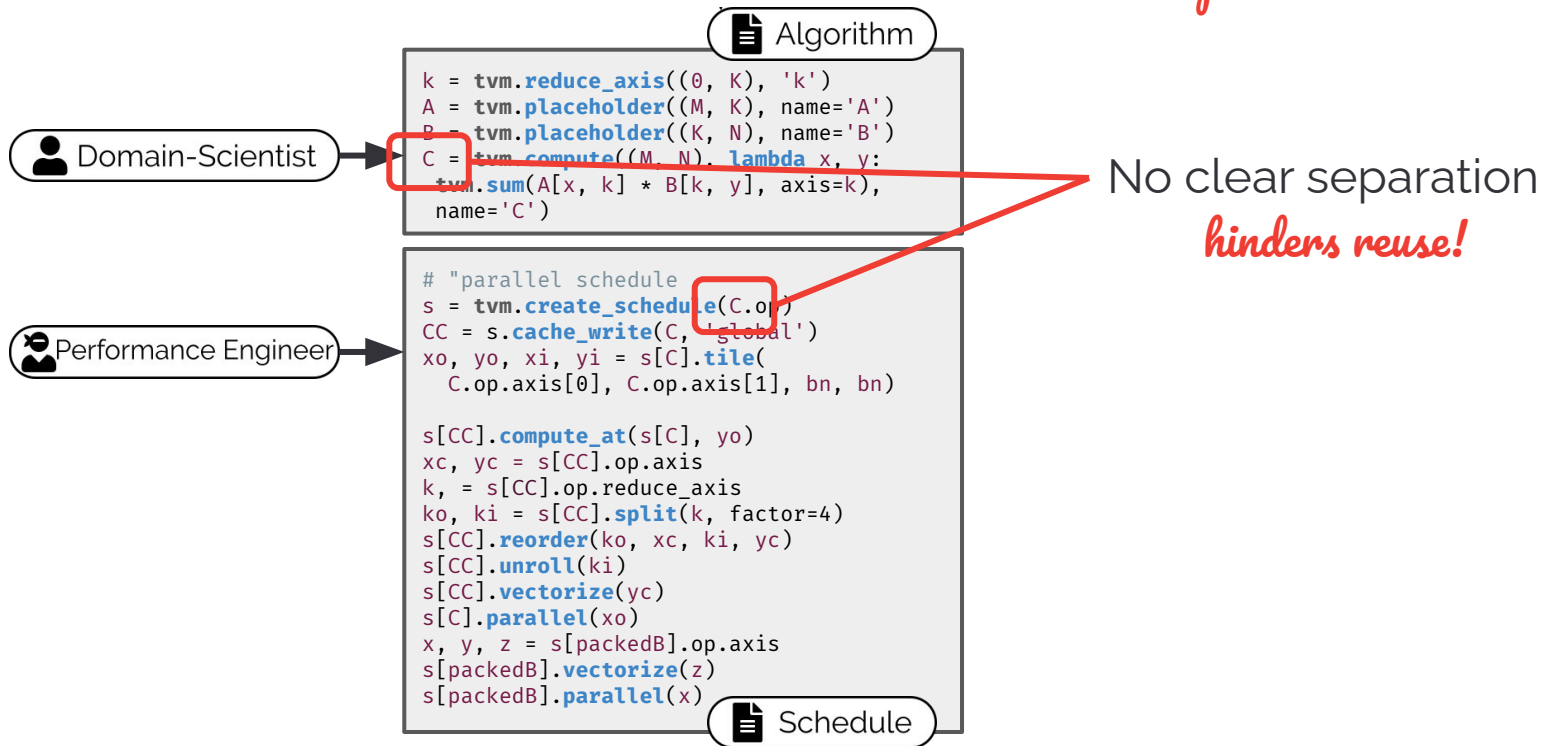
HIGH-PERFORMANCE

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



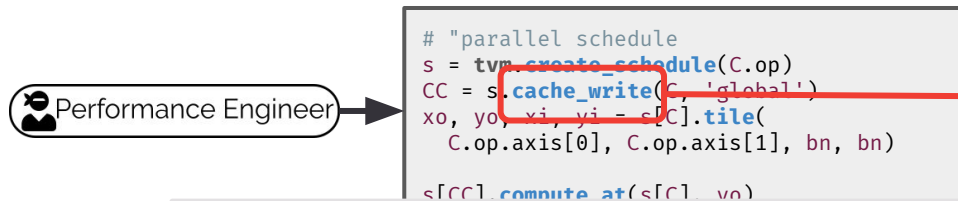
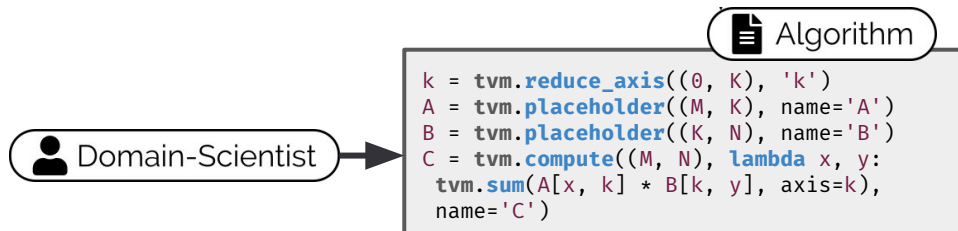
HIGH-PERFORMANCE

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



HIGH-PERFORMANCE

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



No clear separation

hinders reuse!

