

# HOW TO DESIGN THE NEXT 700 OPTIMIZING COMPILERS

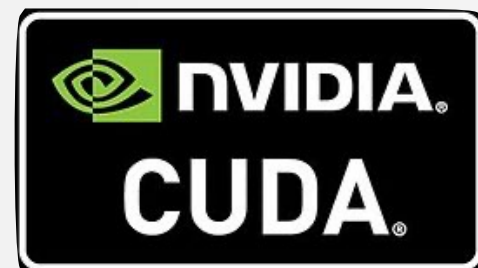
**A framework for designing optimising domain-specific compilers  
for specialised hardware in the era of ML and AI**

Michel Steuwer



THE UNIVERSITY *of* EDINBURGH

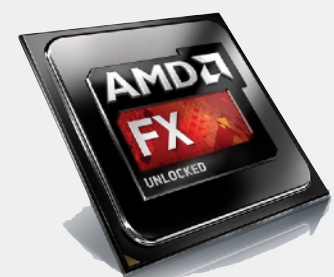
General purpose



# Software



Specialised



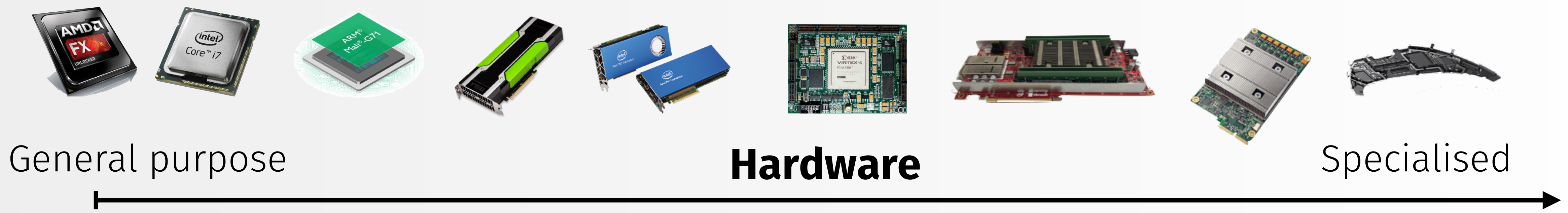
General purpose

# Hardware

Specialised



**How do we build compilers to (automatically) optimise specialised software for specialised hardware?**



*COMPUTATION*

*OPTIMISATION*

## Domain Specific Example: TensorFlow

- > 500 different type of nodes in the TF IR
- > 50 different type of nodes in the XLA IR
- > 2.500.000 lines of code
- Support for custom hardware: TPU
- **Hughe effort to build still highly specialised**
- Problem solved?



# Machine Learning Systems are Stuck in a Rut

Paul Barham  
Google Brain

Michael Isard  
Google Brain

[HotOS'19]

## Abstract

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

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

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

### ACM Reference Format:

Paul Barham and Michael Isard. 2019. Machine Learning Systems are Stuck in a Rut. In *Workshop on Hot Topics in Operating Systems (HotOS '19)*, May 13–15, 2019, Bertinoro, Italy. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3317550.3321441>

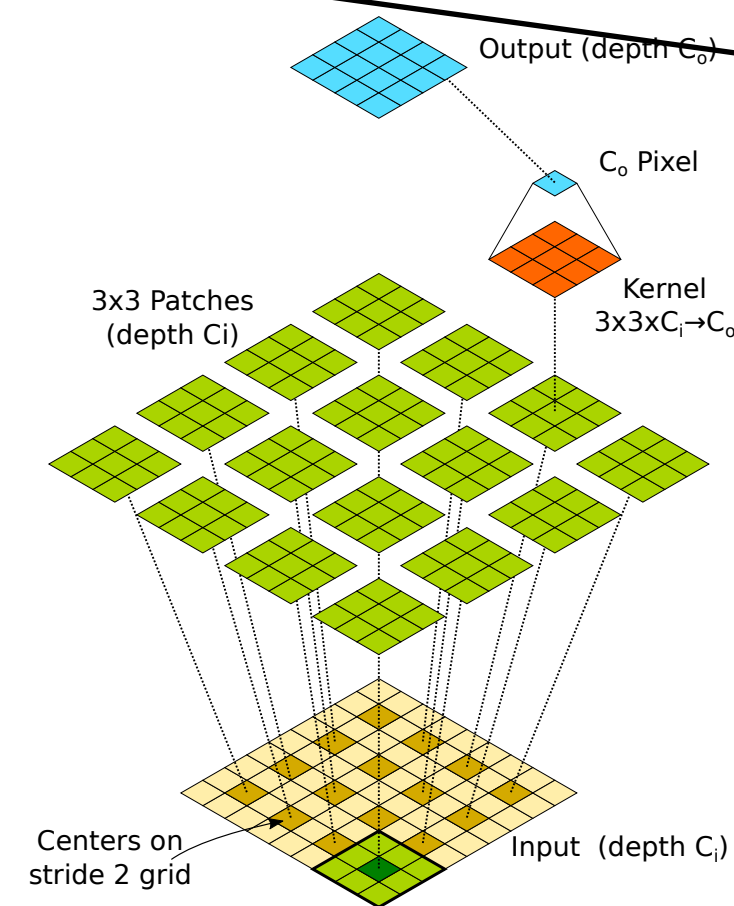


Figure 1. Conv2D operation with  $3 \times 3$  kernel, stride=2

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

### 1.1 New ideas often require new primitives

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

Original authors  
of TensorFlow



# Machine Learning Systems are Stuck in a Rut

Paul Barham  
Google Brain

Michael Isard  
Google Brain

## Abstract

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

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

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

### ACM Reference Format:

Paul Barham and Michael Isard. 2019. Machine Learning Systems are Stuck in a Rut. In *Workshop on Hot Topics in Operating Systems (HotOS '19)*, May 13–15, 2019, Bertinoro, Italy. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3317550.3321441>

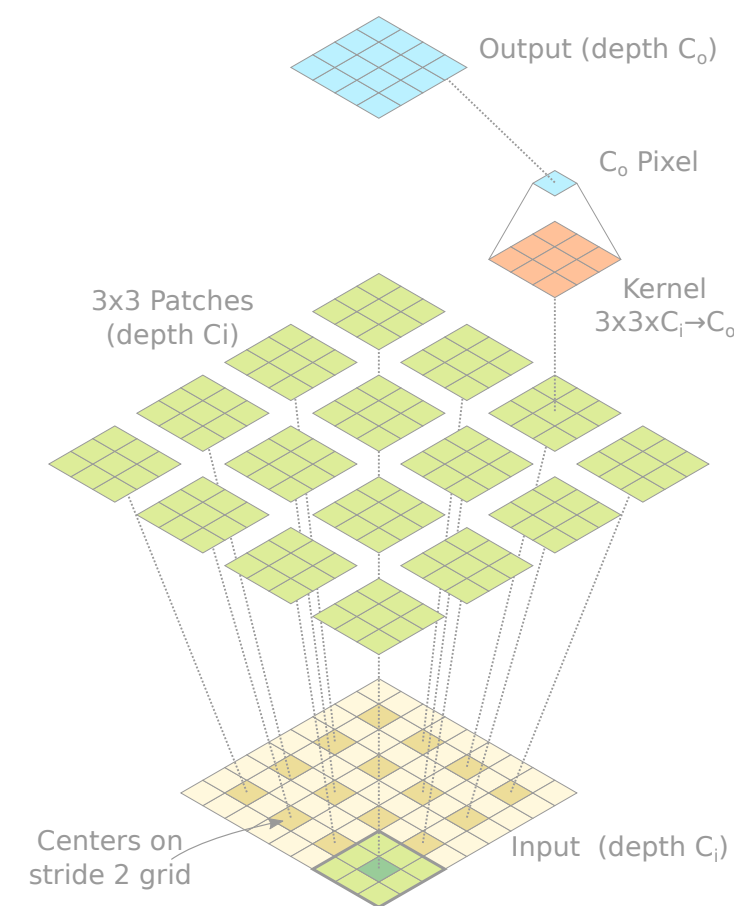


Figure 1. Conv2D operation with 3×3 kernel, stride=2

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

### 1.1 New ideas often require new primitives

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

# Machine Learning Systems are Stuck in a Rut

Paul Barham  
Google Brain

Michael Isard  
Google Brain

## Abstract

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

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

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

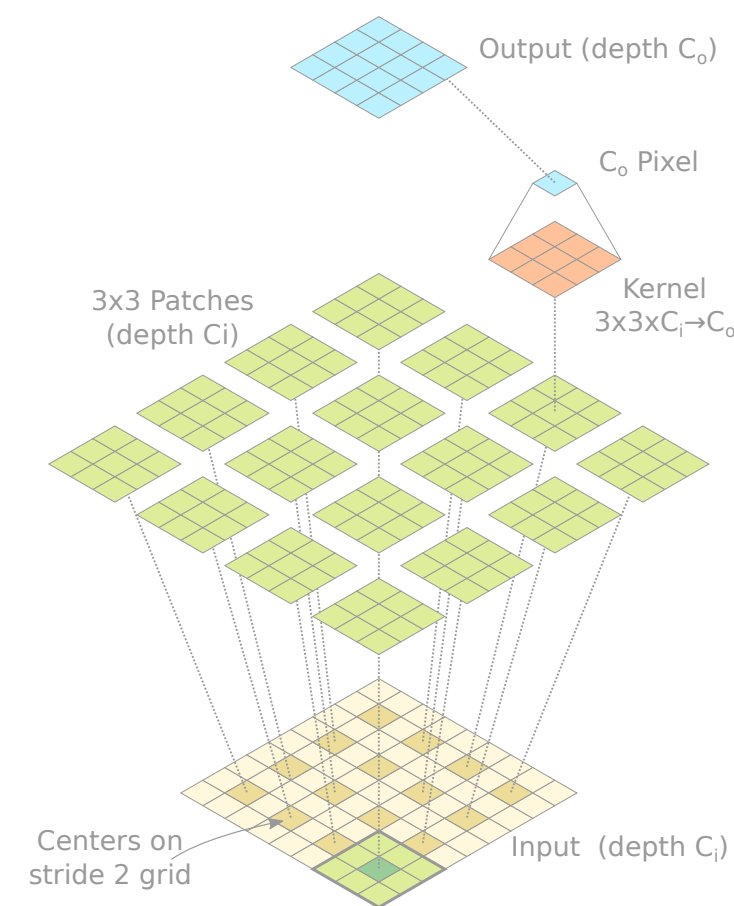


Figure 1. Conv2D operation with 3x3 kernel, stride=2 with 16 times fewer training parameters than the convo-

**We should aim for more principled higher level intermediate representations**

Paul Barham and Michael Isard. 2019. Machine Learning Systems are Stuck in a Rut. In *Workshop on Hot Topics in Operating Systems (HotOS '19)*, May 13–15, 2019, Bertinoro, Italy. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3317550.3321441>

## 1.1 New ideas often require new primitives

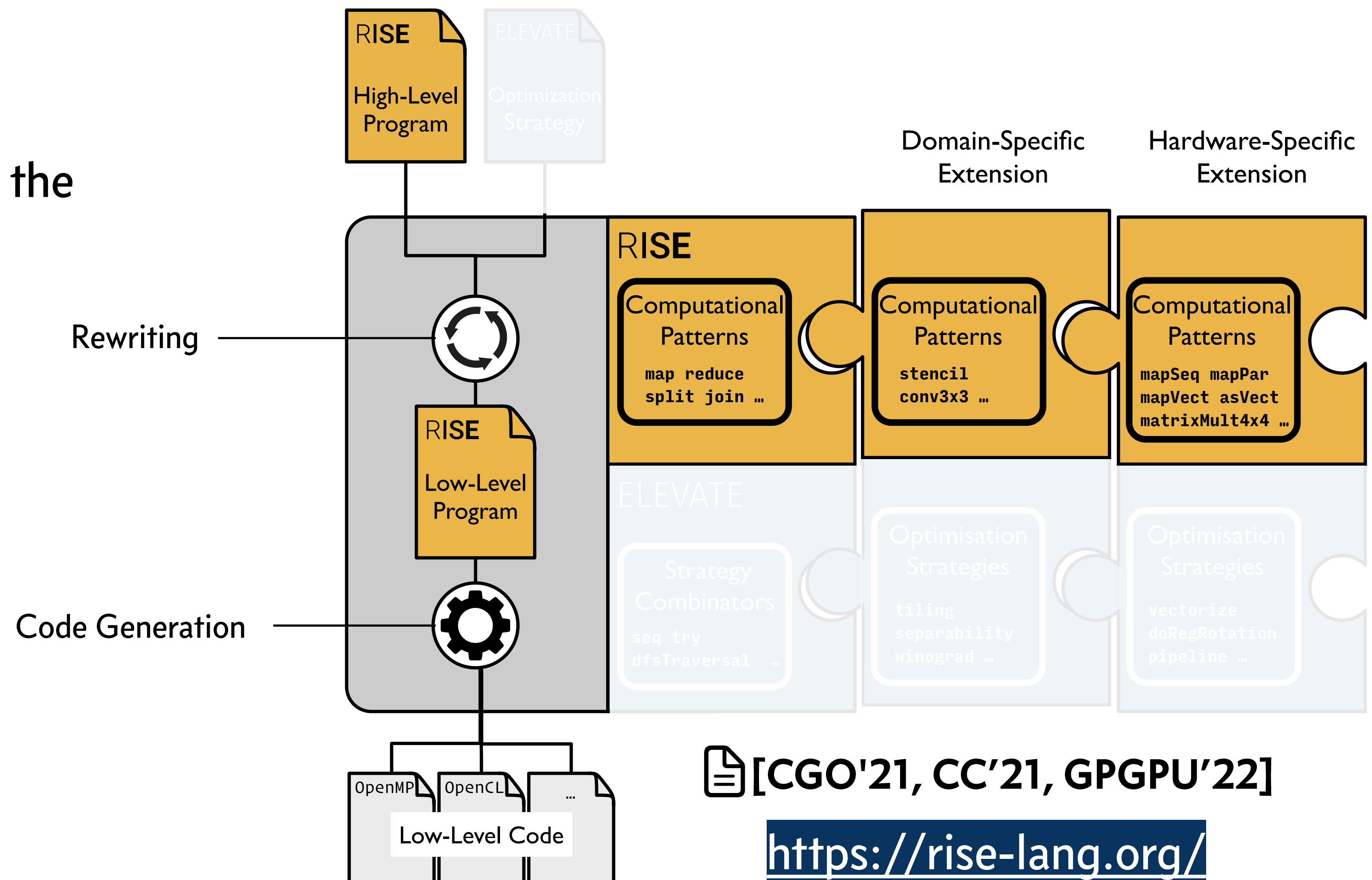
We won't discuss the full details of Capsule networks in this paper<sup>1</sup>, but for our purposes it is sufficient to consider a simplified form of the inner loop, which is



*COMPUTATION*

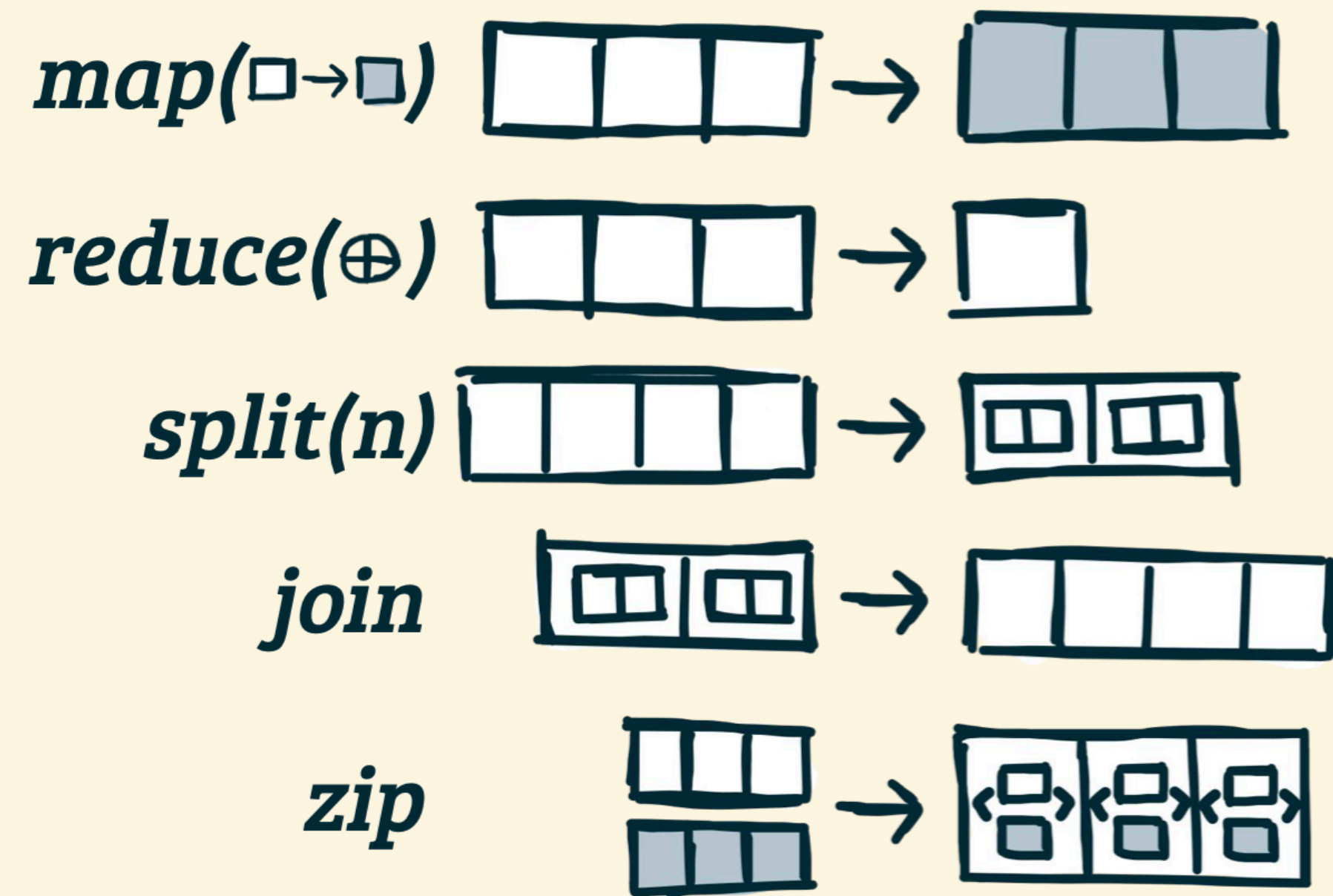
# RISE & Shine an extensible compiler design

- Spiritual successor to the LIFT project
- Functional language as foundation
- Computations are expressed by computational patterns