No well-defined semantics

hinders understanding!

`cache_write(tensor, scope)`

?

Create a cache write of original tensor, before storing into tensor.

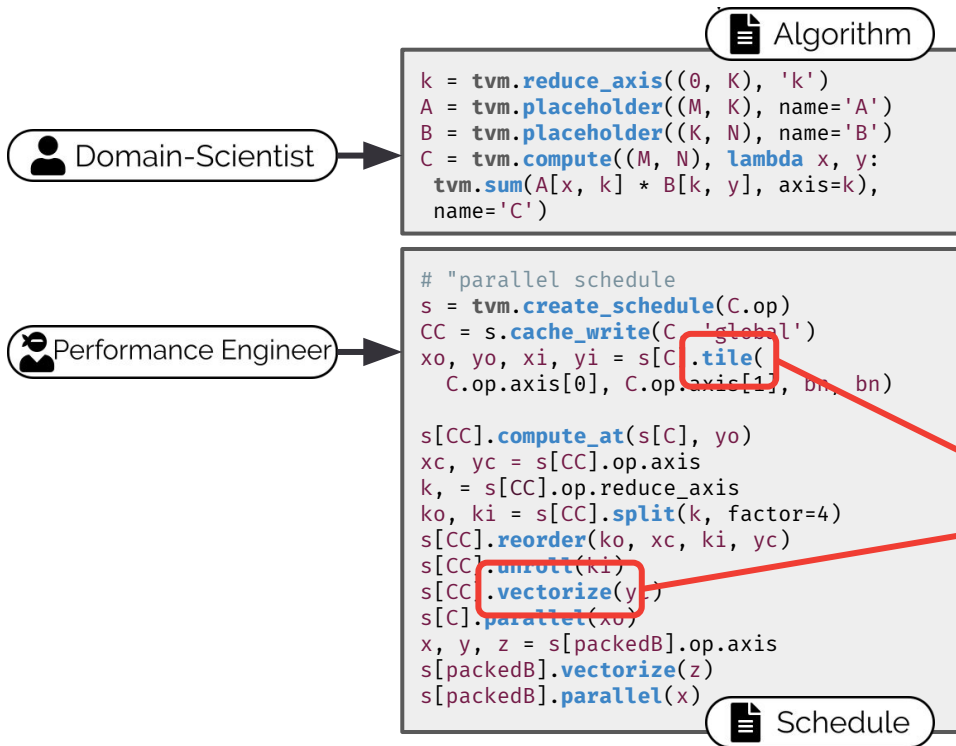
```
s[C].parallel(xo)
x, y, z = s[packedB].op.axis
s[packedB].vectorize(z)
s[packedB].parallel(x)
```



Schedule

HIGH-PERFORMANCE

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



No clear separation

hinders reuse!

No well-defined semantics

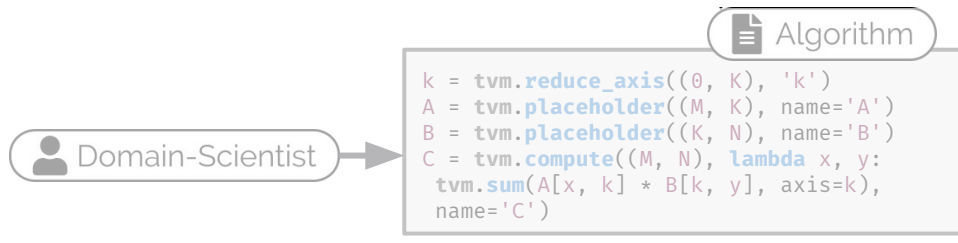
hinders understanding!

Optimizations are built-in

no extensibility!

HIGH-PERFORMANCE

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



No clear separation

hinders reuse!



We aim for a *more principled* way to *describe and apply optimizations*

```
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(y)
s[C].parallel(xo)
x, y, z = s[packedB].op.axis
s[packedB].vectorize(z)
s[packedB].parallel(x)
```

Schedule

Optimizations are built-in

no extensibility!

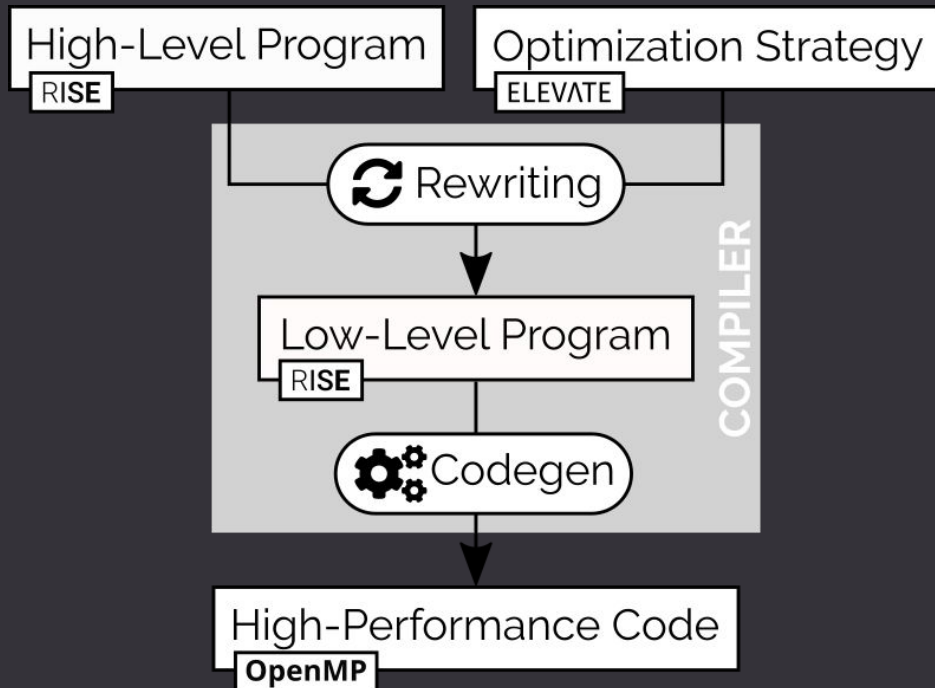
OUR GOALS

A Principled Way to Separate, Describe, and Apply Optimizations

- (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.

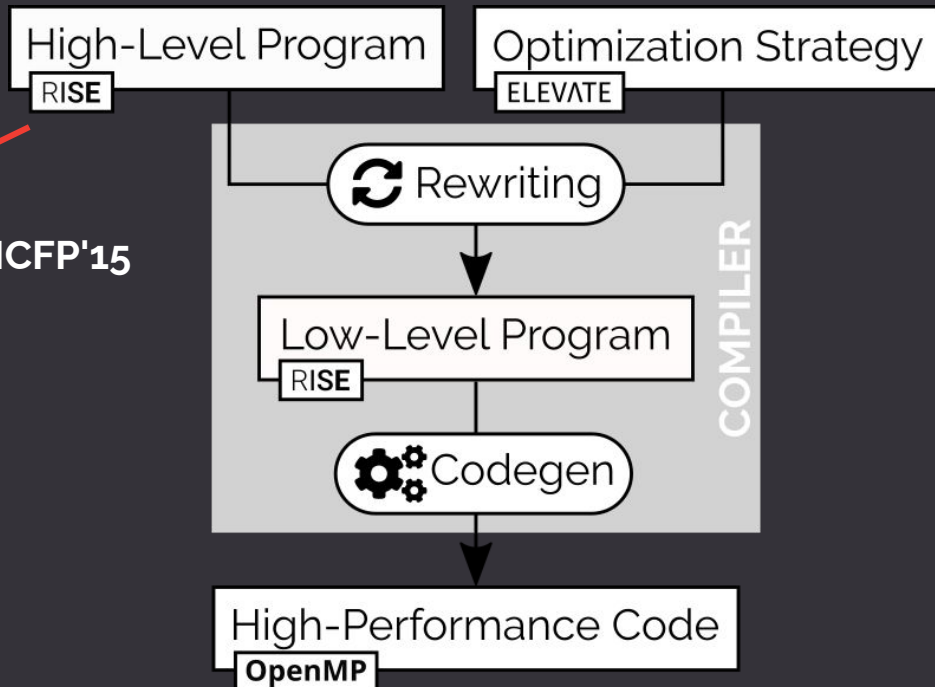
Fundamentally we argue that a more principled high-performance code generation approach should be holistic by considering computation and optimization strategies equally important. As a consequence, a strategy language should be built with the same standards as a language describing computation.