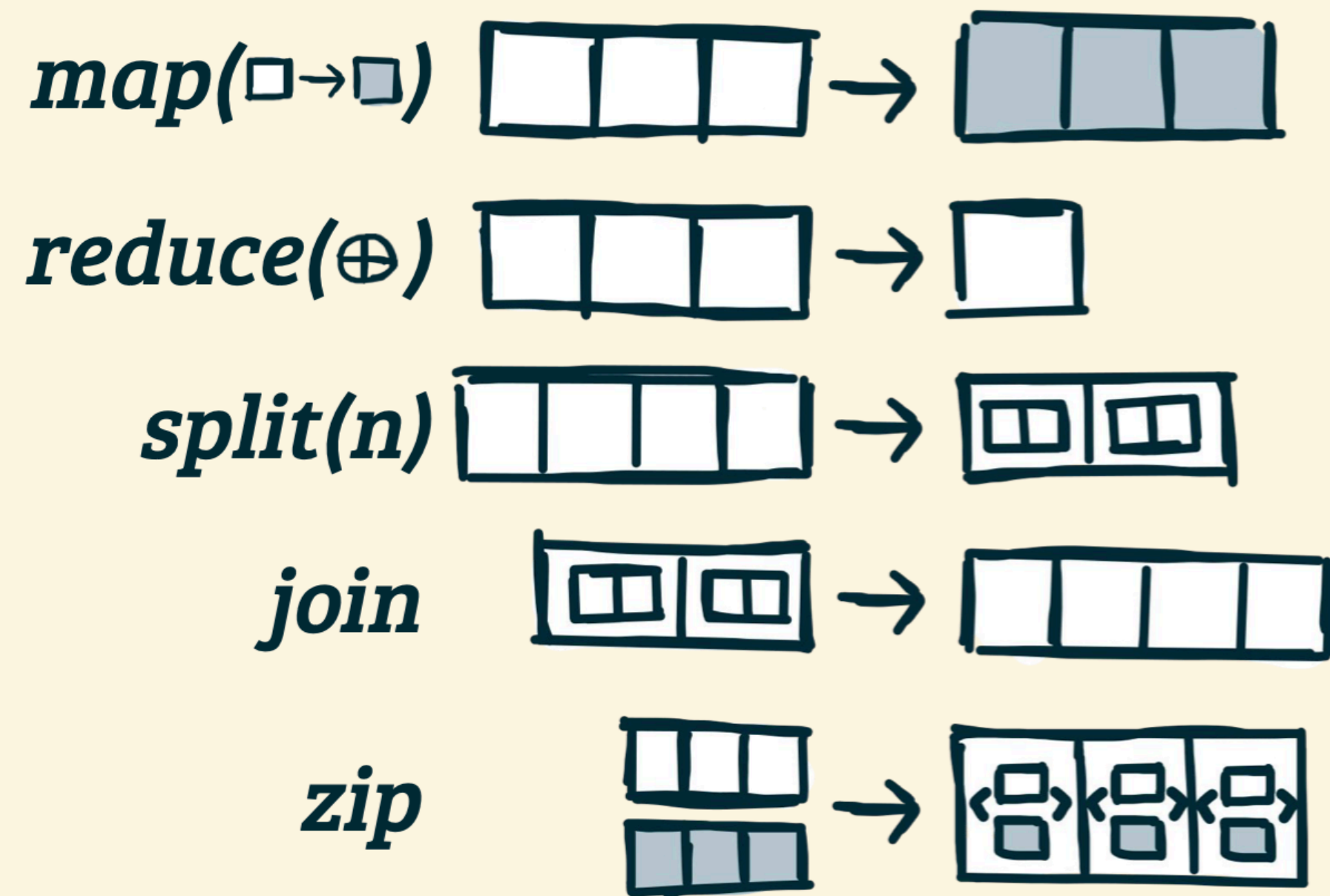
# Computational Patterns

## Data parallel patterns

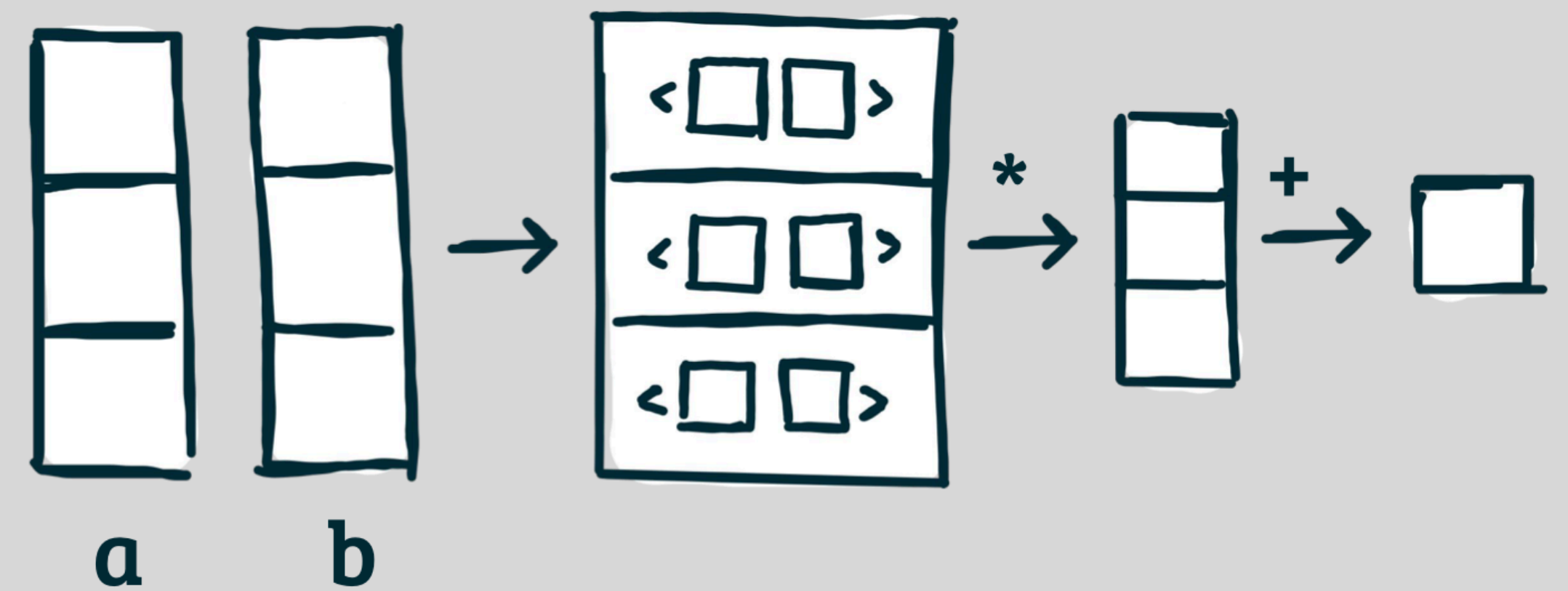


# Computational Patterns

## Data parallel patterns



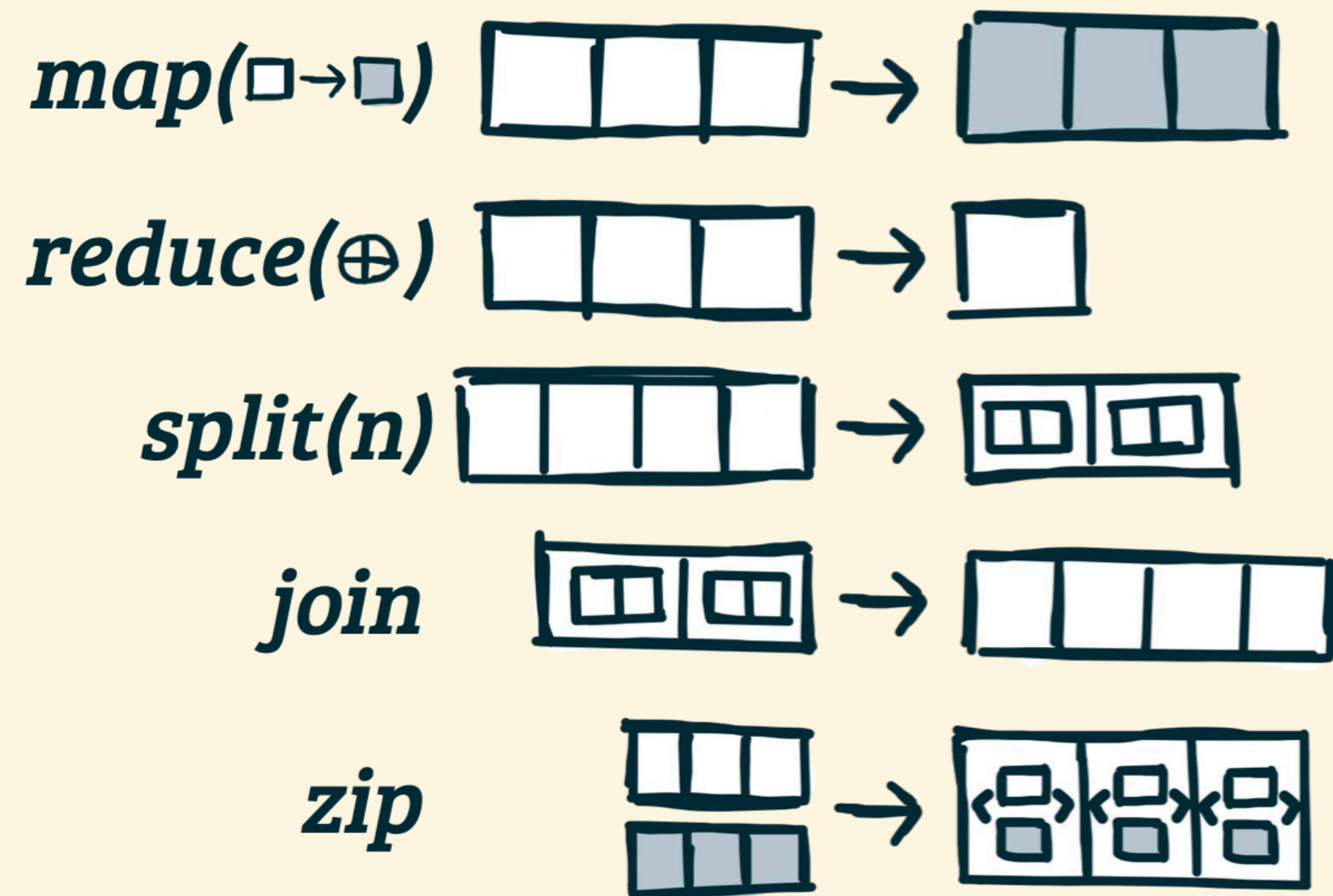
## Dot product



`zip(a, b) ▷ map(*) ▷ reduce(+, 0)`

# Computational Patterns

## Data parallel patterns



## Matrix multiply

```
depFun((m:Nat, n:Nat, k:Nat) =>
  fun((A: Array[m, Array[k, f32]],
      B: Array[k, Array[n, f32]]) =>

    A ▷ map(fun(rowA =>
      B ▷ transpose ▷ map(fun(colB =>
        dot(rowA, colB)

))))))
```

# 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

High-Level  
functional primitives

## 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 += 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(); }
```

# GEMM in RISE

High-Level functional primitives

## 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))) ))))

```

Low-Level functional primitives

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 += 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(); }

```



# GEMM in RISE

High-Level functional primitives

## 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))) ))))

```

Low-Level functional primitives

Translation

Low-Level imperative primitives

## 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 += 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(); }

```

# GEMM in RISE

High-Level functional primitives

## 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))) ))))

```

Low-Level functional primitives

Translation

Low-Level imperative primitives

## 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 += 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 + row
47     }
48   __syncthreads(); }

```

Low-Level imperative code

# 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))) ))))))
```



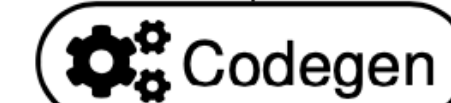
Translation

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()))))
```



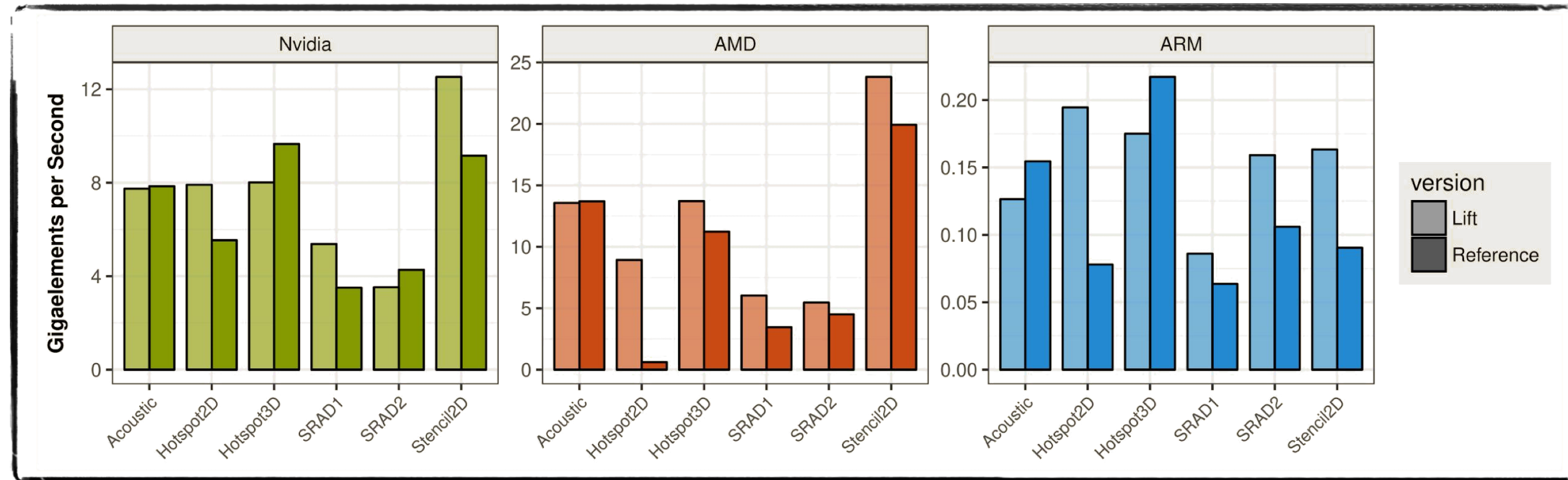
Codegen

Translation

C

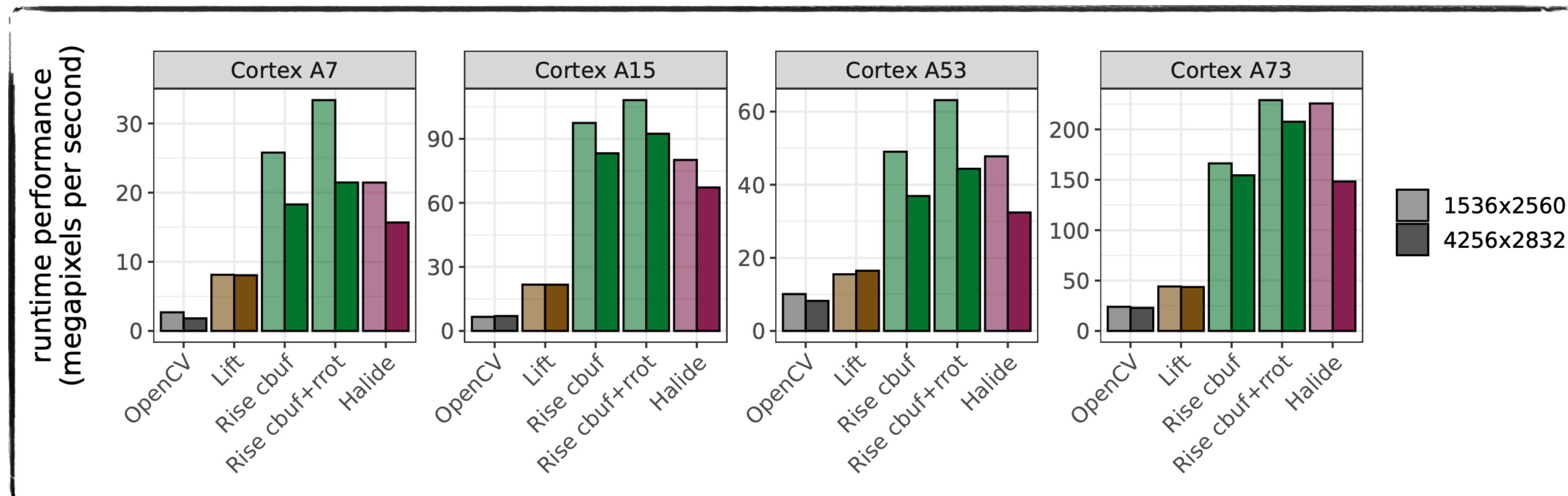
```
33 __global__ void gemm_kernel(float* __restrict__ output,
34   int m, int n, int k, const __half* __restrict__ A,
35   const __half* __restrict__ B,
36   const float* __restrict__ C, float alpha, float beta) {
37   for(int rowIdx=blockIdx.x;
38     blockIdx.x<m; rowIdx += gridDim.x) {
39     for(int 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(); }
```

# Performance Results



Same performance as hand-optimised code!

 **[CGO 2018]**



Outperform Halide with two optimizations added as new patterns.

 **[CGO 2021]**

# Extensibility!

- New patterns can be added at each abstraction layer:

- *Low-level imperative primitives* to capture hardware details
- *Low-level functional primitives* to lift these abstractions into the functional world
- High-level functional primitives to make these abstractions available to rewriting

```
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);
```

## Low-level imperative primitives

```
Fragment[m: Nat, n: Nat, k: Nat, t: DataType, f: FragmKind]

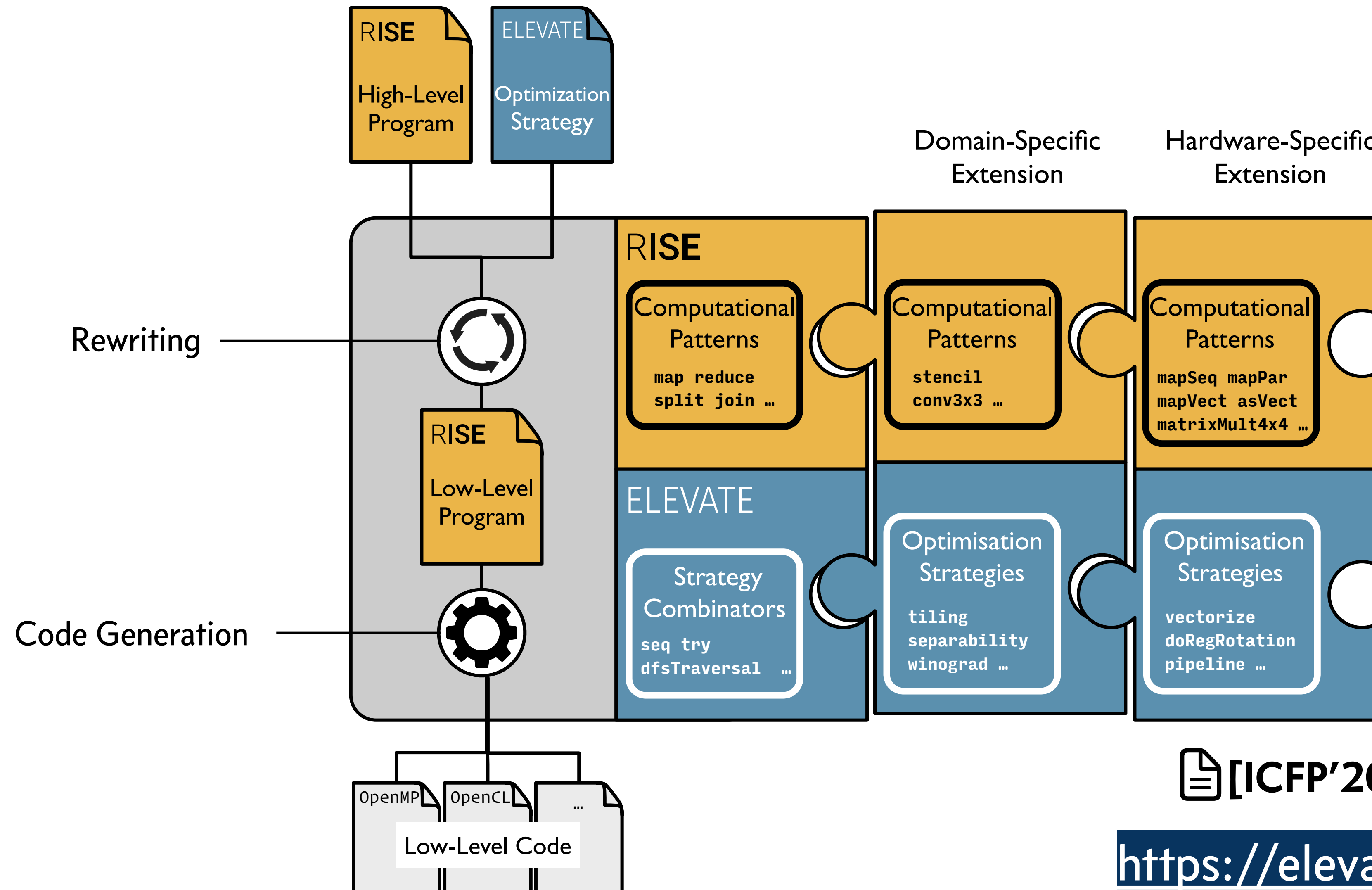
def mmaFragment(m:Nat, n:Nat, k:Nat, s:DataType, t:DataType,
    A: Exp[Fragment[m,k,n,s,AMatrix], Rd],
    B: Exp[Fragment[k,n,m,s,BMatrix], Rd],
    C: Exp[Fragment[m,n,k,t,Accum], Rd],
    D: Acc[Fragment[m,n,k,t,Accum]]): Comm
def loadFragment(f:FragmKind, m:Nat, n:Nat, k:Nat, t:DataType,
    tile: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

## Low-level 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]
```

*OPTIMISATION*

# Extensible Optimizations via Rewriting



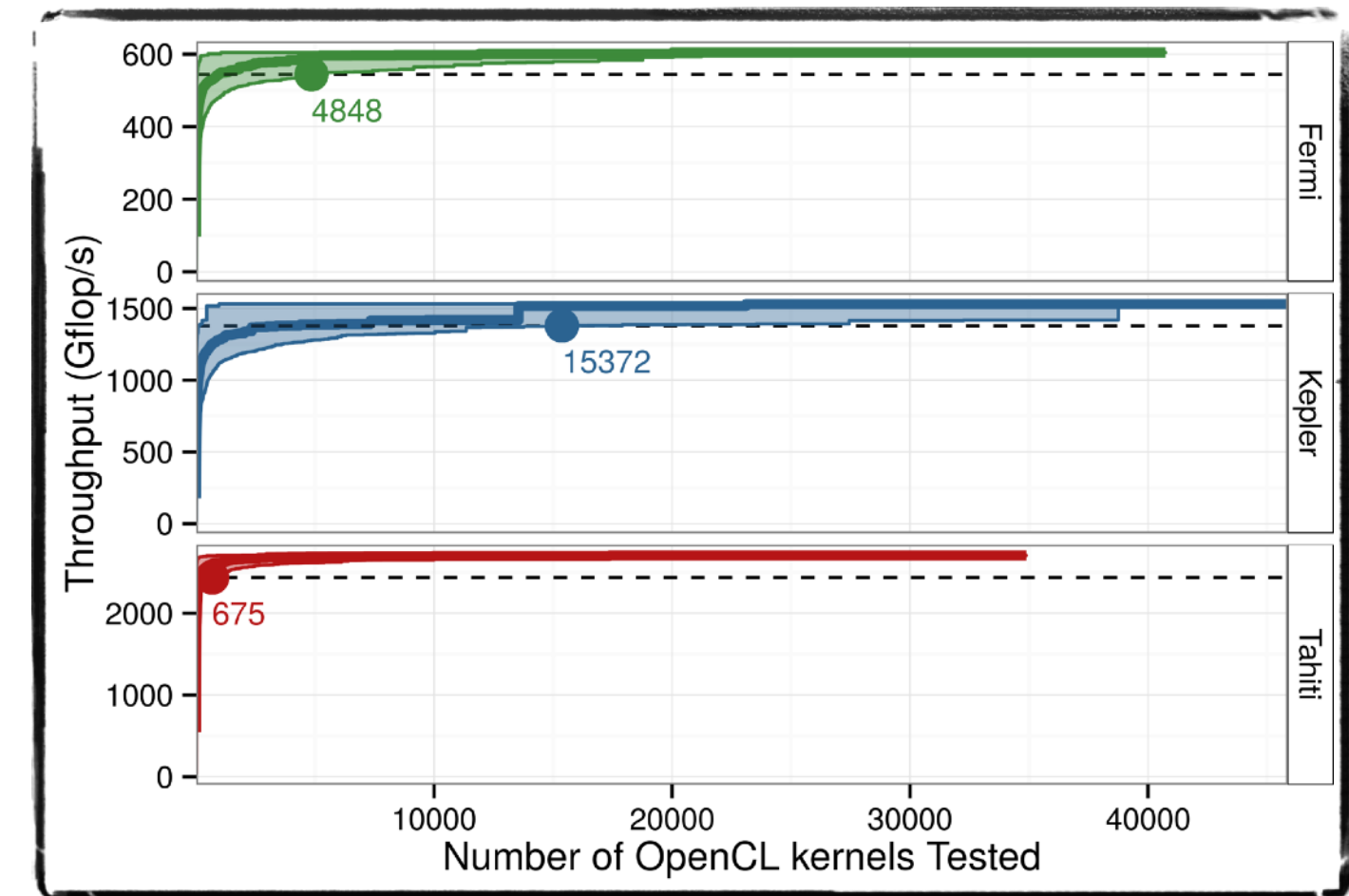
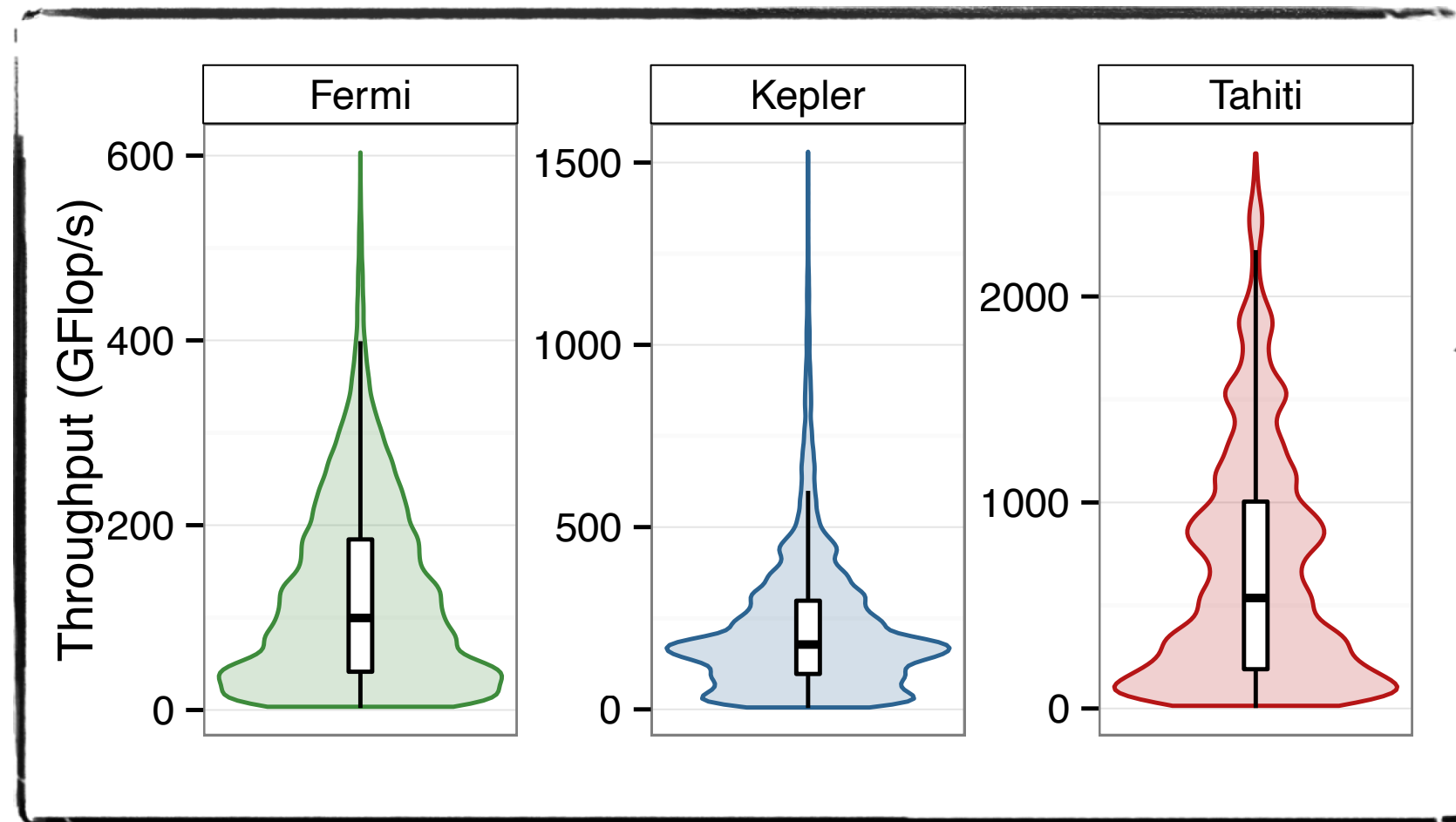


# Tradeoffs when optimizing with rewriting



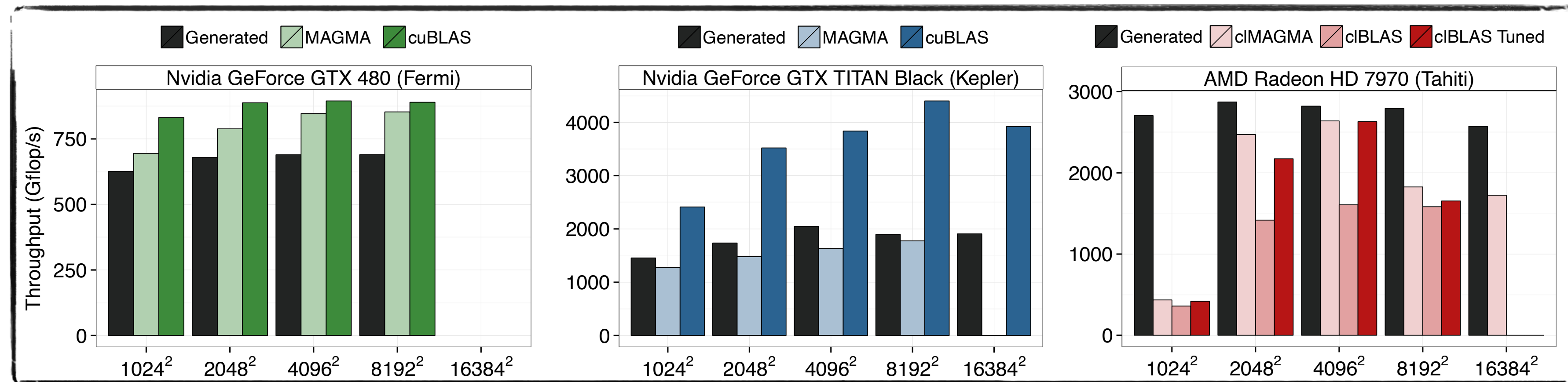
- ✓ No human needed in optimization process

# Automatic Rewriting for Matrix Multiplication



Only few generated code with very good performance

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



Performance close or better than hand-tuned library code

# Tradeoffs when optimizing with rewriting



- ✓ No human needed in optimization process
- ✗ **Costly & Lengthy search process**
- ✗ **Does not (yet) scale to complex programs**

# Tradeoffs when optimizing with rewriting



✓ No human needed in optimization process

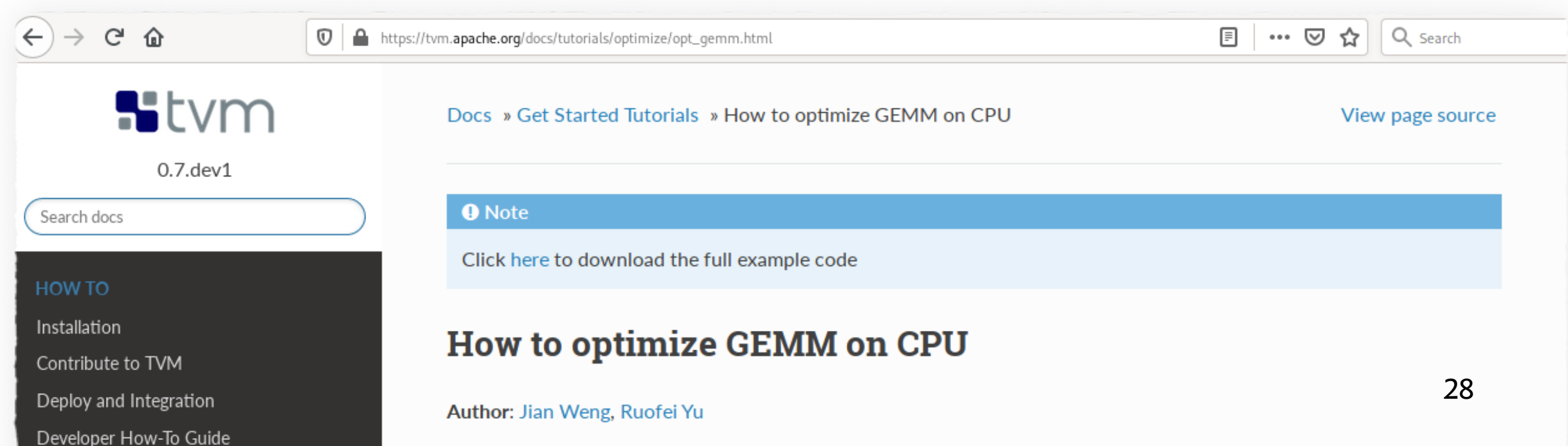
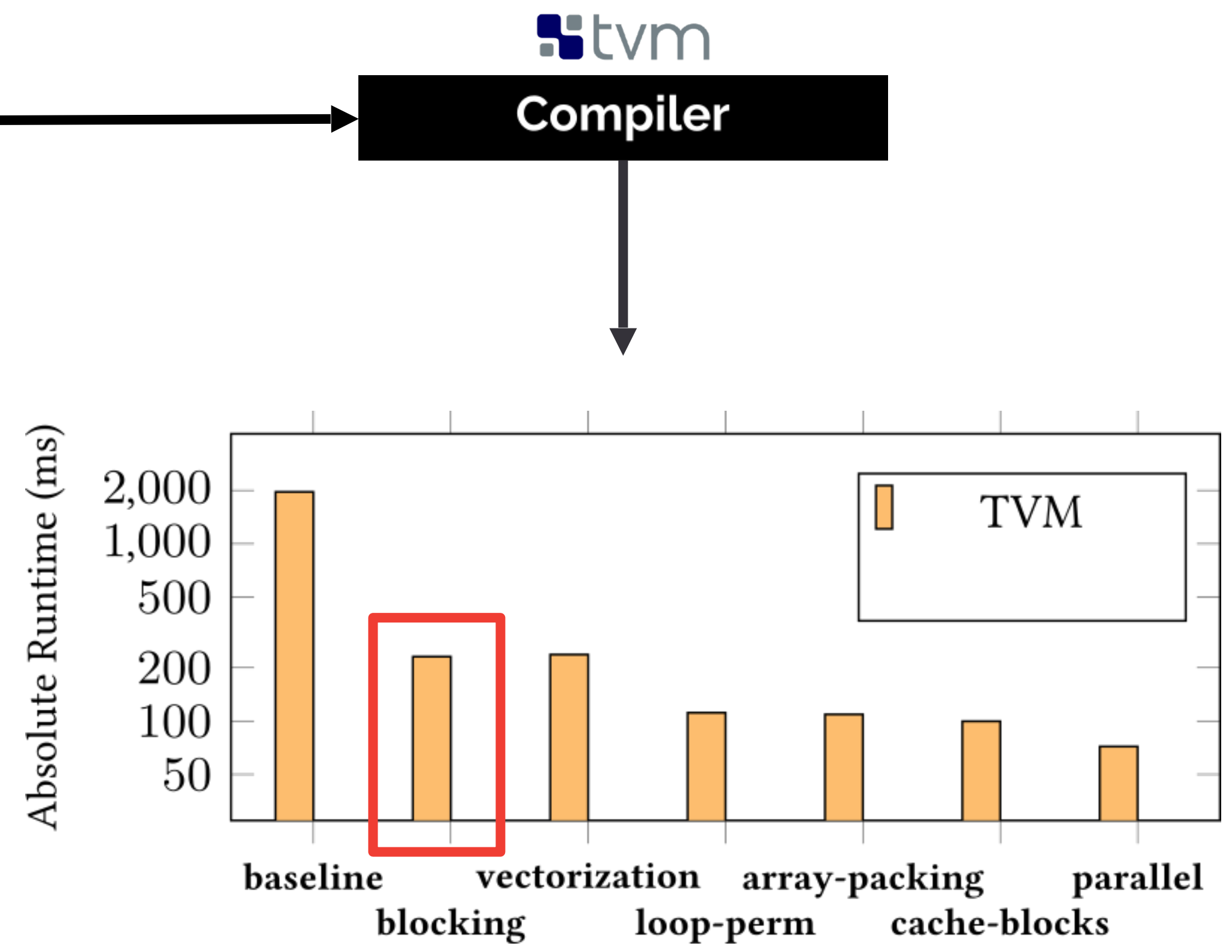
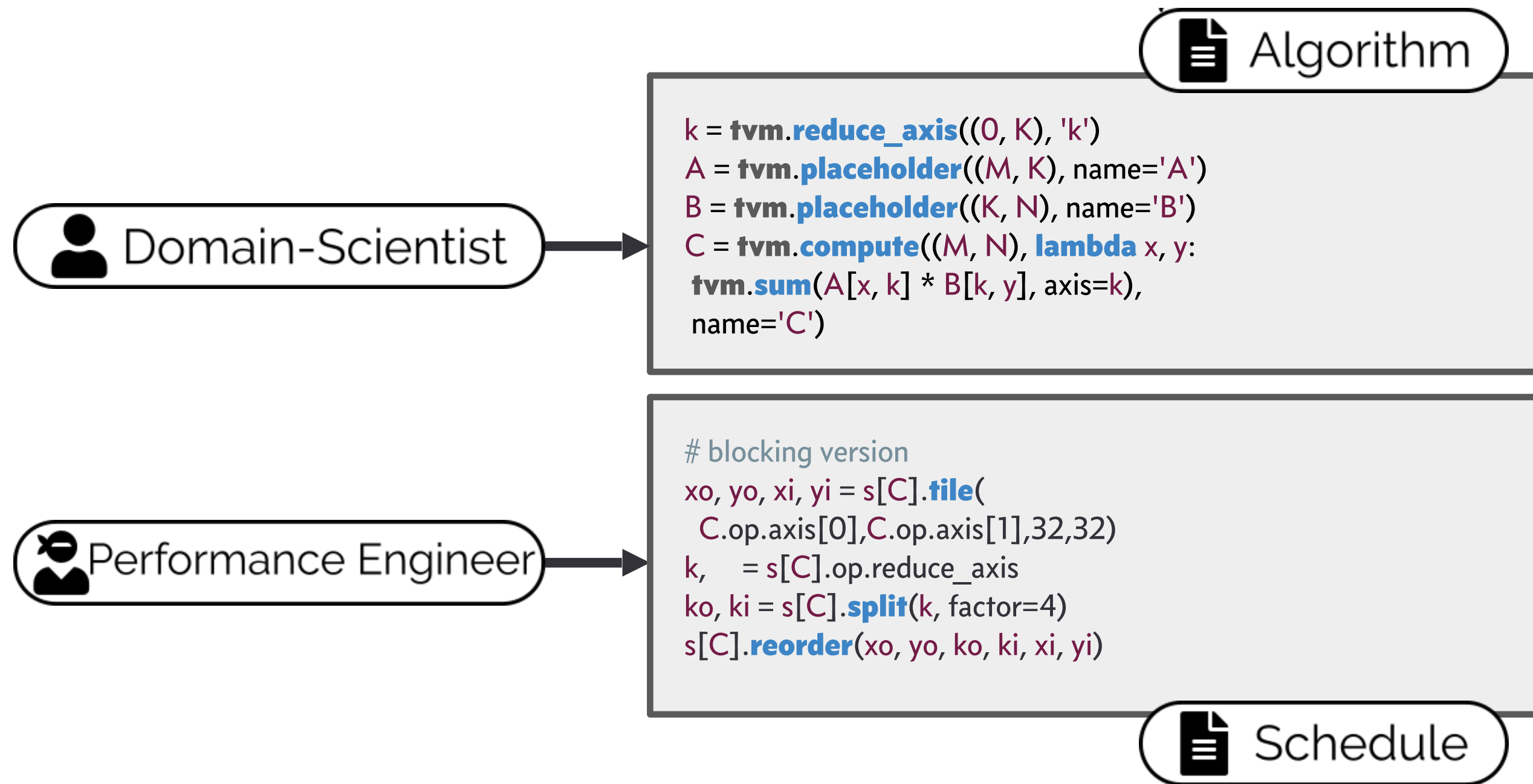
✗ Costly & Lengthy search process

✗ Does not (yet) scale to all programs

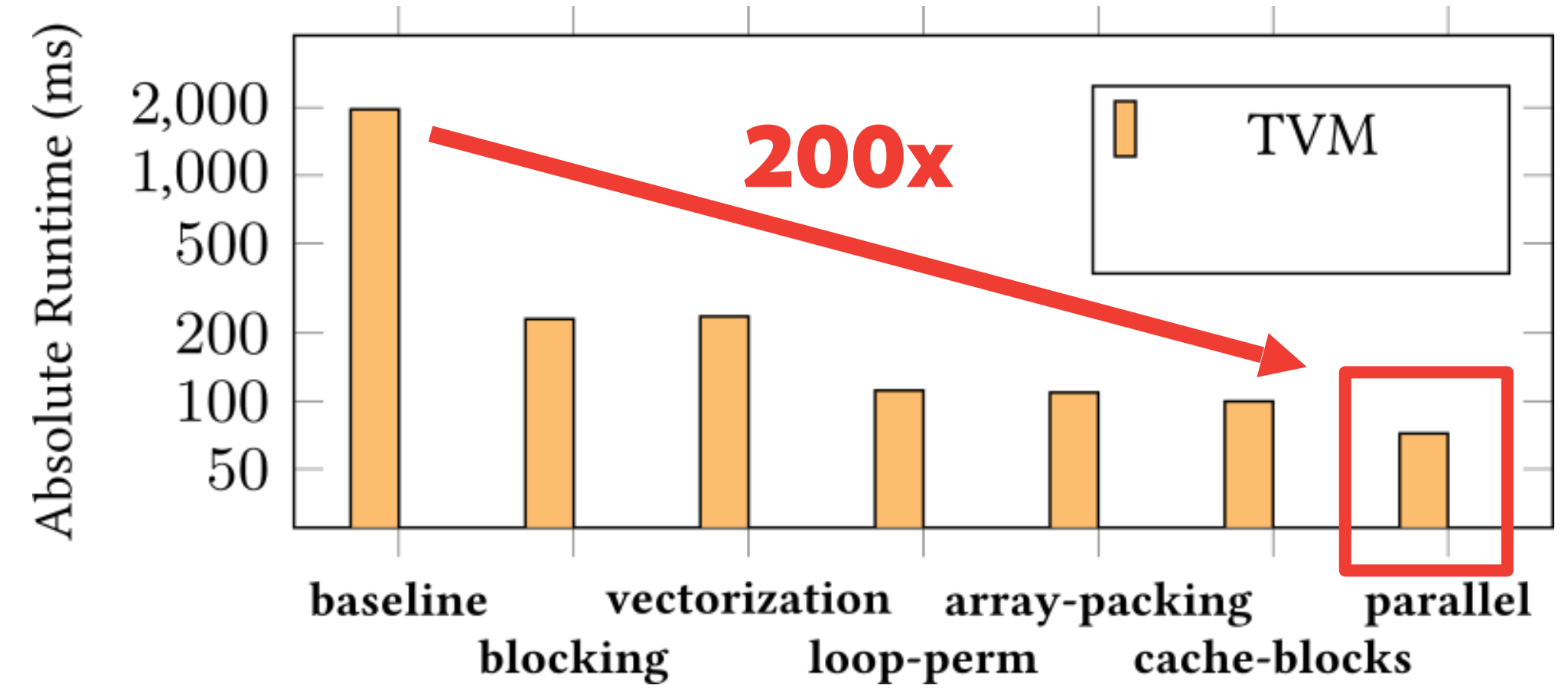
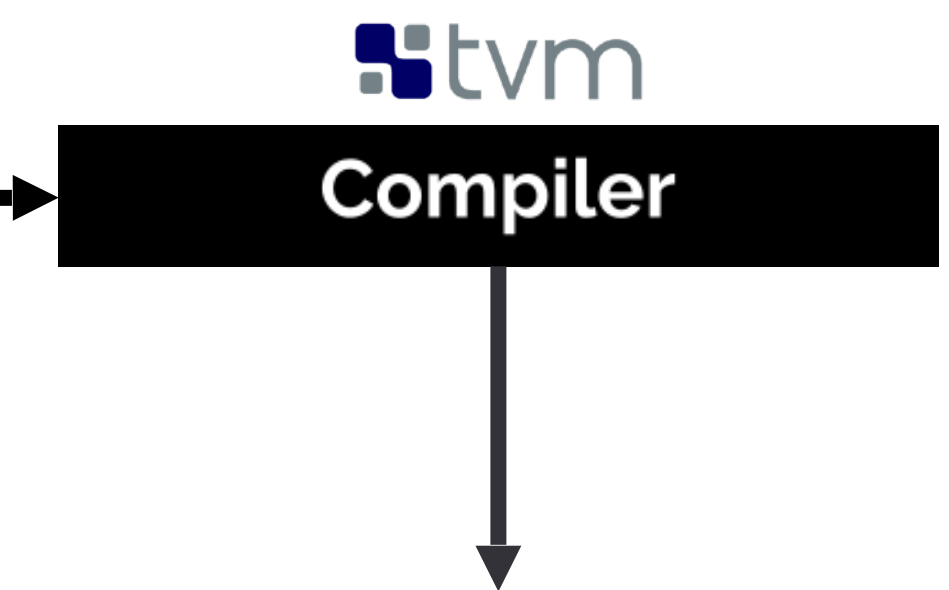
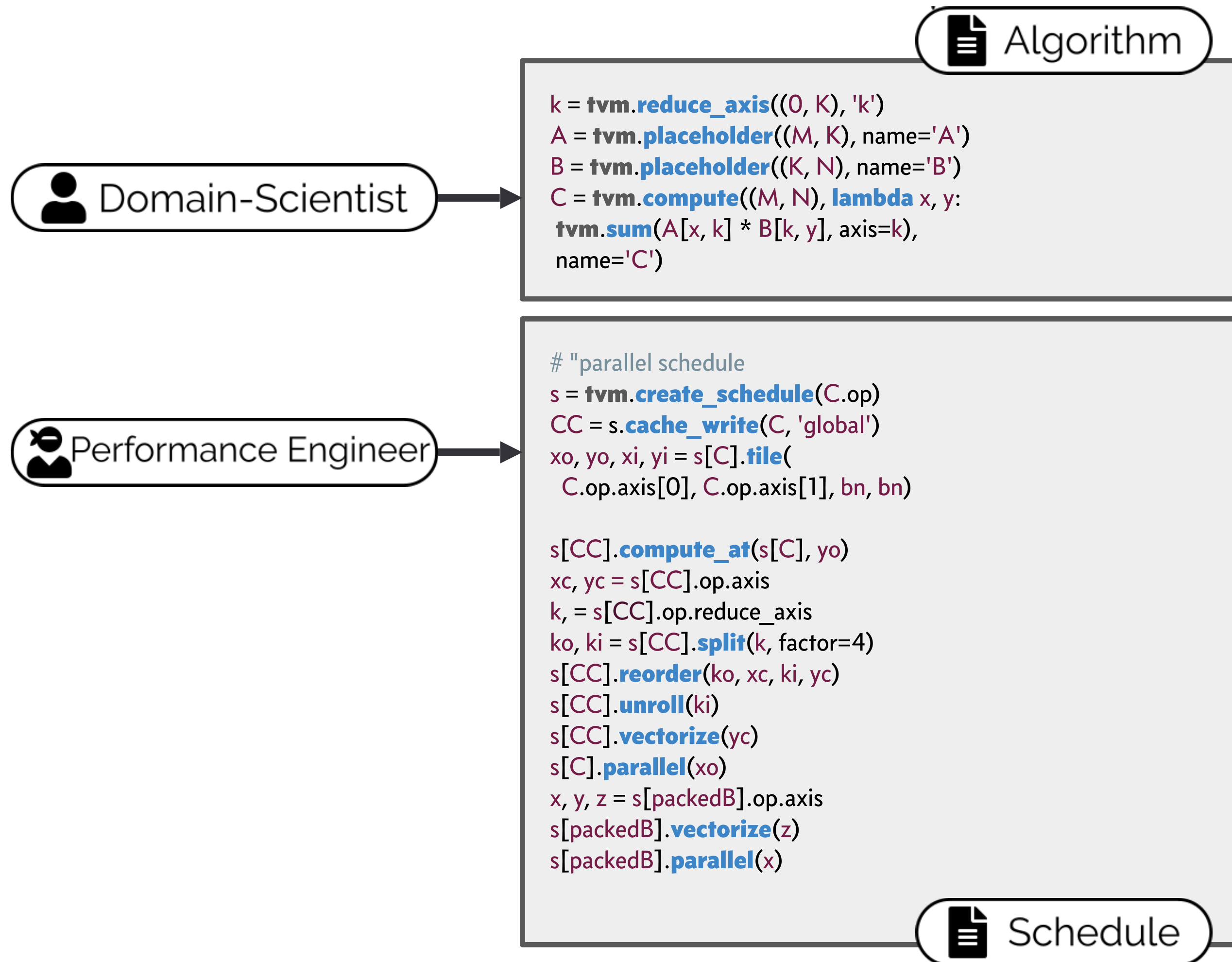
Extensive human effort needed ✗

Expert is in control, no search required ✓

# Compilers with scheduling APIs



# Compilers with scheduling APIs



# Compilers with scheduling APIs

## Compilers with scheduling APIs

**Halide**



Tiramisu-Compiler / tiramisu

Fireiron



Domain-Scientist

Performance Engineer

Program

Algorithm