The *Functional* Way

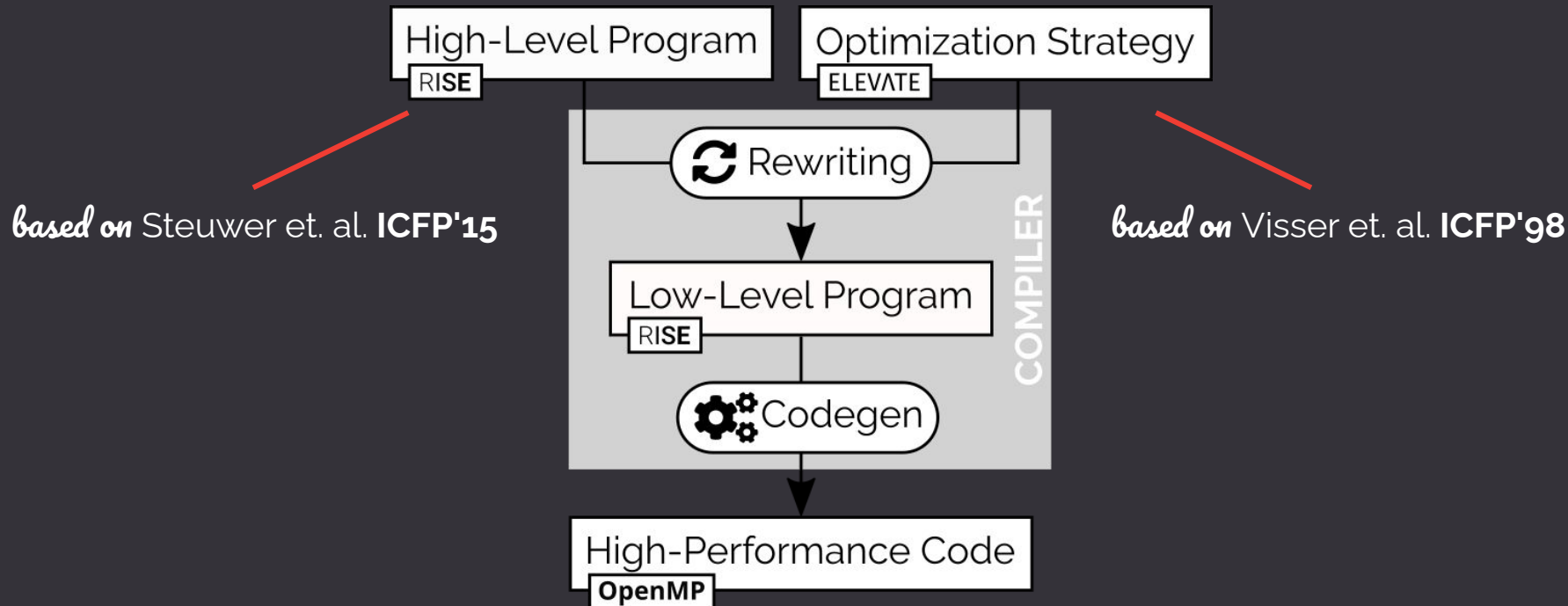


The *Functional* Way

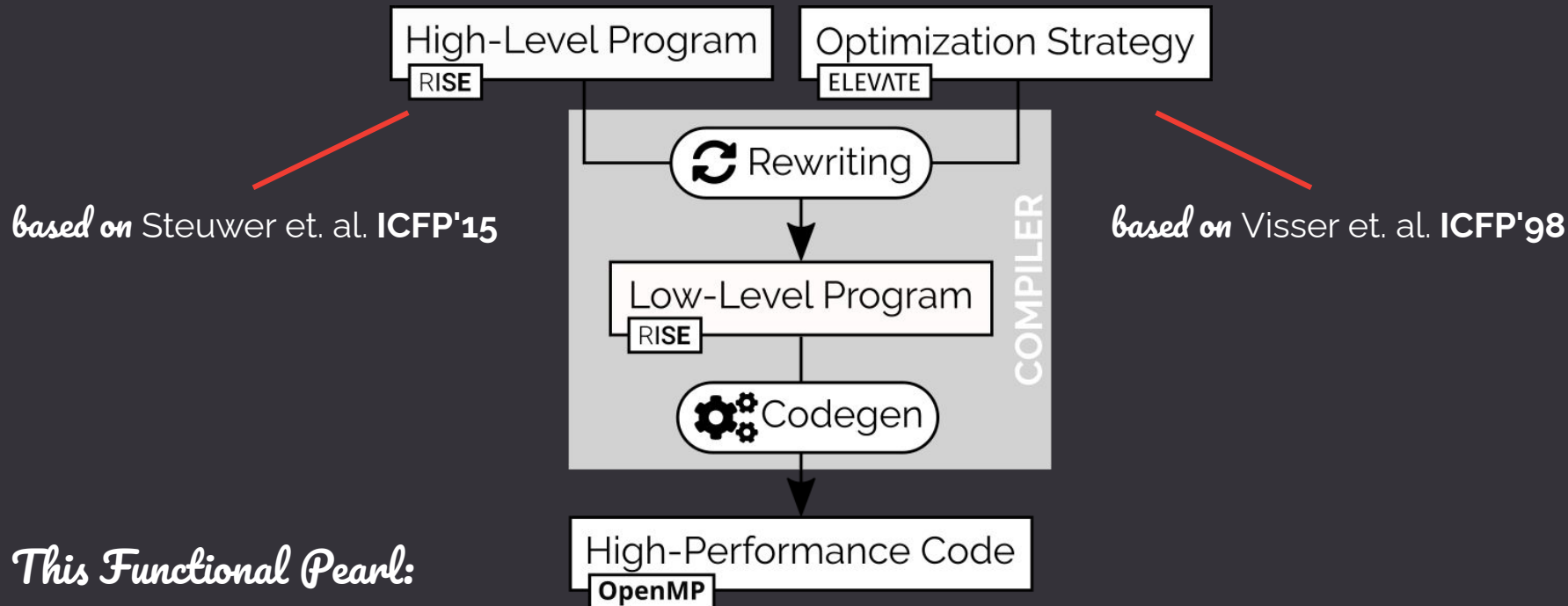
based on Steuwer et. al. ICFP'15



The *Functional* Way



The *Functional* Way



This Functional Pearl:

We apply established *functional programming techniques* for elegantly expressing high-performance program optimizations as composable rewrite strategies

ELEVATE

A Language for Describing Optimization Strategies

A *Strategy* encodes a program transformation:

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

A *RewriteResult* encodes its success or failure:

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

ELEVATE

A Language for Describing Optimization Strategies

A *Strategy* encodes a program transformation:

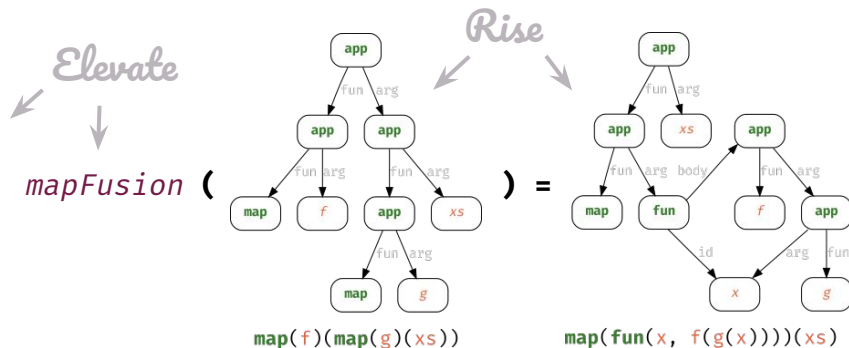
```
type Strategy[P] = P => RewriteResult[P]
```

A *RewriteResult* encodes its success or failure:

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

Rewrite Rules are examples for basic strategies

```
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 )  
  }
```



COMBINATORS