```
// 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 Nvidia 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,rx;
out.bounds(x, 0, size).bounds(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, rx, warp_size)
  .unroll(r.x, r_unroll).reorder(xi, xo, y, rx, 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, rx).vectorize(Ax, vec_size)
  .split(Ax, xo, xi, warp_size).gpu_lanes(xi)
  .unroll(xo).unroll(Ay);
```

Schedule

Optimisation Schedule



**Halide**  
compiler

Optimised Code

# Problems with Scheduling APIs

## Program

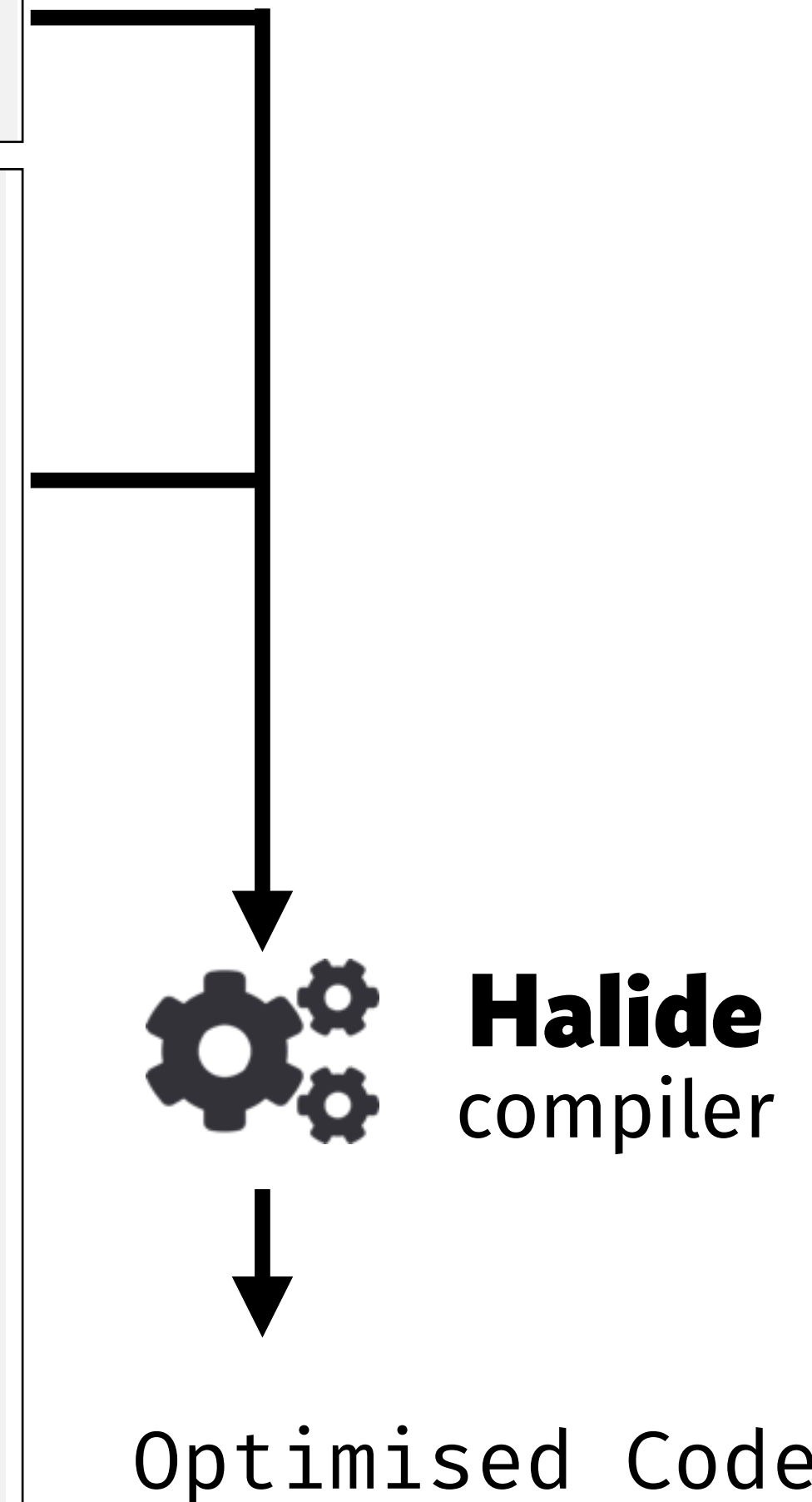
```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
out(x, y) = prod(x, y);
```

```
// schedule for Nvidia 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





# Problems with Scheduling APIs

No clear separation

## 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 Nvidia 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

# Problems with Scheduling APIs

**Hinders reuse**

## Program

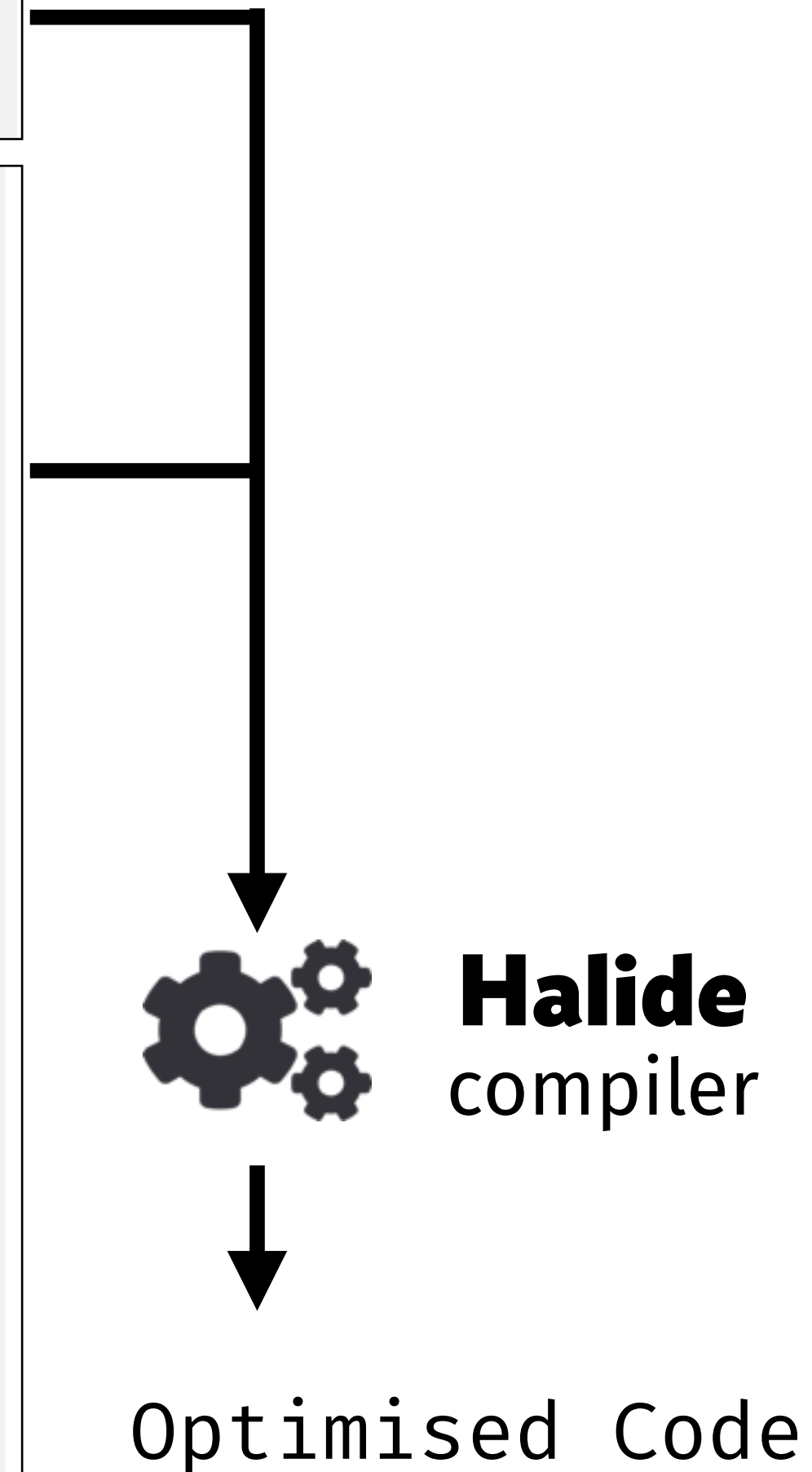
```
// functional description of matrix multiplication
Var x("x"), y("y"); Func prod("prod"); RDom r(0, size);
prod(x, y) += A(x, r) * B(r, y);
out(x, y) = prod(x, y);
```

```
// schedule for Nvidia 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



# 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 Nvidia GPUs
const int warp_size = 32; const int vec_size = 2;
const int x_tile = 3; const int y_tile = 4;
const int r_unroll = 8; const int r_unroll = 1;
```

**Hinders reuse**

```
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);
```

Not well defined semantics

```
.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);
```

```
// schedule for Nvidia 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;
```

**Hinders reuse**

**Not well understood  
Hinders understanding**

```
.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

# 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 Nvidia 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;
```

**Hinders reuse**

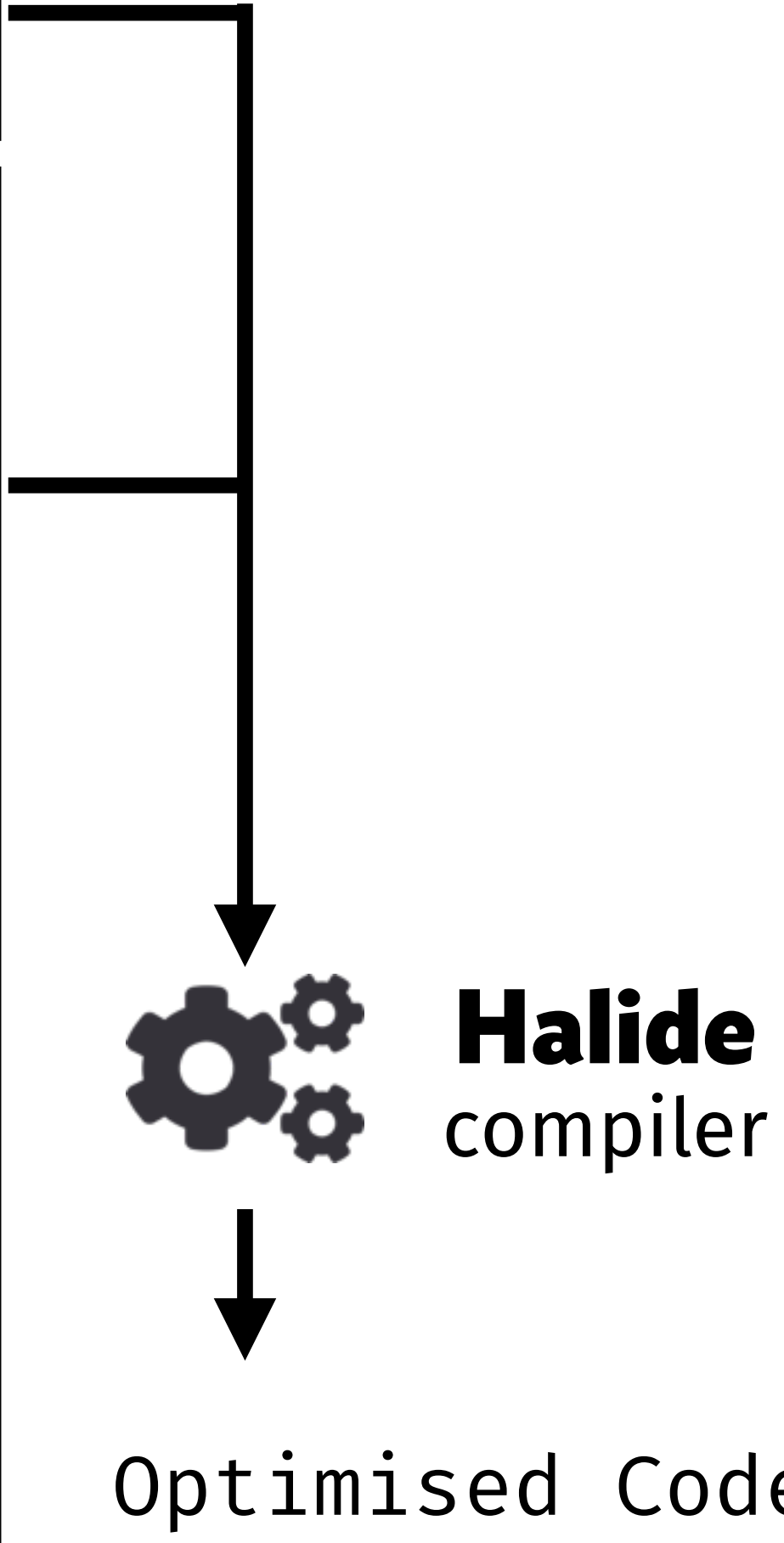
**Hinders understanding**

Only fixed built-in optimisations

```
.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



# 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 Nvidia 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;
```

**Hinders reuse**

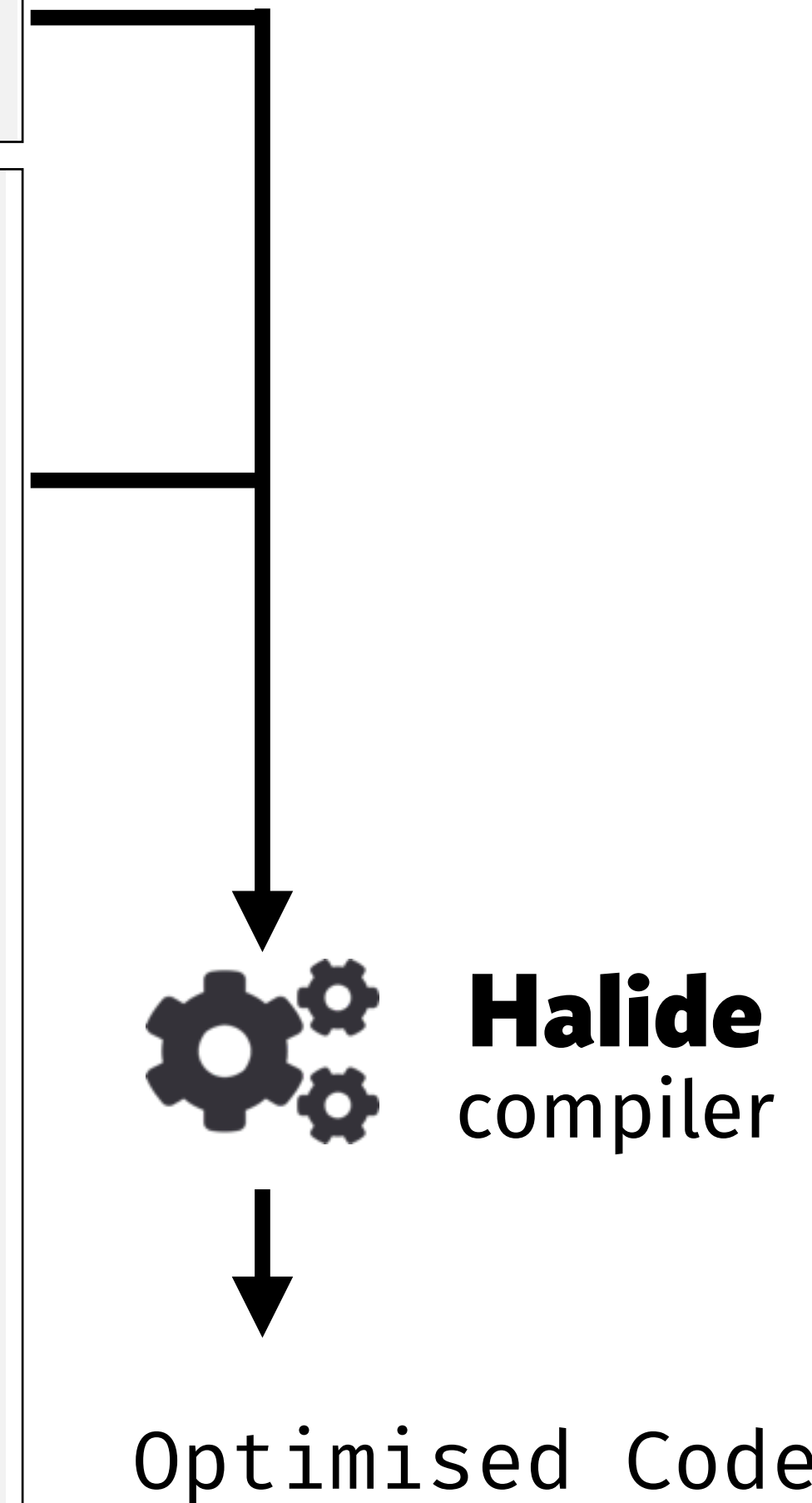
**Not well understood  
Hinders understanding**

**Only fixed combinations  
No extensibility**

```
.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



# 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 Nvidia GPUs  
const int warp_size = 32; const int vec_size = 2;  
const int x_tile = 3; const int y_tile = 4;  
const int r_unroll = 8; const int r_unroll = 1;
```

**Hinders reuse**

**Hinders understanding**

**No extensibility**

```
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];
```

```
A.in().in().compute_at(prod, rx1).vectorize(Ax, vec_size)  
.split(Ax, xo, xi, warp_size).gpu_lanes(xi)  
.unroll(xo).unroll(Ay);
```



**Halide**  
compiler

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

## Optimisation Schedule

# 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.*



# 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.

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**

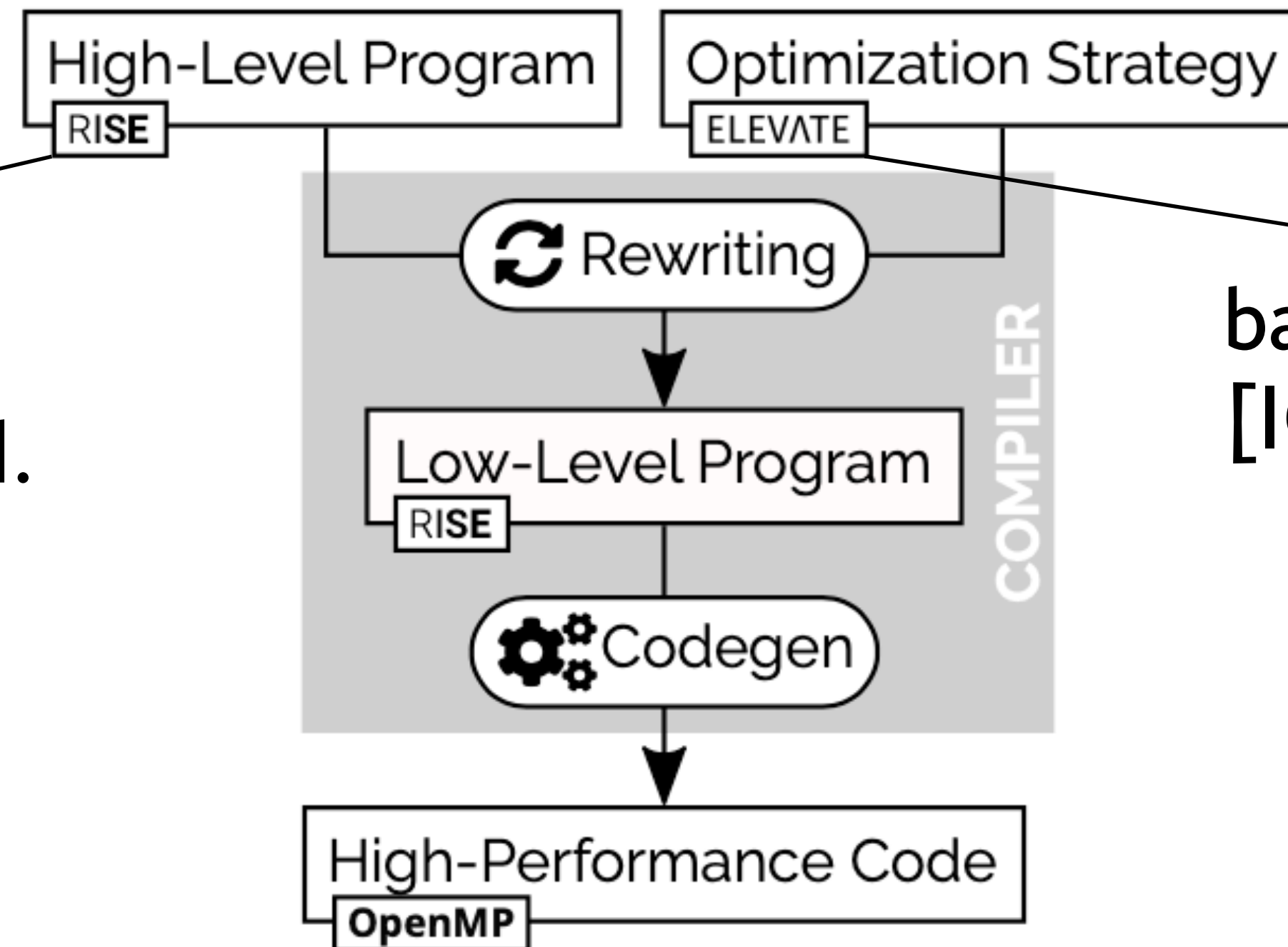
*Implicit default behavior should be avoided to empower users to be in control.*

# 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) 'a' outermost(mapNest(2)) ';;'
  fissionReduceMap 'a' outermost(appliedReduce) ';;'
  split(4) 'a' innermost(appliedReduce) ';;'
  reorder(Seq(1,2,5,3,6,4)) ';;'
  vectorize(32) 'a' 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])
```

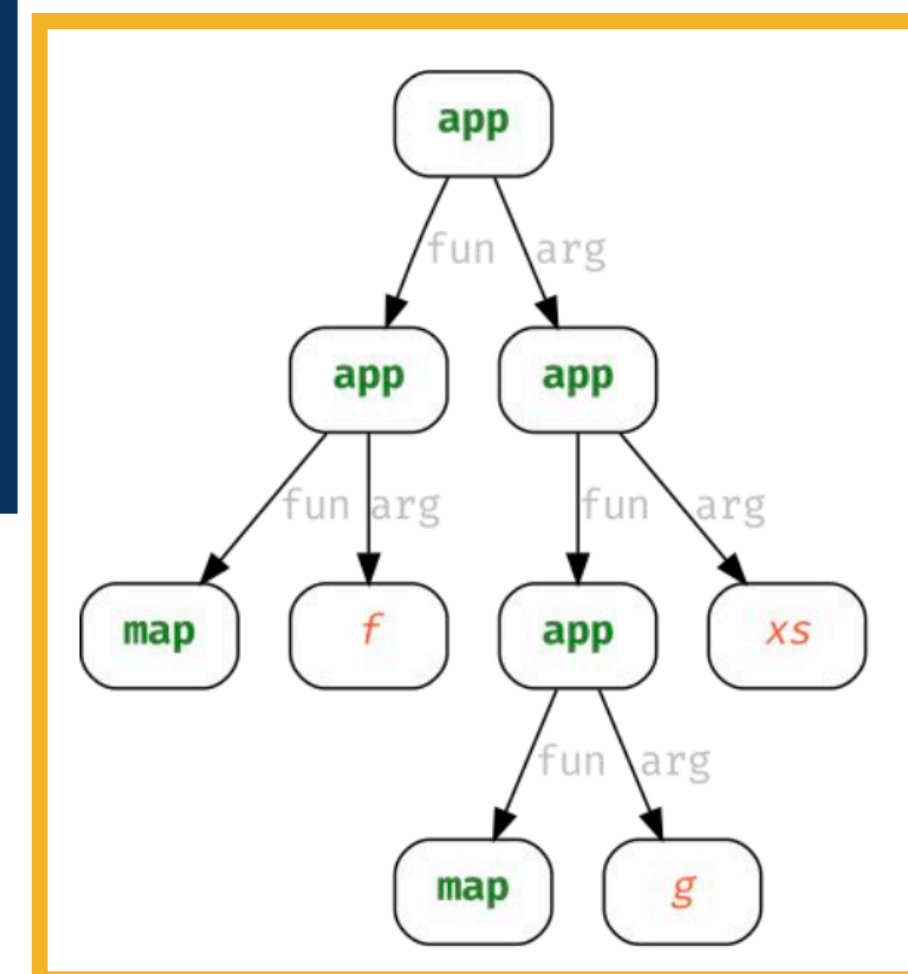
# Rewrite Rules in ELEVATE

- Rewrite rules are basic strategies

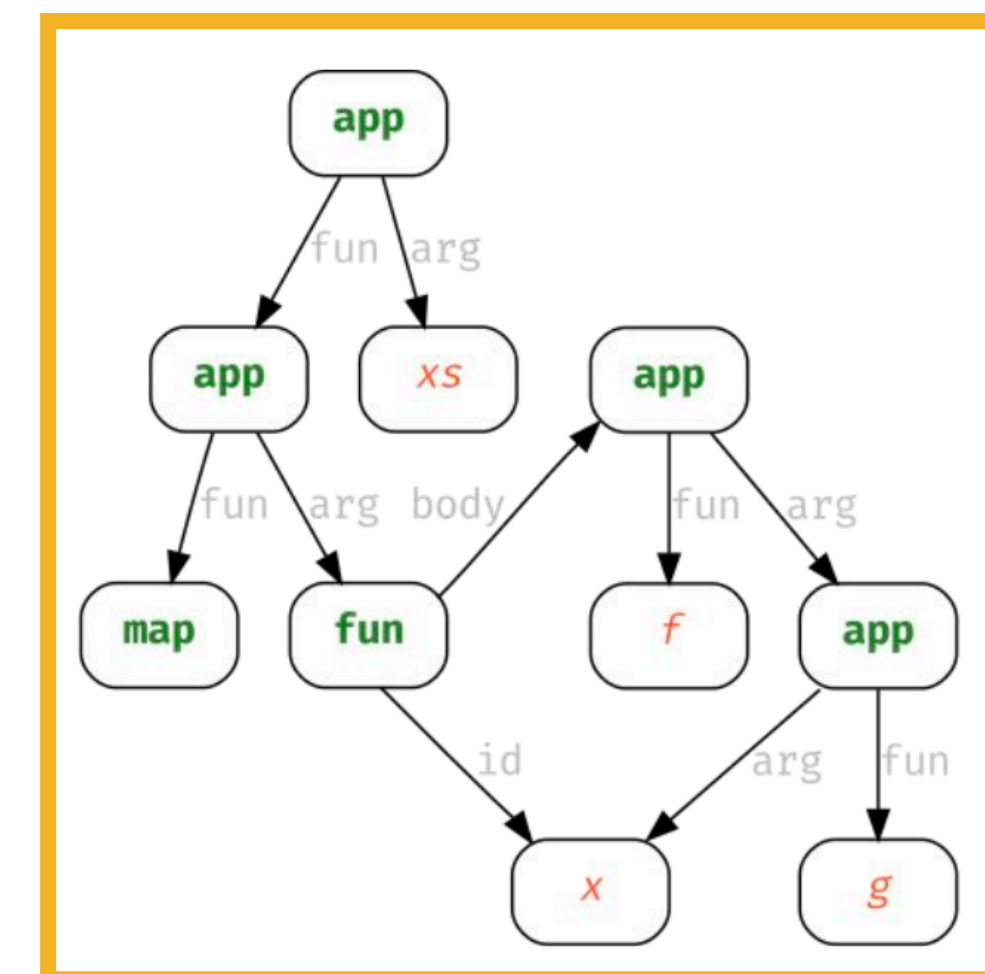
$$\text{map}(f) \ll \text{map}(g) \rightsquigarrow \text{map}(f \ll g)$$

```
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 ones

- 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

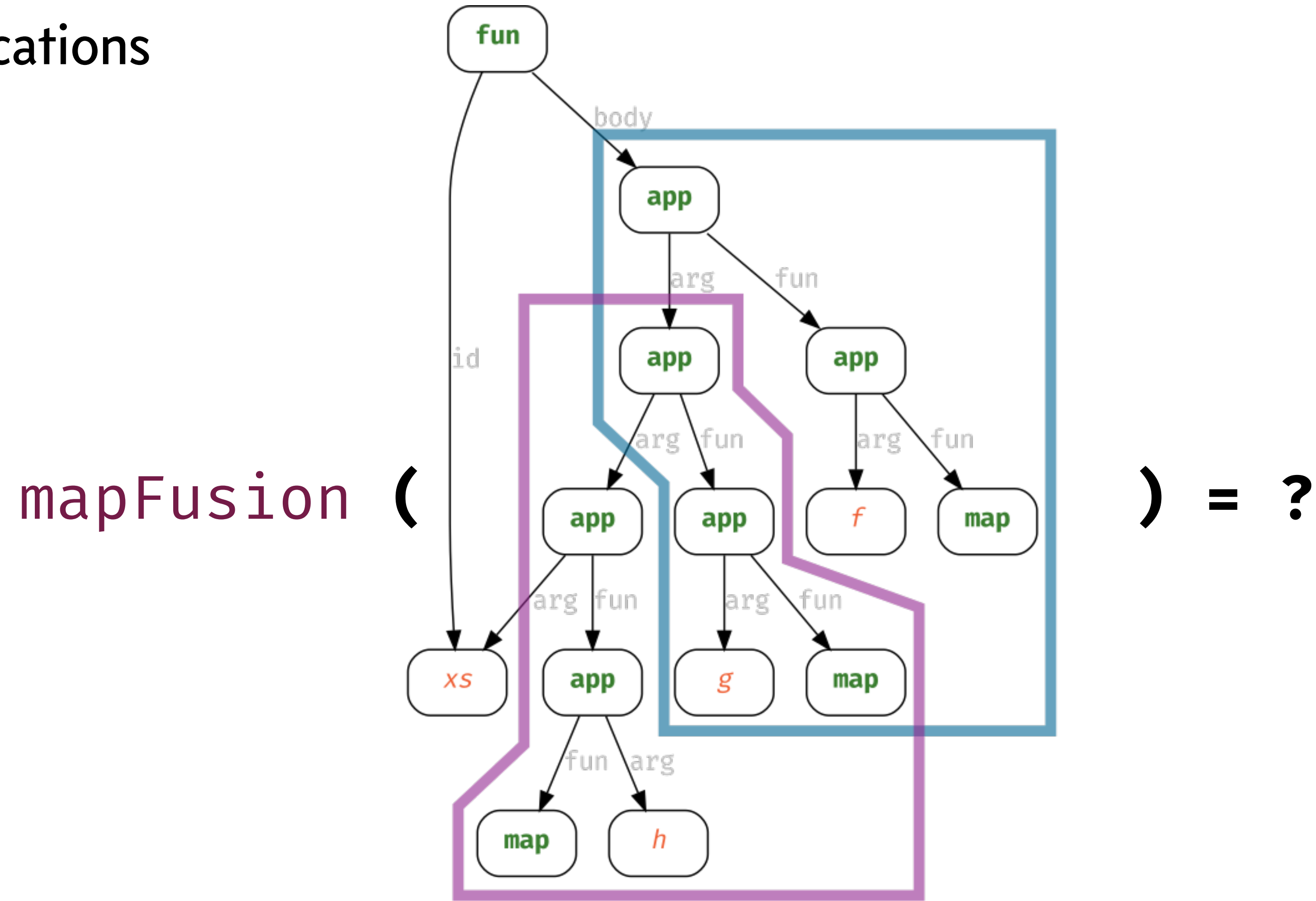
```
def try[P]: Strategy[P] => Strategy[P] =  
  s => p => (s <+ id)(p)
```

- Repeat

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

# Traversals in ELEVATE

- Describing Precise Locations



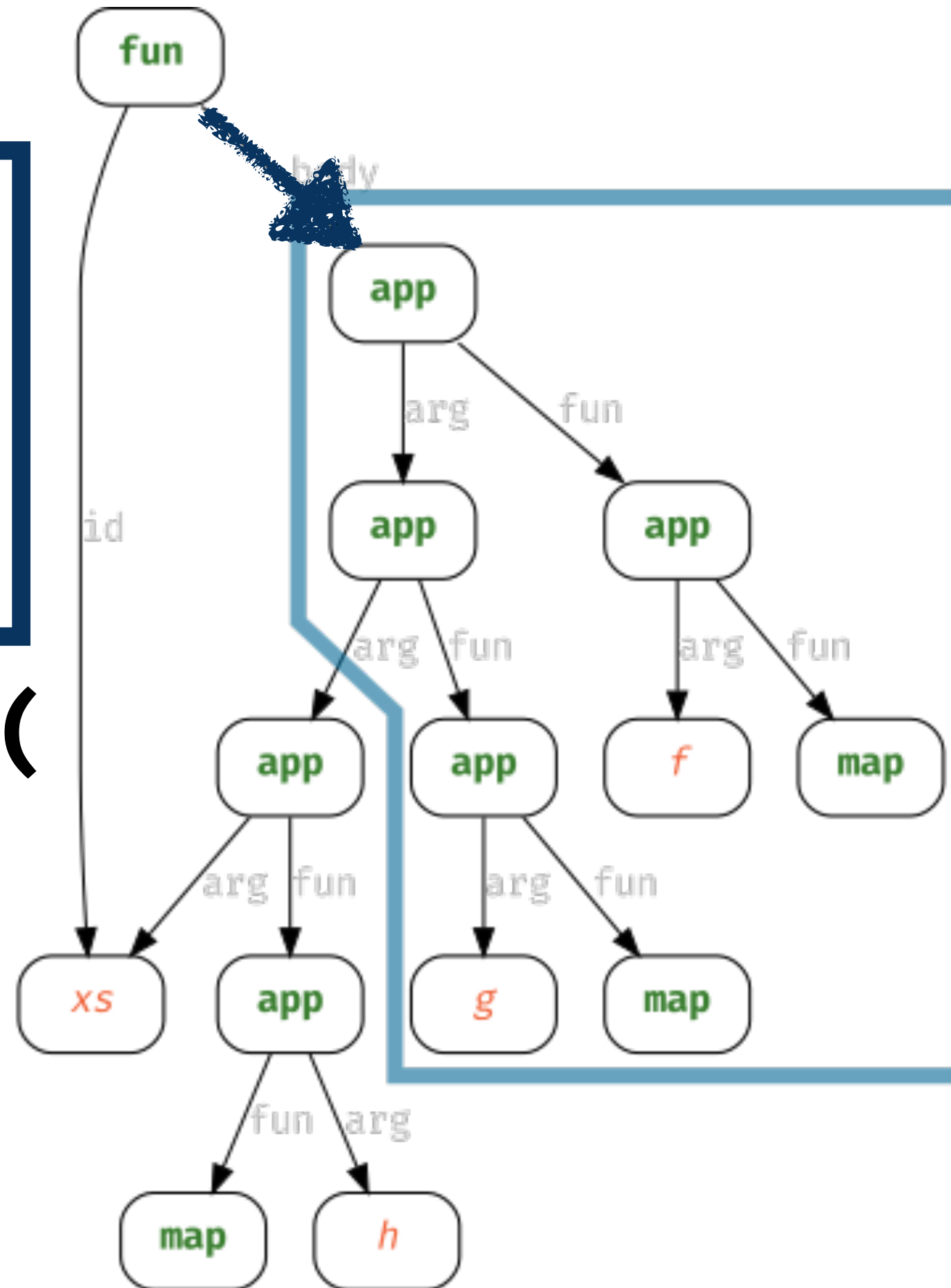
threemaps = fun(xs ⇒ map(f)(map(g)(map(h)(xs)))))

# 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(mapFusion)



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

# Traversals in ELEVATE

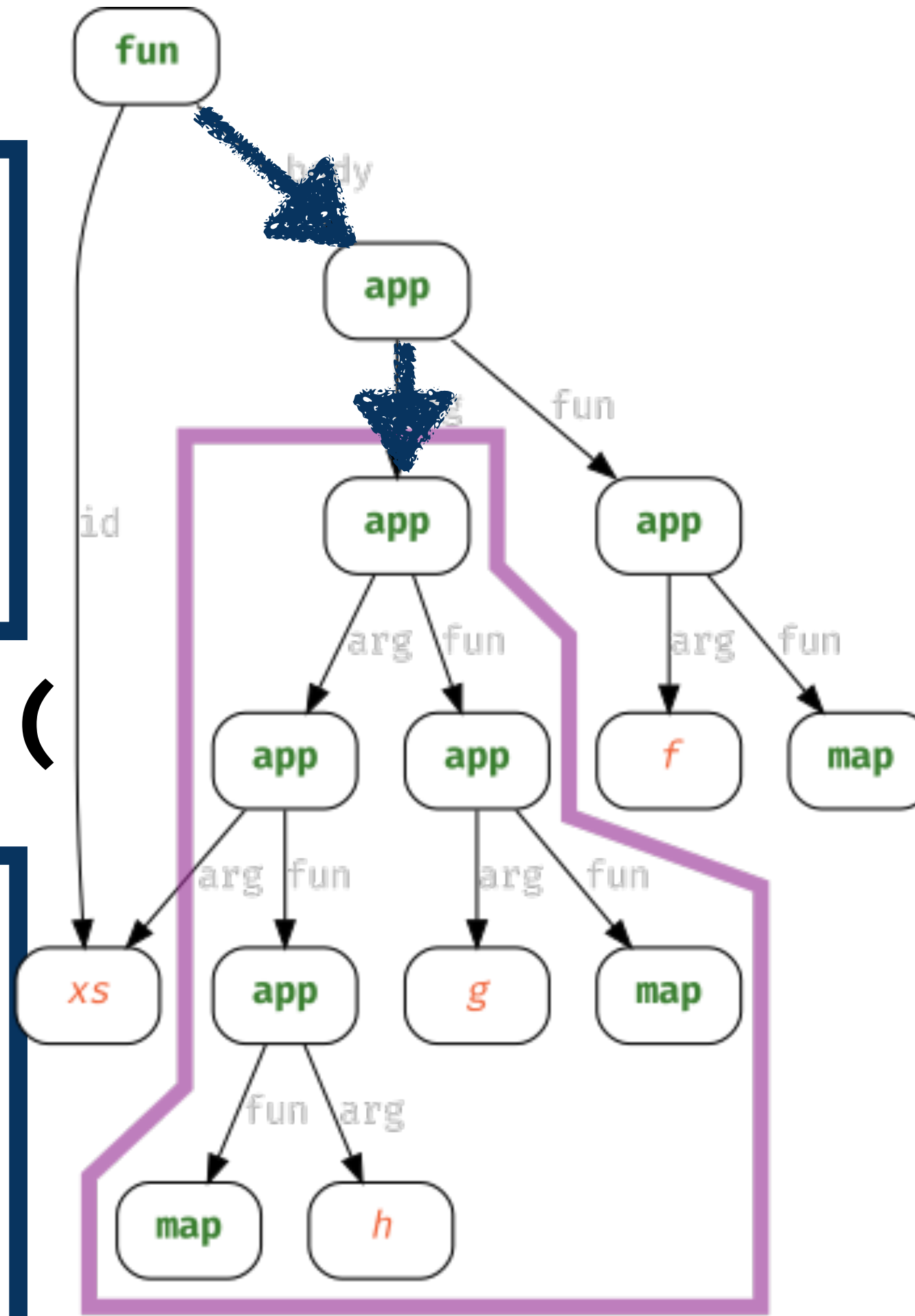
- Describing Precise Locations

```
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]
```

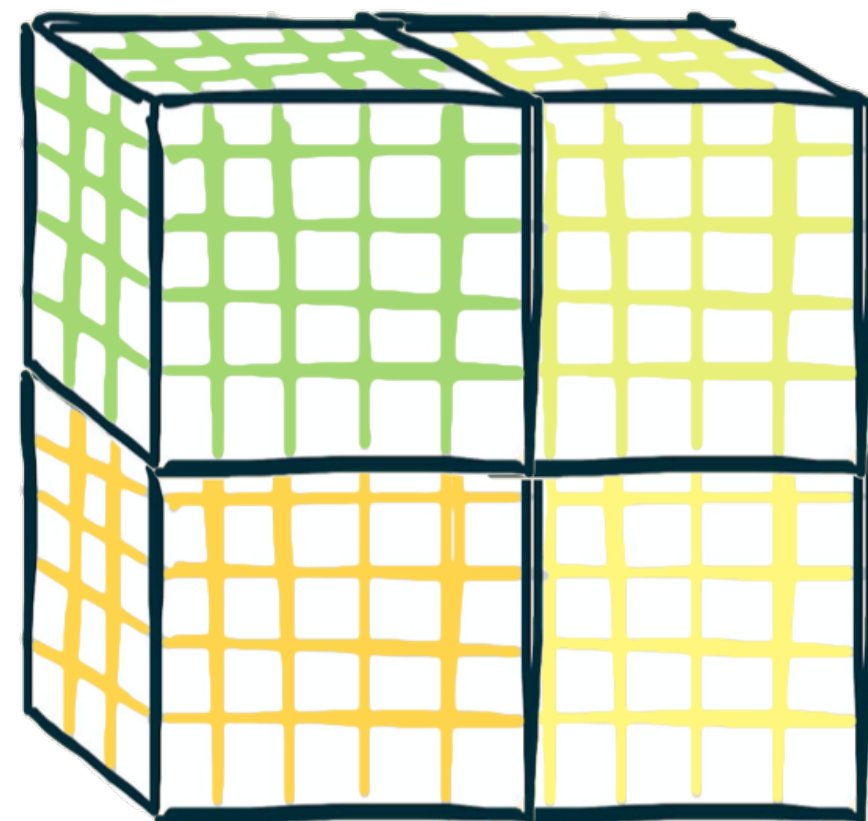
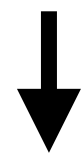
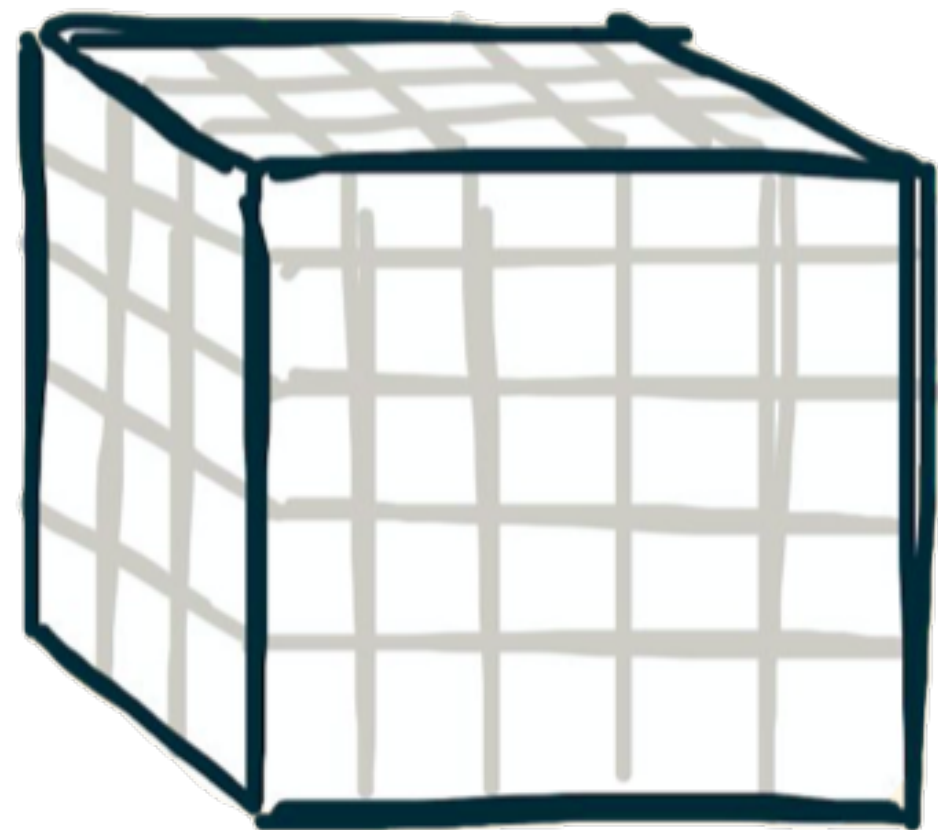
- we define more complex traversals:

```
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

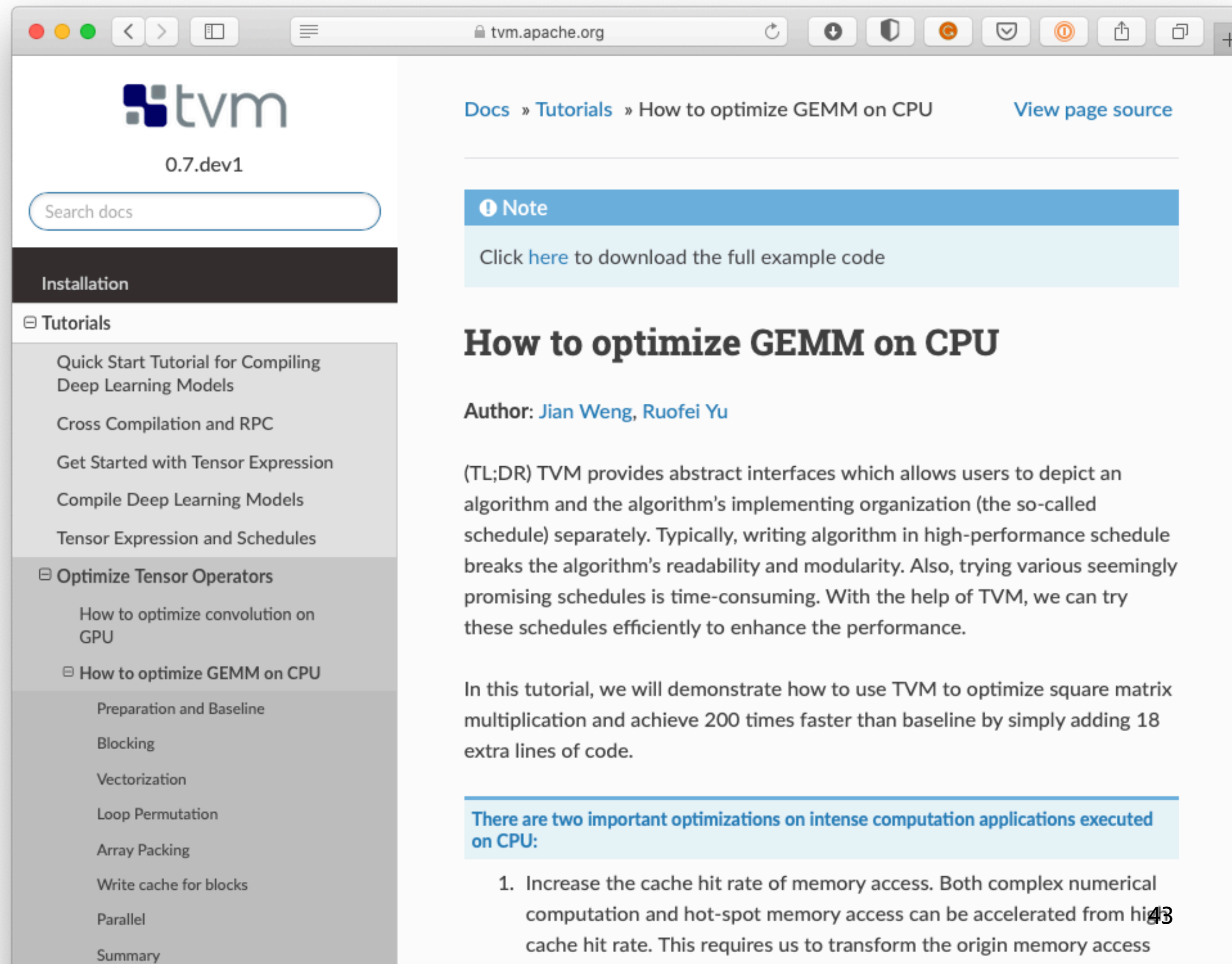