How to Build More Powerful Strategies

Sequential Composition (;)

```
def seq[P]: Strategy[P] => Strategy[P] => Strategy[P] =  
  fs => ss => p => fs(p) >>= ss
```

Left Choice (<+)

```
def lChoice[P]: Strategy[P] => Strategy[P] => Strategy[P] =  
  fs => ss => p => fs(p) <|> 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)
```

Describing Precise Locations



There are *two possible locations* for successfully applying the rule

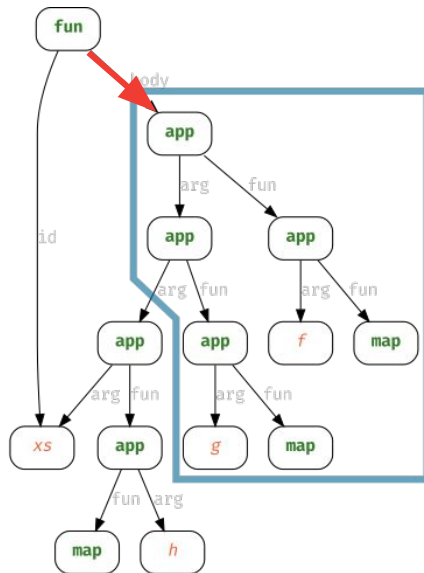
TRAVERSALS

Describing Precise Locations

```
def body: Traversal[Rise] = s => p => p match {  
  case fun(x,b) => (nb => fun(x,nb) <$> s(b))  
  case _ => Failure( body(s) )  
}
```

apply s at $body$ of function abstraction

$body(mapFusion)$ (



threemaps = $\text{fun}(xs, \text{map}(f)(\text{map}(g)(\text{map}(h)(xs))))$

There are *two possible locations* for successfully applying the rule

TRAVERSALS

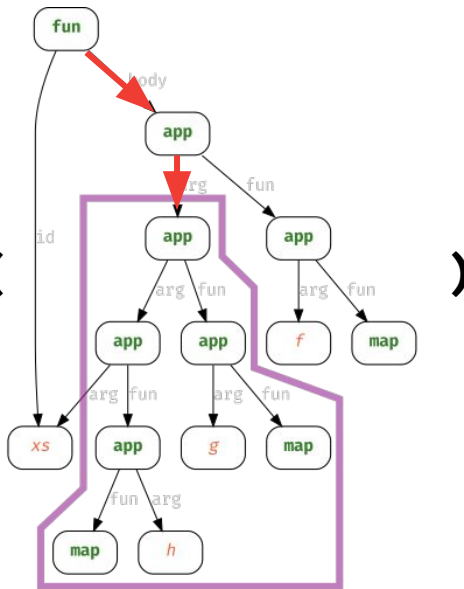
Describing Precise Locations

```
def body: Traversal[Rise] = s => p => p match {  
  case fun(x,b) => (nb => fun(x,nb) <$> s(b))  
  case _ => Failure( body(s) )  
}
```