```
def tile: Int → Int → Strategy =  
  (dim) ⇒ (n) ⇒ dim match {  
    case 1 = function(splitJoin(n))  
    case 2 = fmap(function(splitJoin(n))) ;  
              function(splitJoin(n)) ; interchange(2)  
    case i = fmap(tile(dim-1, n)) ;  
              function(splitJoin(n)) ; interchange(n)  
  }
```

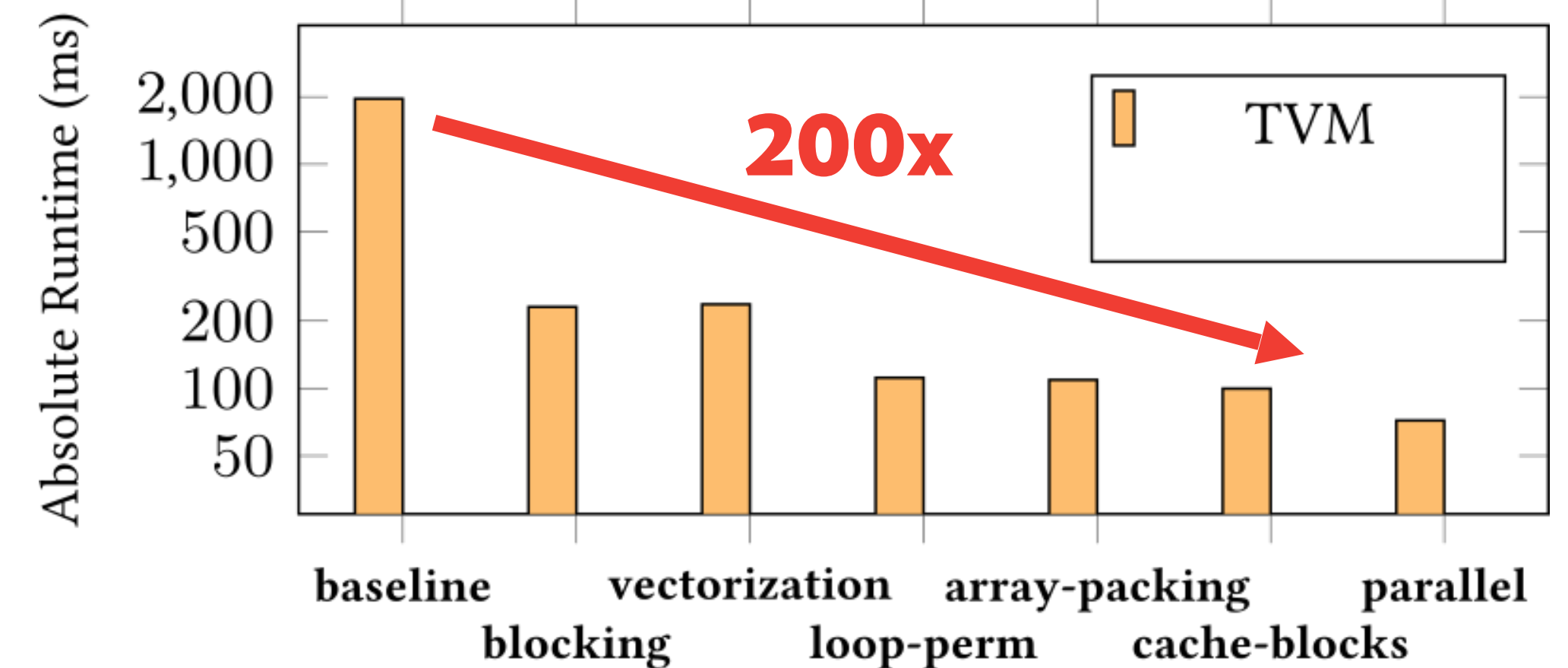
Tiling defined as composition of rewrites not a built-in!

# Case Study: Implementing TVM's Scheduling API

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



The screenshot shows the TVM documentation website (tvm.apache.org) for version 0.7.dev1. The page is titled "How to optimize GEMM on CPU" and is part of the "Tutorials" section. A note at the top suggests downloading the full example code. The main content area contains the title "How to optimize GEMM on CPU" by Jian Weng and Ruofei Yu. The text explains that TVM provides abstract interfaces for algorithm and schedule, and that using TVM can significantly improve performance. A sidebar on the left lists various tutorials, including "Optimize Tensor Operators" and "How to optimize GEMM on CPU".



# 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(row, transpose(b) |>
7     map( fun(col,
8       dot(row)(col) )))) ) )
```

```
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)
```

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(row, transpose(b) |>
7     map( fun(col,
8       dot(row)(col) )))) ) )
```

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

## ELEVATE

Be explicit

Enable composability

Baseline Strategy



```
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 behavior

# 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)
```

```
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

## 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

vs

No clear separation of concerns

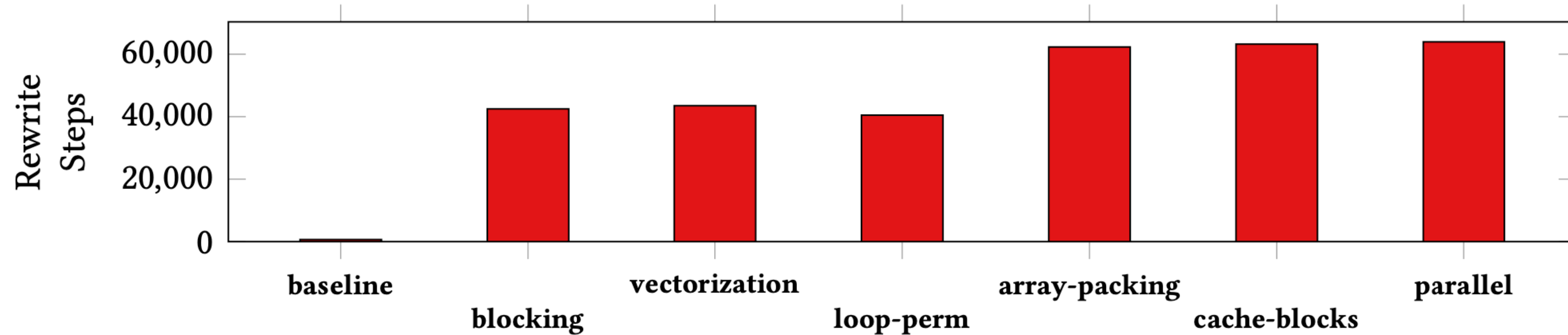


```
1 # Modified algorithm
2 bn = 32
3 k = tvm.reduce_axis((0, K), 'k')
4 A = tvm.placeholder((M, K), name='A')
5 B = tvm.placeholder((K, N), name='B')
6 pB = tvm.compute((N / bn, K, bn),
7   lambda x, y, z: B[y, x * bn + z], name='pB')
8 C = tvm.compute((M, N), lambda x, y:
9   tvm.sum(A[x, k] * pB[y//bn, 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

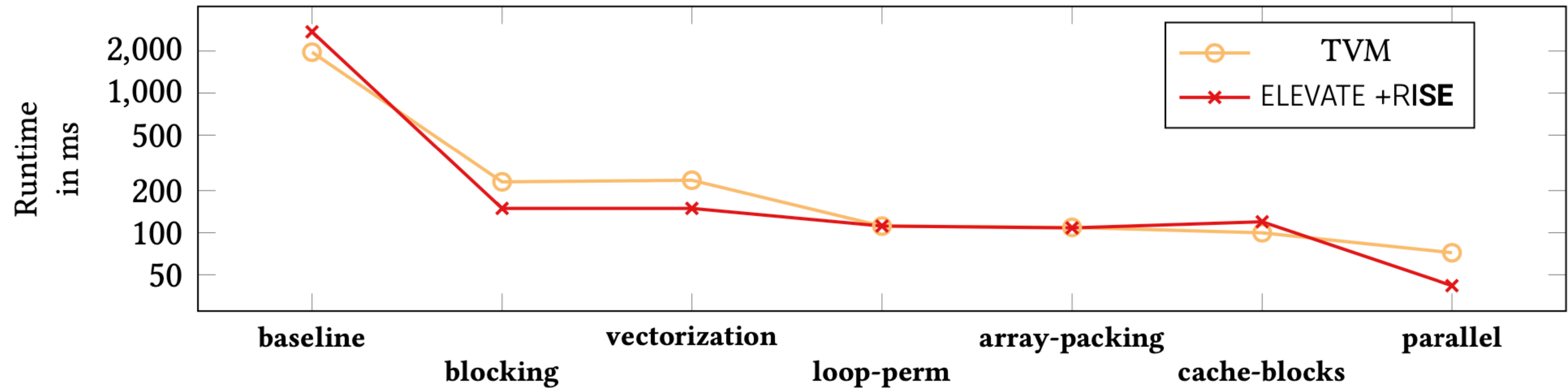


Rewriting took less than 2 seconds with our unoptimised implementation

**Rewrite based approach scales to complex optimizations**

# Optimizing Matrix Matrix Multiplication with ELEVATE Strategies

Performance of generated code



**Competitive performance compared to TVM compiler**

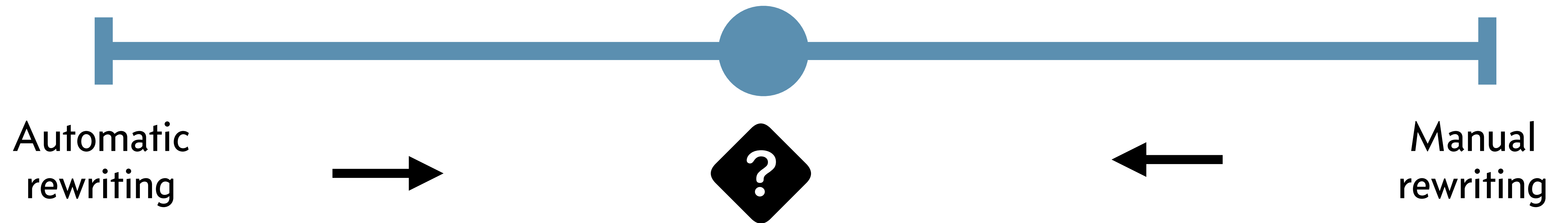
# Tradeoffs when optimizing with rewriting



- ✓ No human needed in optimization process
- ✗ Costly & Lengthy search process
- ✗ Does not (yet) scale to all programs

- Extensive human effort needed ✗
- Expert is in control, no search required ✓
- Strategies are too sensitive**  
⇒ **don't scale across applications** ✗

# Tradeoffs when optimizing with rewriting



✓ No human needed in optimization process

✗ Costly & Lengthy search process

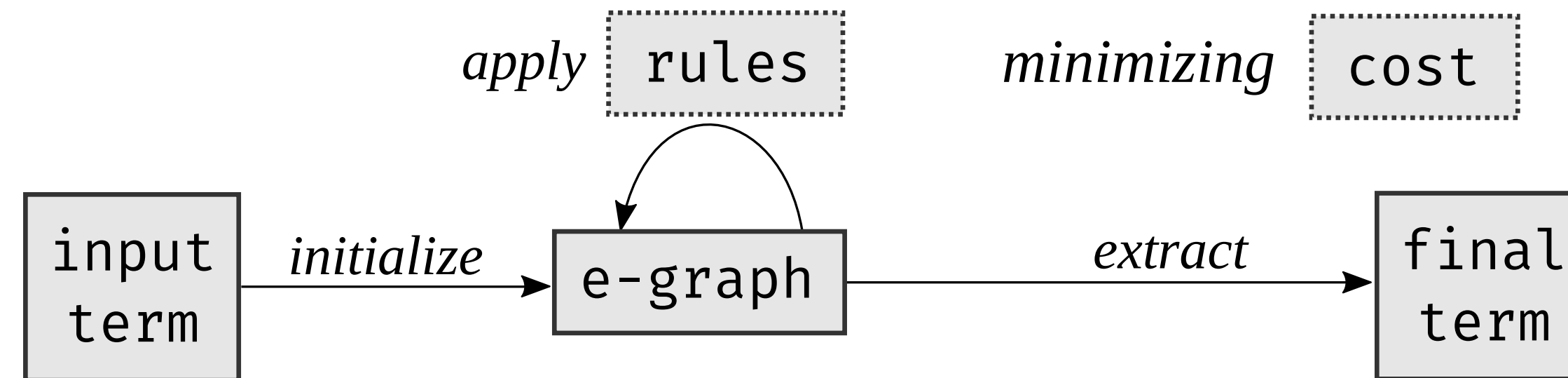
✗ Does not (yet) scale to all programs

Extensive human effort needed ✗

Human is in control, no search required ✓

Strategies are too sensitive  
⇒ don't scale across applications ✗

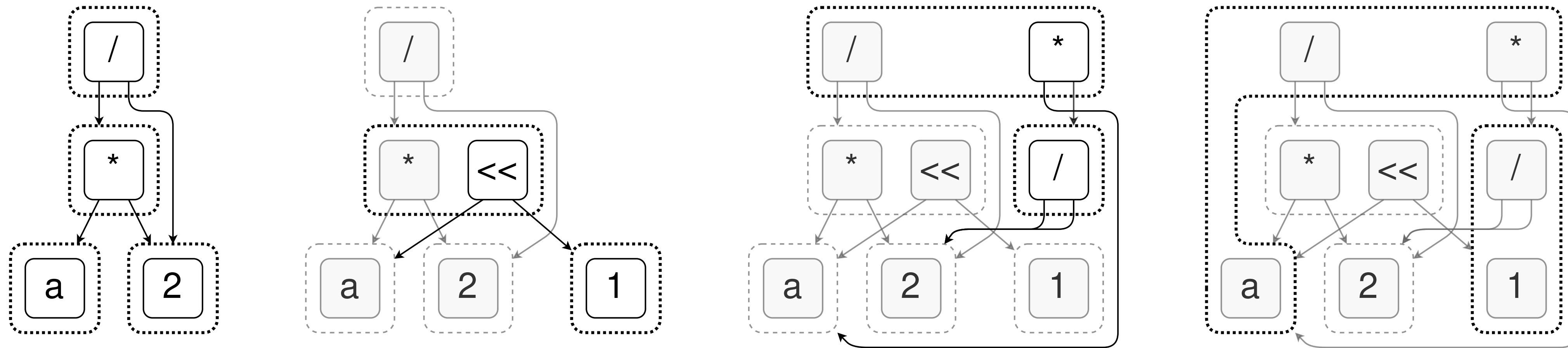
# Equality Saturation



- ▶ Optimize programs by efficiently exploring many possible rewrites
- ▶ Many successful applications sparked from the recent egg library

# E-Graph

Expression  
 $(a*2)/2$



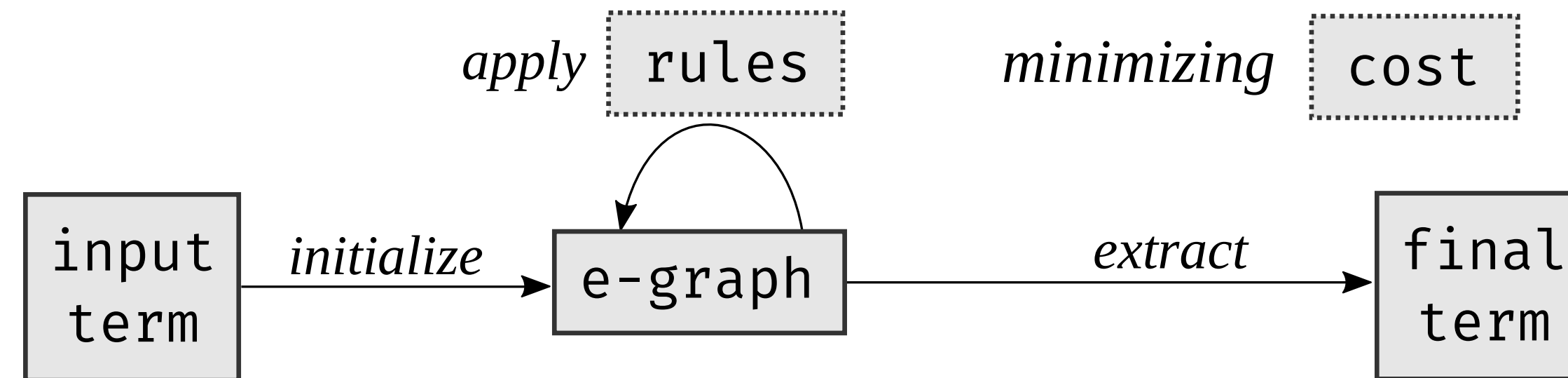
$$x*2 \rightarrow x<<1$$

$$(x*y)/z \rightarrow x*(y/z)$$

$$\begin{aligned} x/x &\rightarrow 1 \\ 1*x &\rightarrow x \end{aligned}$$

After applying Rewrites

# Equality Saturation



- ▶ Optimize programs by efficiently exploring many possible rewrites
- ▶ Many successful applications sparked from the recent egg library

**Some optimizations remain out of reach as the e-graph grows too big**



# Case Study

## Matrix Multiplication Optimizations for CPU:

- ▶ transform loops  
*blocking, permutation, unrolling*
- ▶ change data layout
- ▶ add parallelism  
*vectorization, multi-threading*

# Case Study

## Matrix Multiplication Optimizations for CPU:

- ▶ transform loops  
*blocking, permutation, unrolling*
- ▶ change data layout
- ▶ add parallelism  
*vectorization, multi-threading*

**Space of equivalent programs to consider is huge**

# Case Study

- Rewritten language: **RISE**, a functional array language

## Matrix Multiplication in **RISE**:

---

```
def mm a b =  
  map (λaRow.  
    map (λbCol.  
      dot aRow bCol)  
    (transpose b)) a  
  | for aRow in a:  
  |   for bCol in transpose(b):  
  |     ... = dot(aRow, bCol)  
  
def dot xs ys =  
  reduce + 0  
  (map (λ(x, y). x × y)  
  (zip xs ys))  
  | for (x, y) in zip(xs, ys):  
  |   acc += x × y
```

---

# Case Study

- Rewritten language: **RISE**, a functional array language

## Matrix Multiplication in **RISE**:

---

```
def mm a b =  
  map (λaRow.  
    map (λbCol.  
      dot aRow bCol)  
      (transpose b)) a  
  | for aRow in a:  
  |   for bCol in transpose(b):  
  |     ... = dot(aRow, bCol)  
  