body(argument(mapFusion)) (

```
def argument: Traversal[Rise] = s => p => p match {  
  case app(f,a) => (na => app(f,na) <$> s(a))  
  case _ => Failure( argument(s) )  
}
```

apply s at argument of function application



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

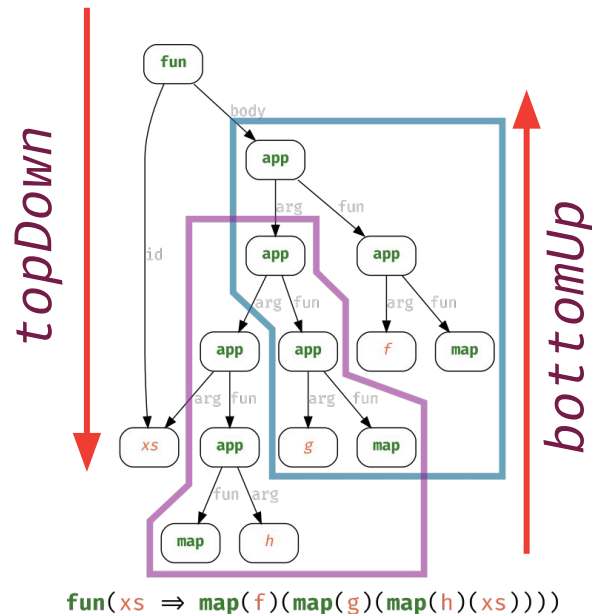
There are *two possible locations* for successfully applying the rule

NORMALIZATION

More Complex Traversals

Generic Tree Traversals...

```
def topDown: Traversal[Rise] = s => p => (s <+ one(topDown(s)))(p)
def bottomUp: Traversal[Rise] = s => p => (one(topDown(s)) <+ s)(p)
...
```



NORMALIZATION

More Complex Traversals

Generic Tree Traversals...

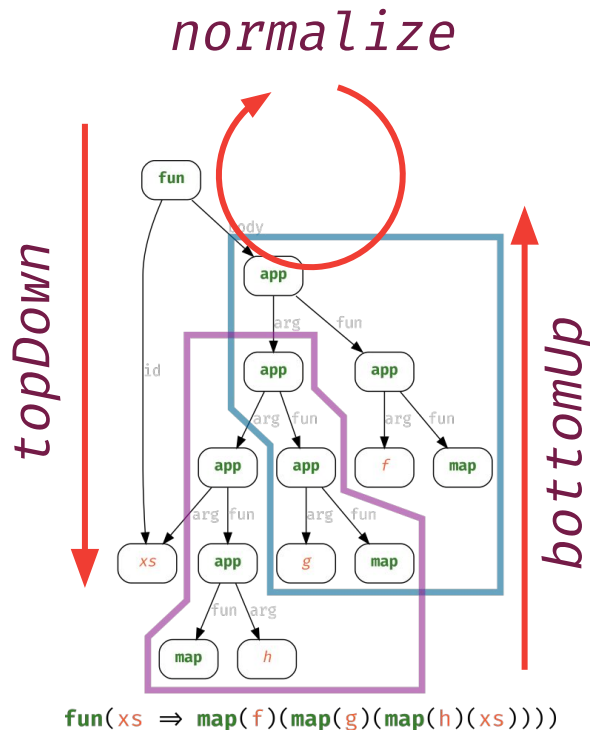
```
def topDown: Traversal[Rise] = s => p => (s <+ one(topDown(s)))(p)
def bottomUp: Traversal[Rise] = s => p => (one(topDown(s)) <+ s)(p)
...
```

... and a strategy for normalization

```
def normalize: Traversal[Rise] = s => p => repeat(topDown(s))(p)
```

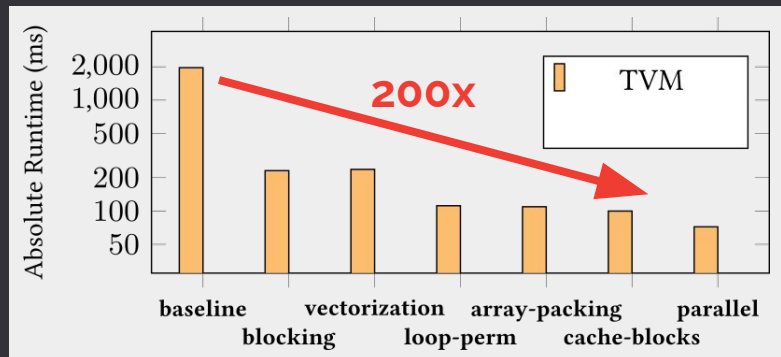
With these, we define normal-forms like $\beta\eta$ -normal-form