def dot xs ys =  
  reduce + 0  
  (map (λ(x, y). x × y)  
    (zip xs ys))  
  | for (x, y) in zip(xs, ys):  
  |   acc += x × y
```

---

**RISE is designed for optimization via term rewriting**

# Case Study

- ▶ Achieve the same 7 optimization goals with equality saturation?<sup>1</sup>

goal	found?	runtime	RAM
<i>baseline</i>	✓	0.5s	0.02 GB
<i>blocking</i>	✓	>1h	35 GB
<i>vectorization</i>	✗	>1h	>60 GB
<i>loop-perm</i>	✗	>1h	>60 GB
<i>array-packing</i>	✗	35mn	>60 GB
<i>cache-blocks</i>	✗	35mn	>60 GB
<i>parallel</i>	✗	35mn	>60 GB

- ▶ Most goals are not found before exhausting 60 GB.
- ▶ For comparison, rewriting strategies take <2s and <1GB.

---

<sup>1</sup>on Intel Xeon E5-2640 v2

# Case Study

- ▶ Achieve the same 7 optimization goals with equality saturation?<sup>1</sup>

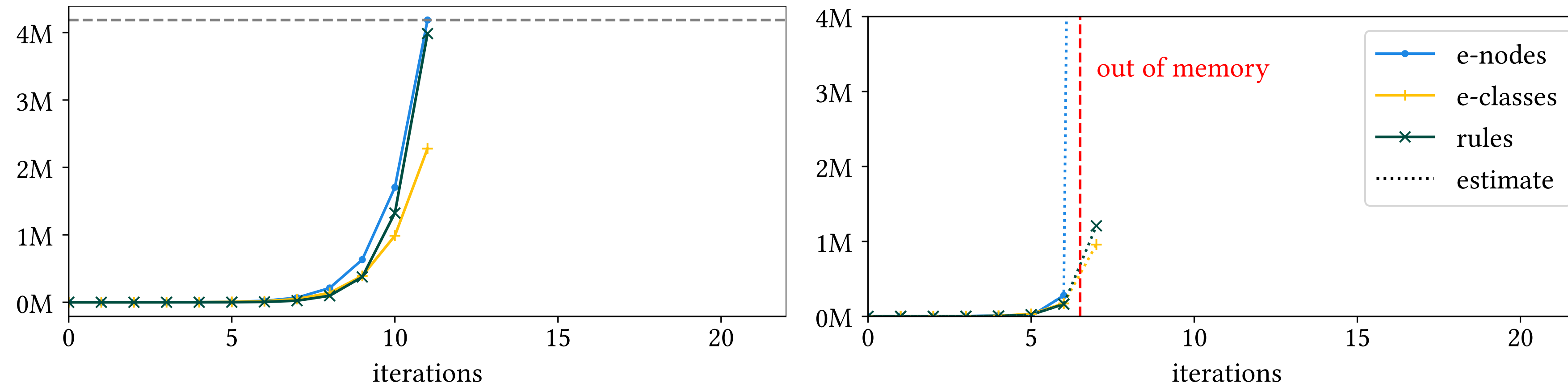
goal	found?	runtime	RAM
<i>baseline</i>	✓	0.5s	0.02 GB
<i>blocking</i>	✓	>1h	35 GB
<i>vectorization</i>	✗	>1h	>60 GB
<i>loop-perm</i>	✗	>1h	>60 GB
<i>array-packing</i>	✗	35mn	>60 GB
<i>cache-blocks</i>	✗	35mn	>60 GB
<i>parallel</i>	✗	35mn	>60 GB

**Standard equality saturation does not scale to this optimization space**

---

<sup>1</sup>on Intel Xeon E5-2640 v2

# E-Graph Evolution



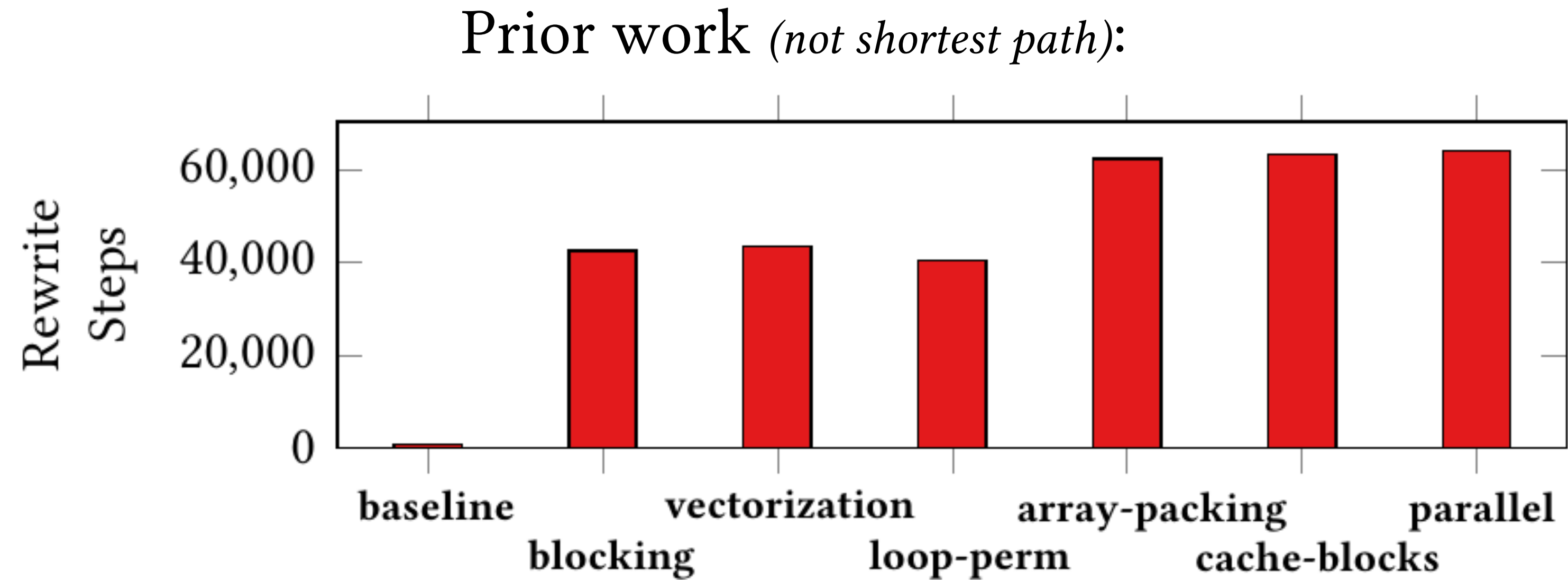
(a) *blocking*, found: ✓

(b) *parallel*, found: ✗

## Two difficulties:

1. Long rewrite sequences  $\implies$  many iterations are required
2. Explosive combination of rewrite rules  $\implies$  exponential growth
  - ▶ millions of e-nodes and e-classes in less than 10 iterations
  - ▶ worse for *parallel*, memory is exhausted in the 7th iteration

# Difficulty 1. Long Rewrite Sequences





# Difficulty 2. Explosive Combinations of Rewrite Rules

Two example rules that quickly generate many possibilities:

*split-join:*

---

<code>map f x</code>	<b>for</b> m:
	... = f(...)
$\mapsto$	
<code>join</code>	<b>for</b> m / n:
<code>(map</code>	<b>for</b> n:
<code>  (map f)</code>	... = f(...)
<code>  (split n x))</code>	

---

*transpose-around-map-map:*

---

<code>map</code>	<b>for</b> m:
<code>  (map f) x</code>	<b>for</b> n:
	... = f(...)
$\mapsto$	
<code>transpose</code>	<b>for</b> n:
<code>(map</code>	<b>for</b> m:
<code>  (map f)</code>	... = f(...)
<code>  (transpose x))</code>	

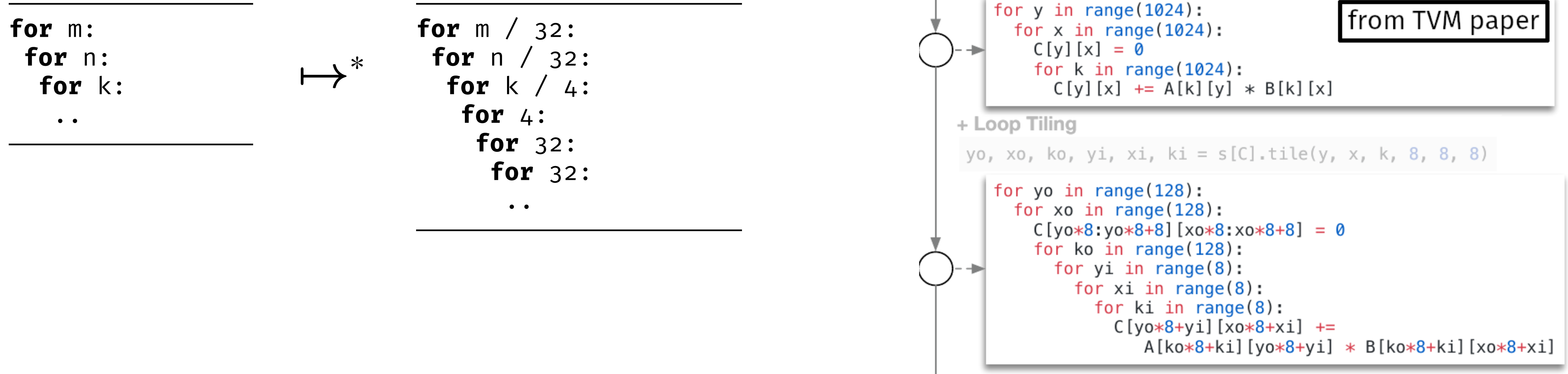
---

To overcome these difficulties, we came up with *sketch-guided equality saturation*

# Sketch-Guided Equality Saturation

Observation:

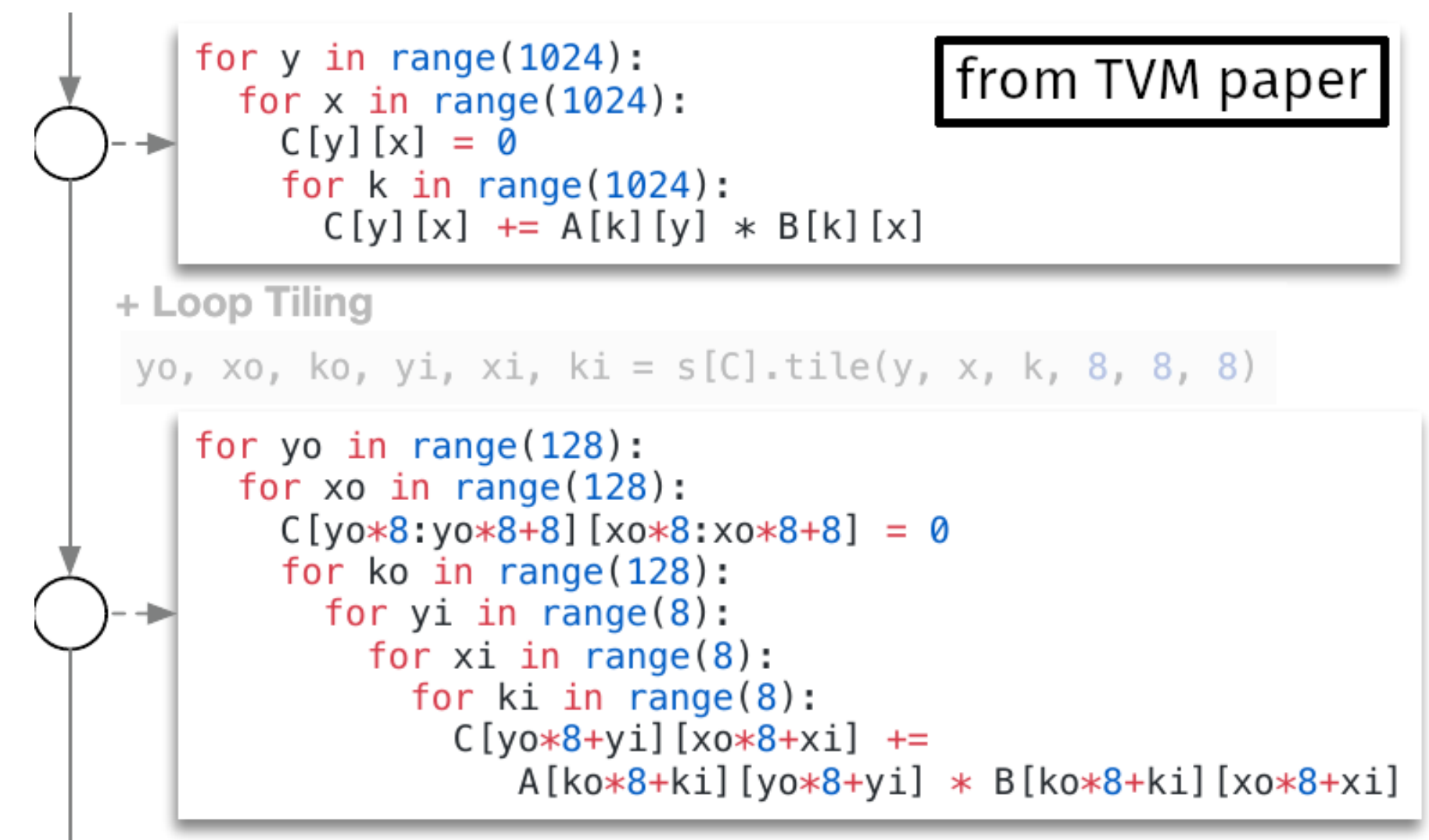
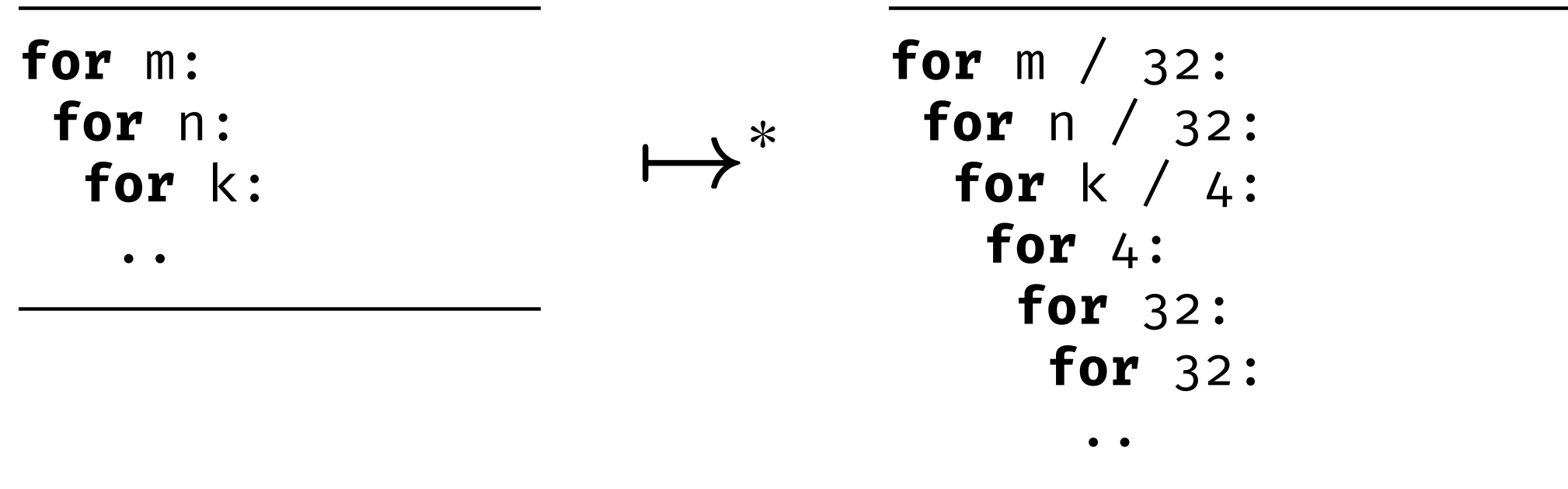
- The *shape* of the optimised program is often used to explain optimizations:



# Sketch-Guided Equality Saturation

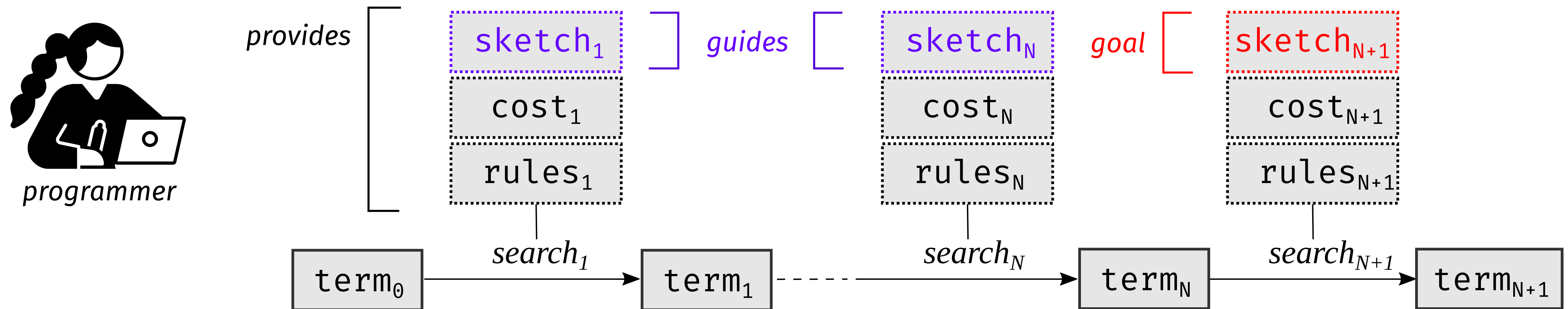
Observation:

- The *shape* of the optimised program is often used to explain optimizations:



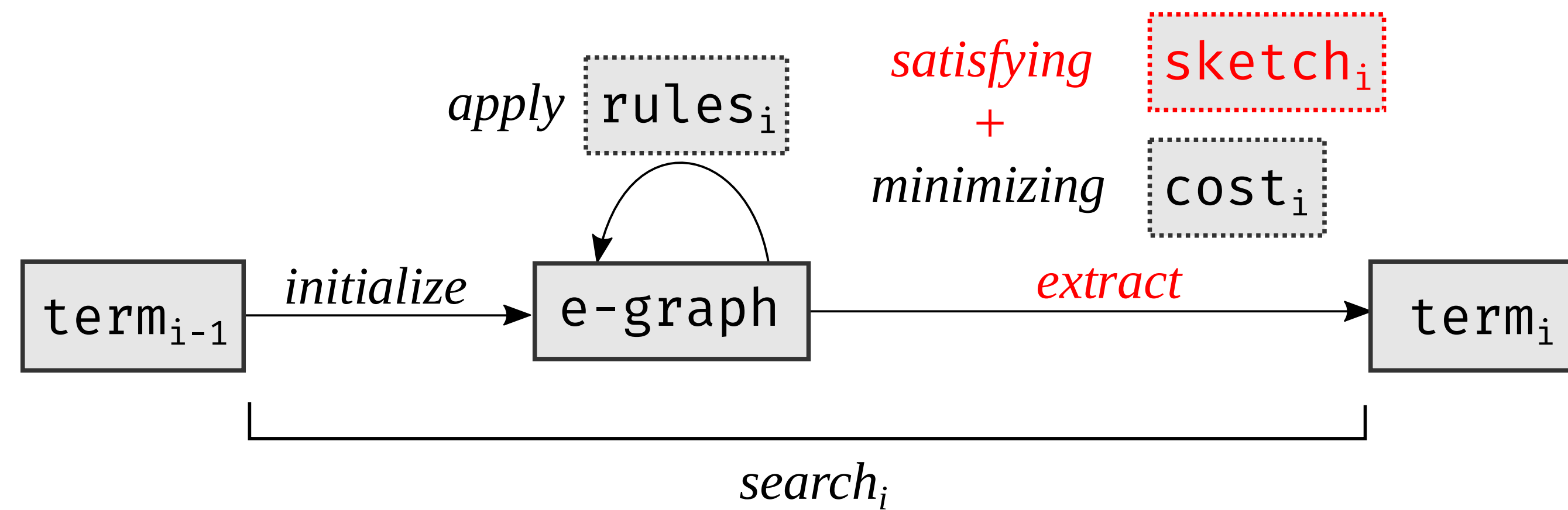
**Explanatory shapes can be formalized as sketches and used to guide rewriting**

# Sketch-Guided Equality Saturation



- Factors an unfeasible search into a sequence of feasible ones:
  1. Break long rewrite sequences
  2. Isolate explosive combinations of rewrite rules

# Sketch-Satisfying Equality Saturation



- ▶ Terminates as soon as a program satisfying the sketch is found

# Sketches

- ▶ *Sketches* are program patterns that leave details unspecified

*baseline* sketch:

---

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

---

- ▶ Abstractions defined in terms of smaller building blocks:

---

```
def containsAddMul: Sketch =  
  contains(app(app(+, ?), contains(×)))
```

---

# Sketches

- ▶ *Sketches* are program patterns that leave details unspecified

*baseline* sketch:

---

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

---

- ▶ A sketch  $s$  is satisfied by a set of terms  $R(s)$ :

---

```
def containsAddMul: Sketch =
  contains(app(app(+, ?), contains(x)))