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



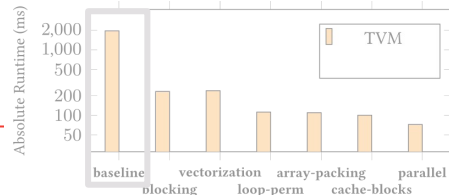
CASE STUDY

Implementing TVM's Scheduling Language



CASE STUDY

Optimizing Matrix Multiplication - Baseline



RISE

What to compute



```
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(arrow, transpose(b) |>
7     map( fun(bcol,
8       dot(arrow)(bcol) ) ) ) ) ) )
```

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

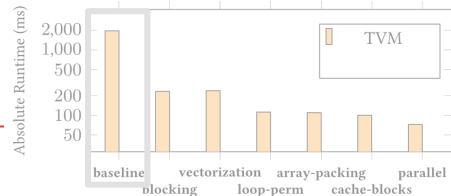
ELEVATE

```
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)
```

How to optimize

CASE STUDY

Optimizing Matrix Multiplication - Baseline



clear separation

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)(@) ) )
5 val mm = fun(a, fun(b, a |>
6   map( fun(arrow, transpose(b) |>
7     map( fun(bcol,
8       dot(arrow)(bcol) )))) ) )
```

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

ELEVATE

composable

explicit

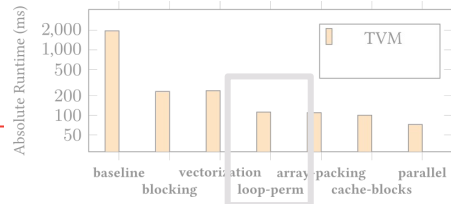


```
1 # Naive matrix multiplication algorithm
2 k = tvm.reduce_axis((0, K), 'k')
3 A = tvm.placeholder((M, K), name='A')
4 B = tvm.placeholder((K, N), name='B')
5 C = tvm.compute((M, N), lambda x, y:
6   tvm.sum(A[x, k] * B[k, y],
7   axis=k), name='C')
8
9
10
11
12 # TVM default schedule
13 s = tvm.create_schedule(C.op)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

implicit

CASE STUDY

Optimizing Matrix Multiplication - Loop Permutation



facilitate reuse

user-defined vs. built-in

```
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)
```

ELEVATE

```
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)
```

tvm

no clear separation of concerns

CASE STUDY

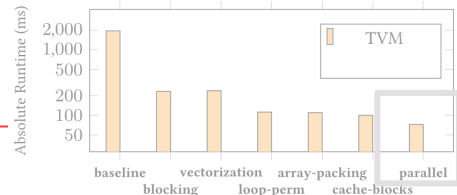
Optimizing Matrix Multiplication - Array Packing

clear separation of concerns vs. no clear separation

facilitate reuse

```
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 par = (
11   packB ';;' loopPerm ';;'
12   (parallel '@' outermost(isMap))
13   '@' outermost(isToMem) ';;'
14   unroll '@' innermost(isReduce))
15
16 (par ';' lowerToC )(mm)
```

```
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 CC = s.cache_write(C, 'global')
14 xo, yo, xi, yi = s[C].tile(
15   C.op.axis[0], C.op.axis[1], bn, bn)
16 s[CC].compute_at(s[C], yo)
17 xc, yc = s[CC].op.axis
18 k, = s[CC].op.reduce_axis
19 ko, ki = s[CC].split(k, factor=4)
20 s[CC].reorder(ko, xc, ki, yc)
21 s[CC].unroll(ki)
22 s[CC].vectorize(yc)
23 s[C].parallel(xo)
24 x, y, z = s[pB].op.axis
25 s[pB].vectorize(z)
26 s[pB].parallel(x)
```

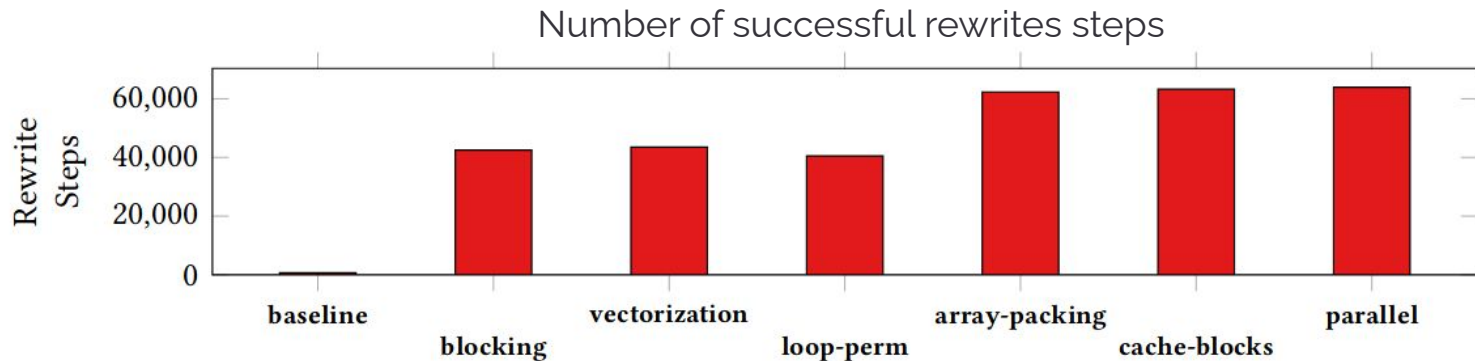


ELEVATE



CASE STUDY

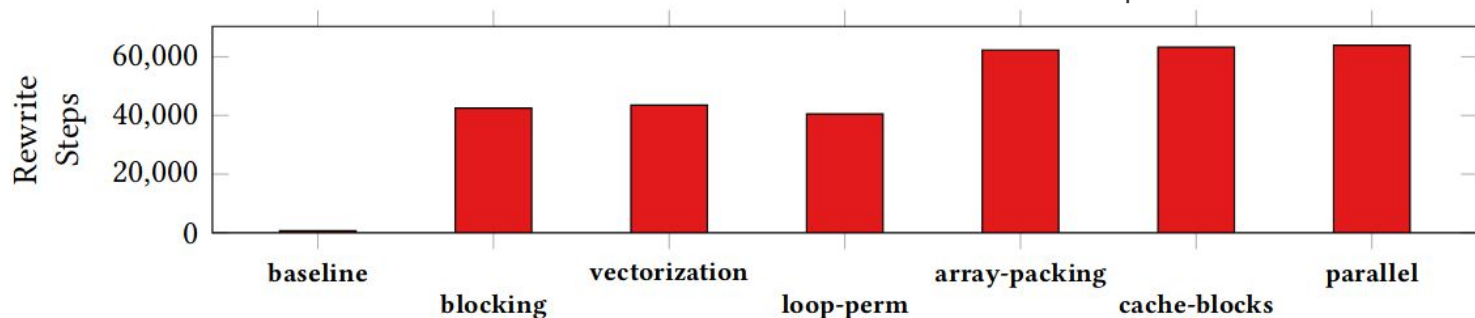
Counting Rewrite Steps and Measuring Performance



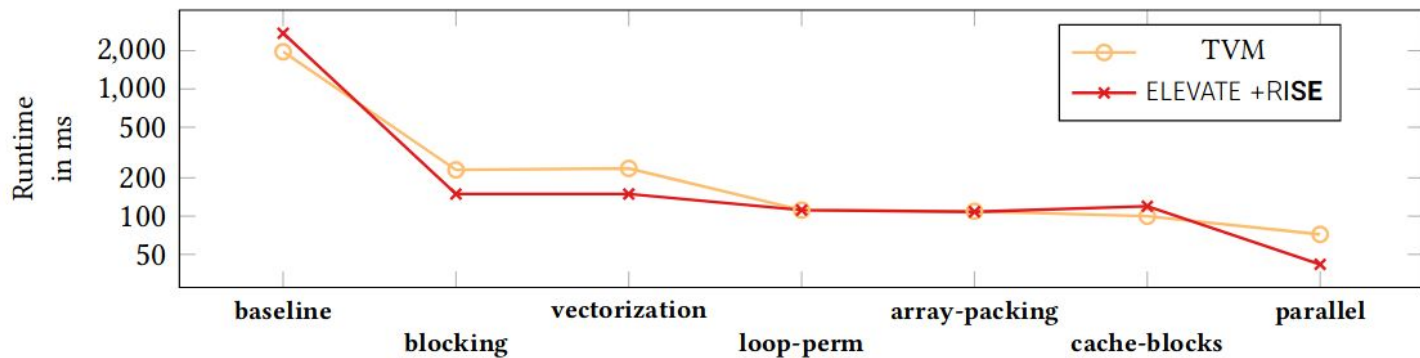
CASE STUDY

Counting Rewrite Steps and Measuring Performance

Number of successful rewrites steps



Performance of the generated code



ACHIEVING HIGH-PERFORMANCE THE *Functional* WAY *...is Open Source!*

RISE

rise-lang.org/
github.com/rise-lang

ELEVATE

elevate-lang.org
github.com/elevate-lang



b.hagedorn@wwu.de
bastianhagedorn.github.io