R(containsAddMul) = { R(app(app(+, ?), contains(x))) } ∪
  { F(t1, .., tn) | ∃ti ∈ R(containsAddMul) }
R(app(app(+, ?), contains(x))) = { app(app(+, t1), t2) | t2 ∈ R(contains(x)) }
R(contains(x)) = { x } ∪ { F(t1, .., tn) | ∃ti ∈ R(contains(x)) }
```

---



# Sketches

- ▶ *Sketches* are program patterns that leave details unspecified

*baseline* sketch:

---

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

---

*blocking* sketch:

---

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

---

# Sketches

- ▶ *Sketches* are program patterns that leave details unspecified

*baseline sketch:*

---

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

---

*sketch guide:*

*how to split the loops before reordering them?*

---

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

---

*blocking sketch:*

---

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

---

# Evaluation

- ▶ Equality Saturation without Sketch Guides<sup>2</sup>:

goal	found?	runtime	RAM
<i>baseline</i>	✓	0.5s	0.02 GB
<i>blocking</i>	✓	>1h	35 GB
+ 5 others	✗	>35mn	>60 GB

- ▶ Sketch-Guided Equality Saturation<sup>3</sup>:

goal	sketch guides	found?	runtime	RAM
<i>baseline</i>	0	✓	0.5s	0.02 GB
<i>blocking</i>	1	✓	7s	0.3 GB
+ 5 others	2-3	✓	≤7s	≤0.5 GB

---

<sup>2</sup>Intel Xeon E5-2640 v2

<sup>3</sup>AMD Ryzen 5 PRO 2500U

# Evaluation

- ▶ Equality Saturation without Sketch Guides<sup>2</sup>:

goal	found?	runtime	RAM
<i>baseline</i>	✓	0.5s	0.02 GB
<i>blocking</i>	✓	>1h	35 GB
+ 5 others	✗	>35mn	>60 GB

- ▶ Sketch-Guided Equality Saturation<sup>3</sup>:

goal	sketch guides	found?	runtime	RAM
<i>baseline</i>	0	✓	0.5s	0.02 GB
<i>blocking</i>	1	✓	7s	0.3 GB
+ 5 others	2-3	✓	≤7s	≤0.5 GB

**Sketch-guided equality saturation finds all 7 optimization goals**

---

<sup>2</sup>Intel Xeon E5-2640 v2

<sup>3</sup>AMD Ryzen 5 PRO 2500U

# Evaluation

► Equality Saturation without Sketch Guides<sup>2</sup>:

goal	found?	runtime	RAM
<i>baseline</i>	✓	0.5s	0.02 GB
<i>blocking</i>	✓	>1h	35 GB
+ 5 others	✗	>35mn	>60 GB

► Sketch-Guided Equality Saturation<sup>3</sup>:

goal	sketch guides	found?	runtime	RAM
<i>baseline</i>	0	✓	0.5s	0.02 GB
<i>blocking</i>	1	✓	7s	0.3 GB
+ 5 others	2-3	✓	≤7s	≤0.5 GB

582x

116x

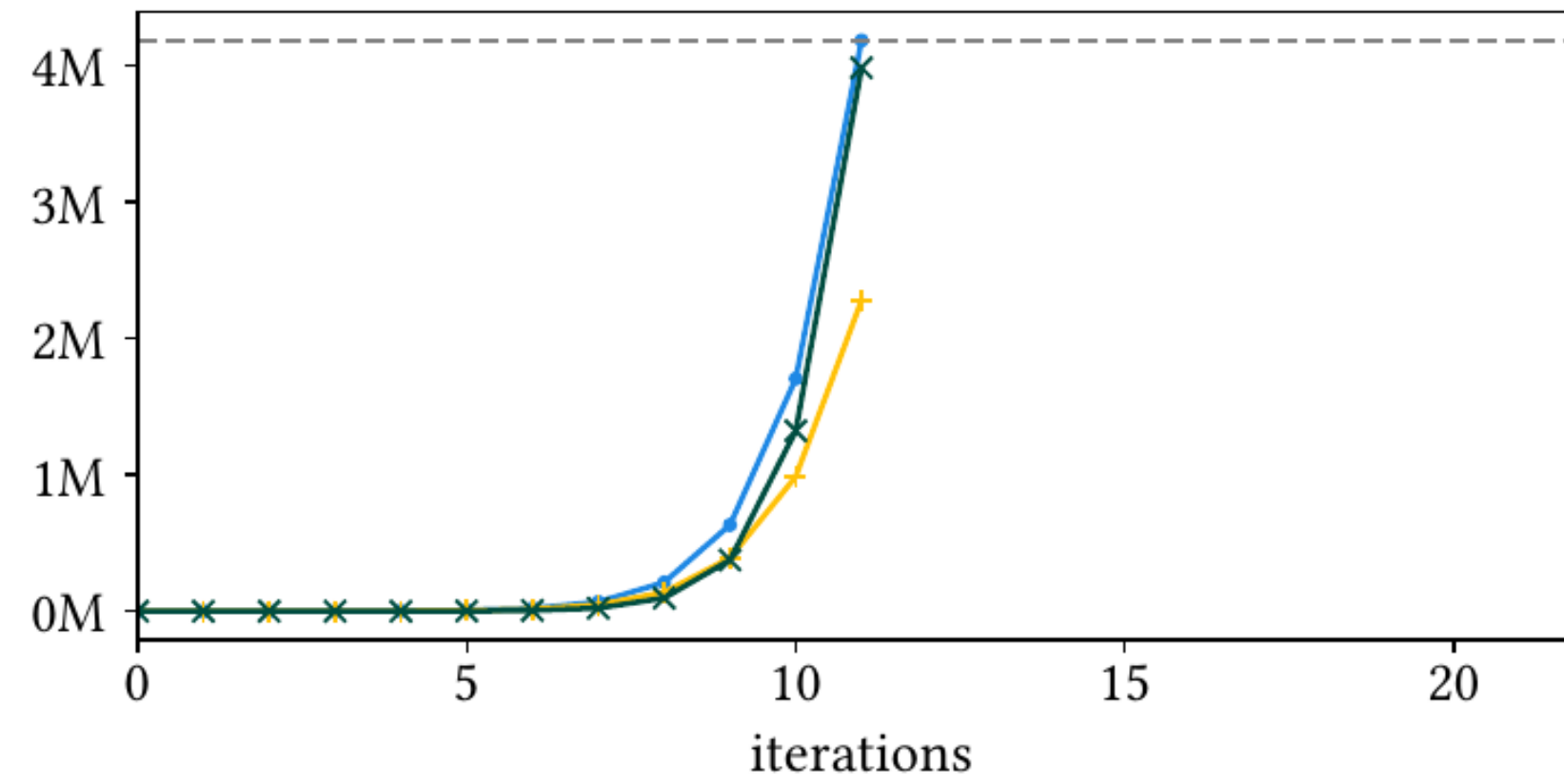
---

<sup>2</sup>Intel Xeon E5-2640 v2

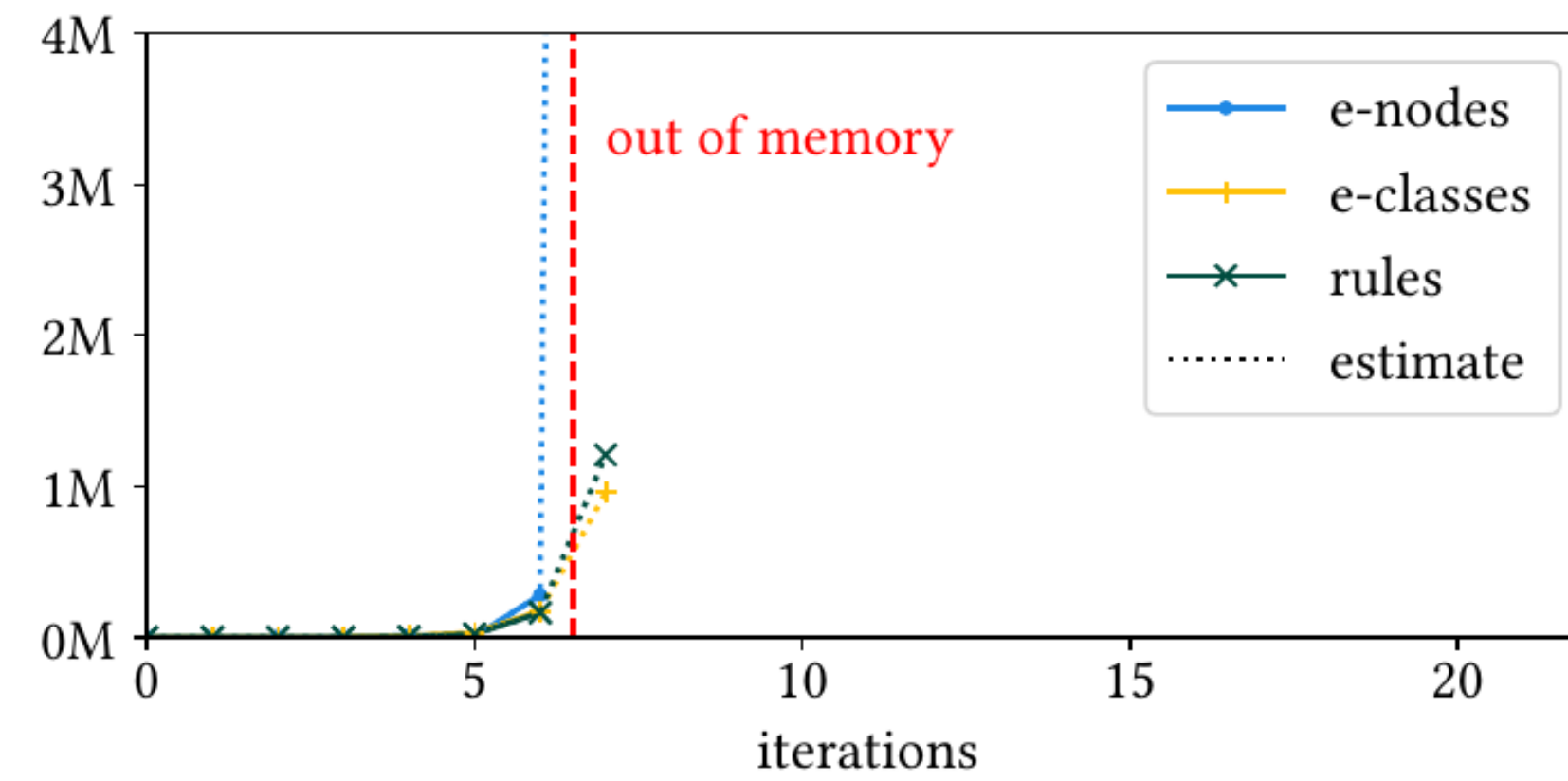
<sup>3</sup>AMD Ryzen 5 PRO 2500U

# Evaluation

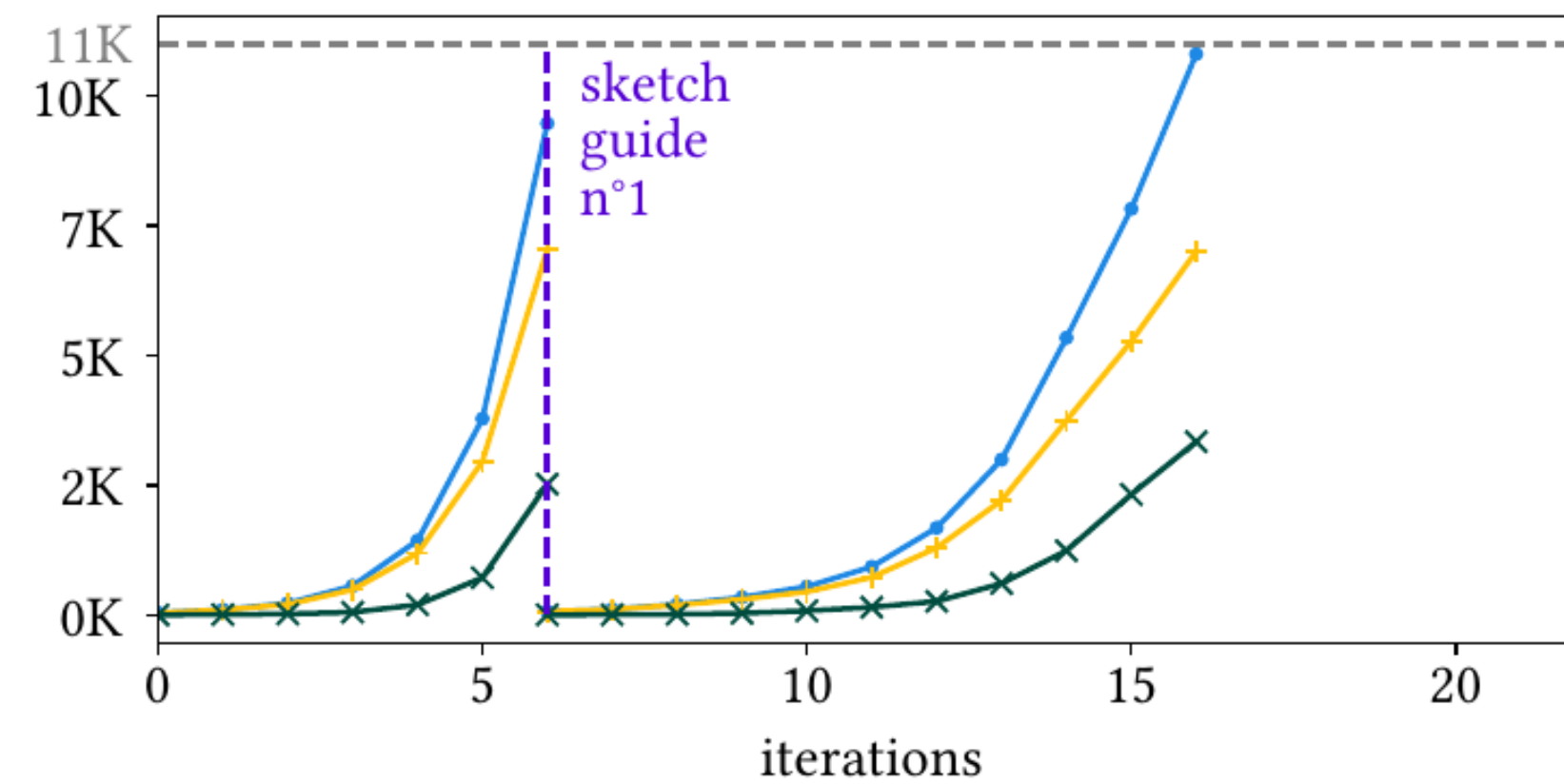
## E-Graph Evolution



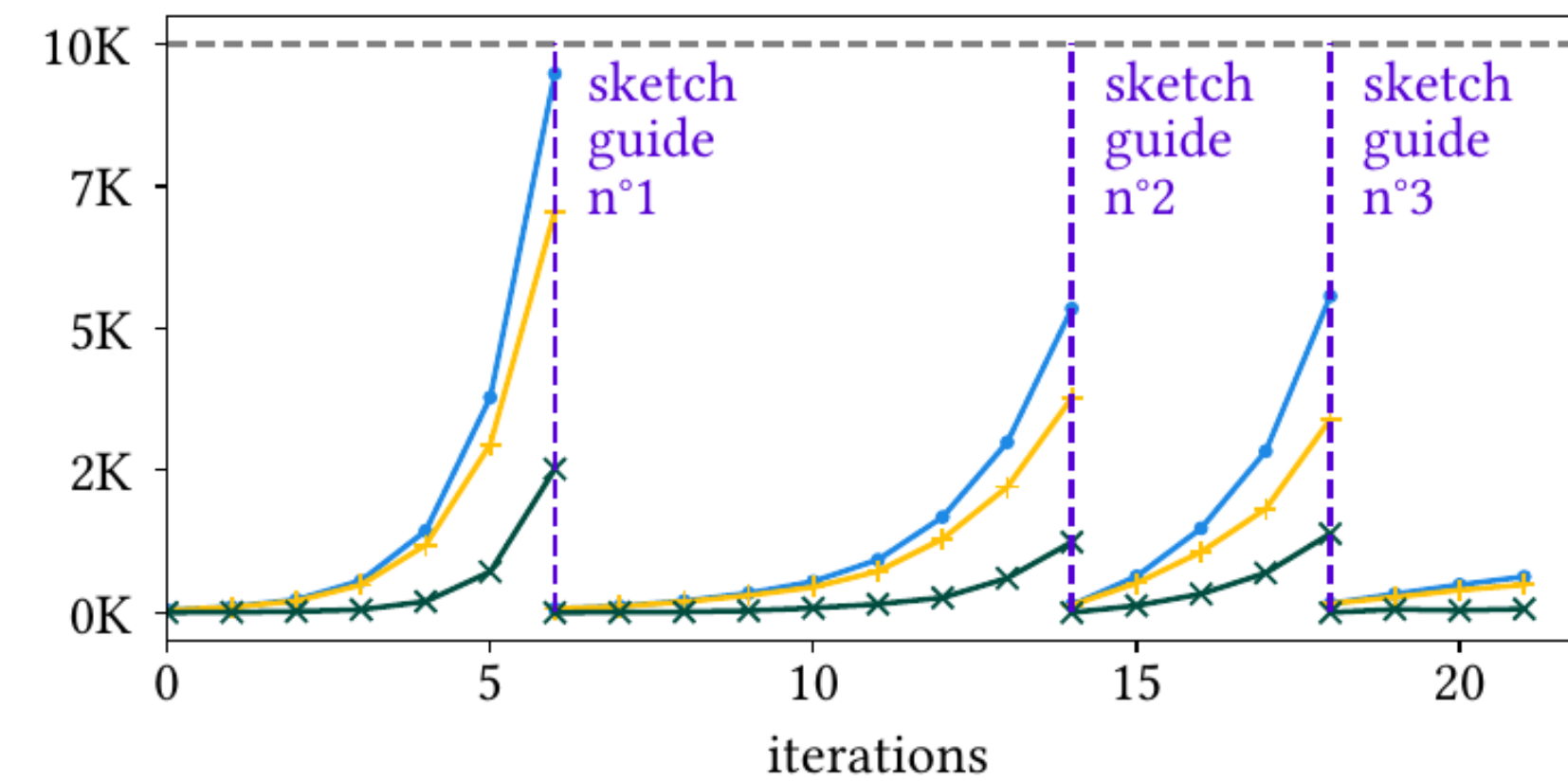
(a) unguided, *blocking*, found: ✓



(b) unguided, *parallel*, found: ✗



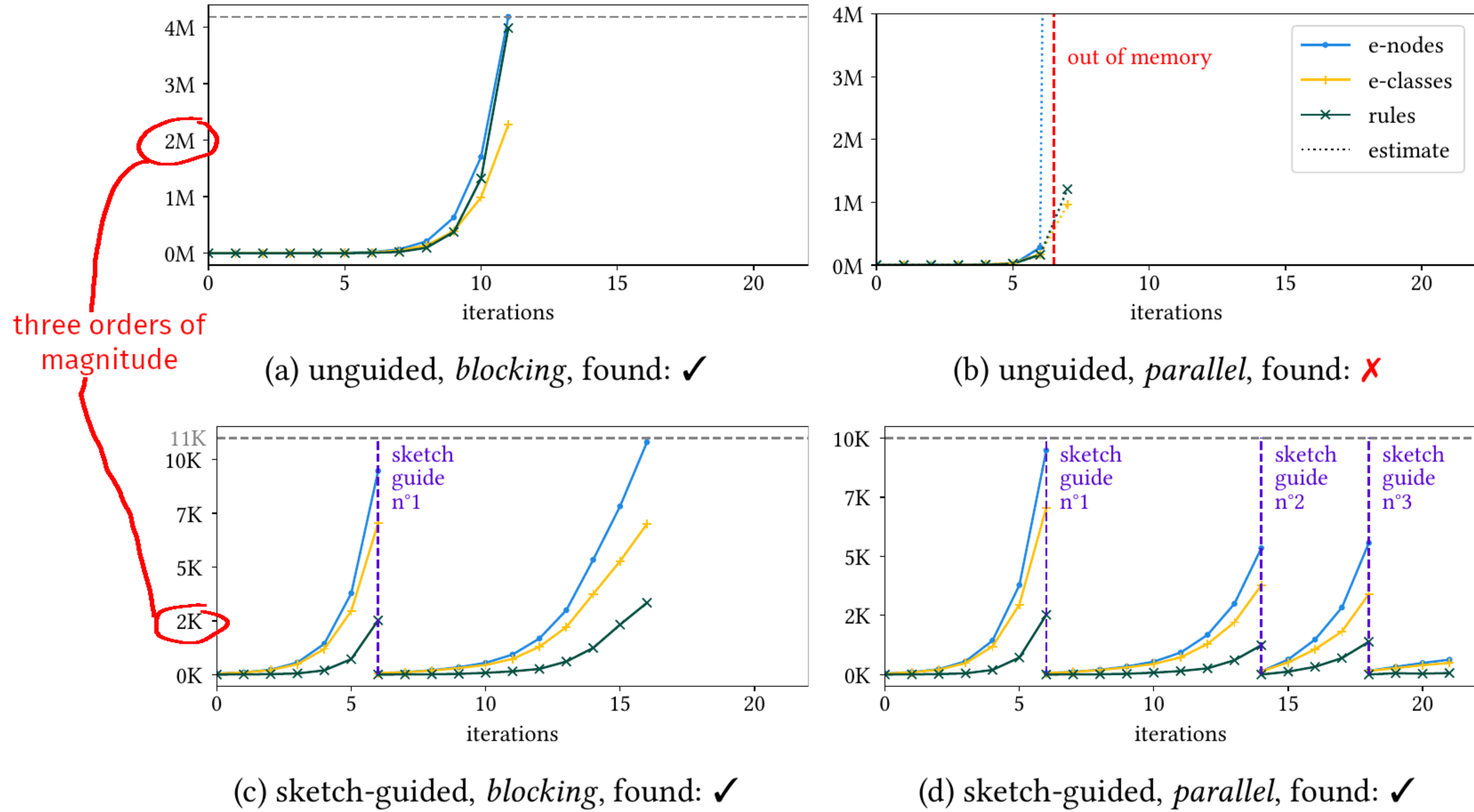
(c) sketch-guided, *blocking*, found: ✓



(d) sketch-guided, *parallel*, found: ✓

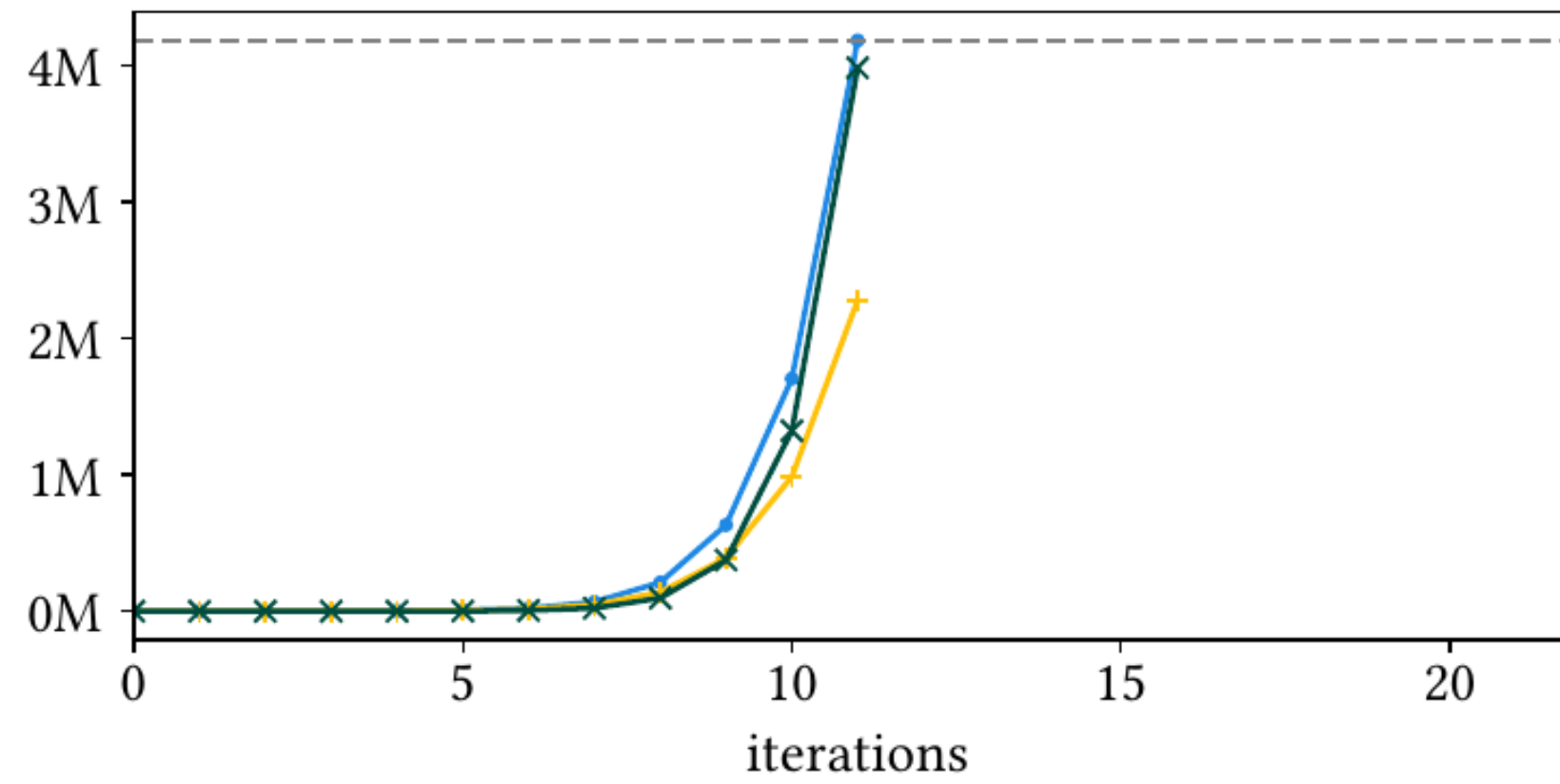
# Evaluation

## E-Graph Evolution

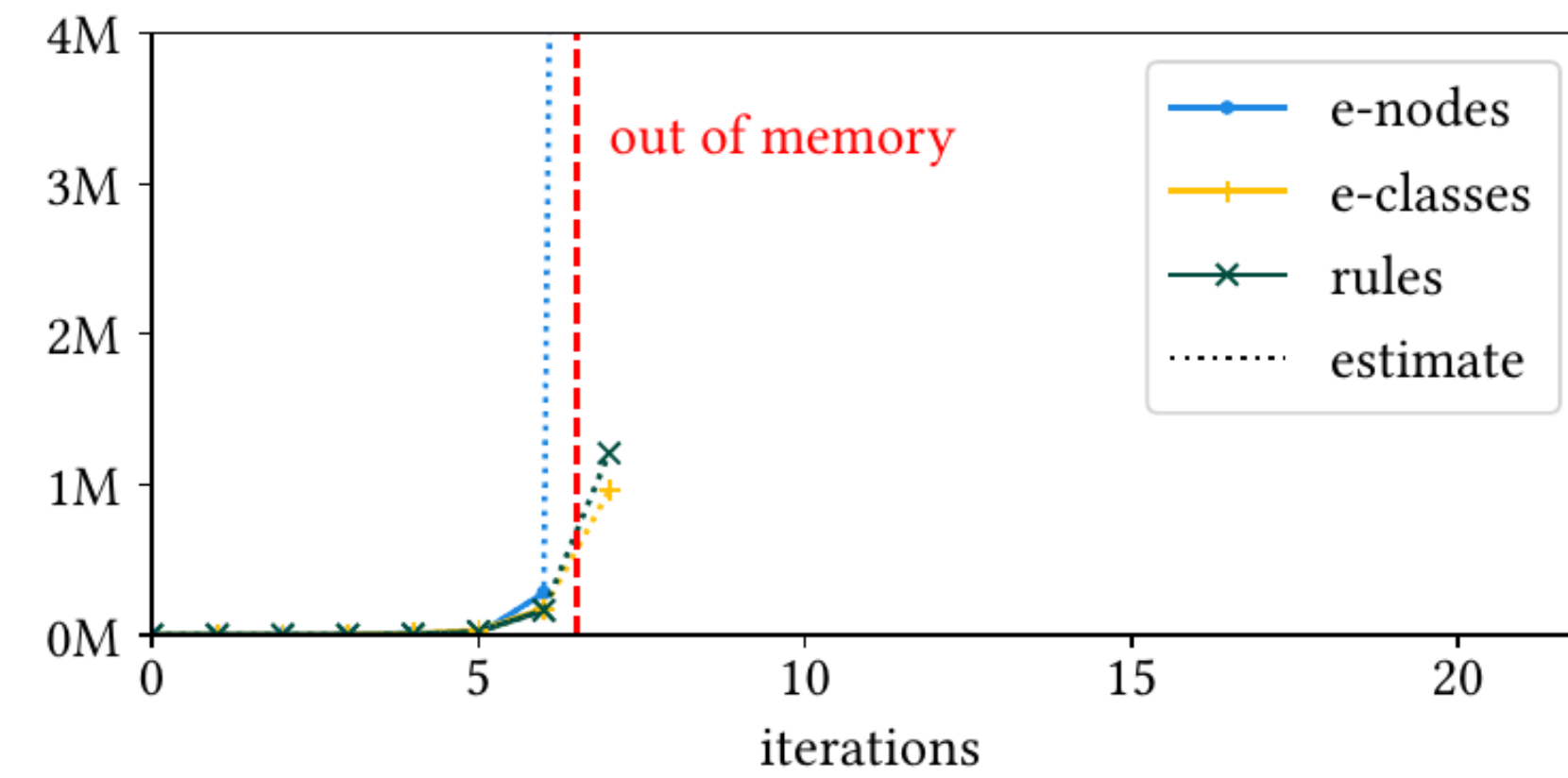


# Evaluation

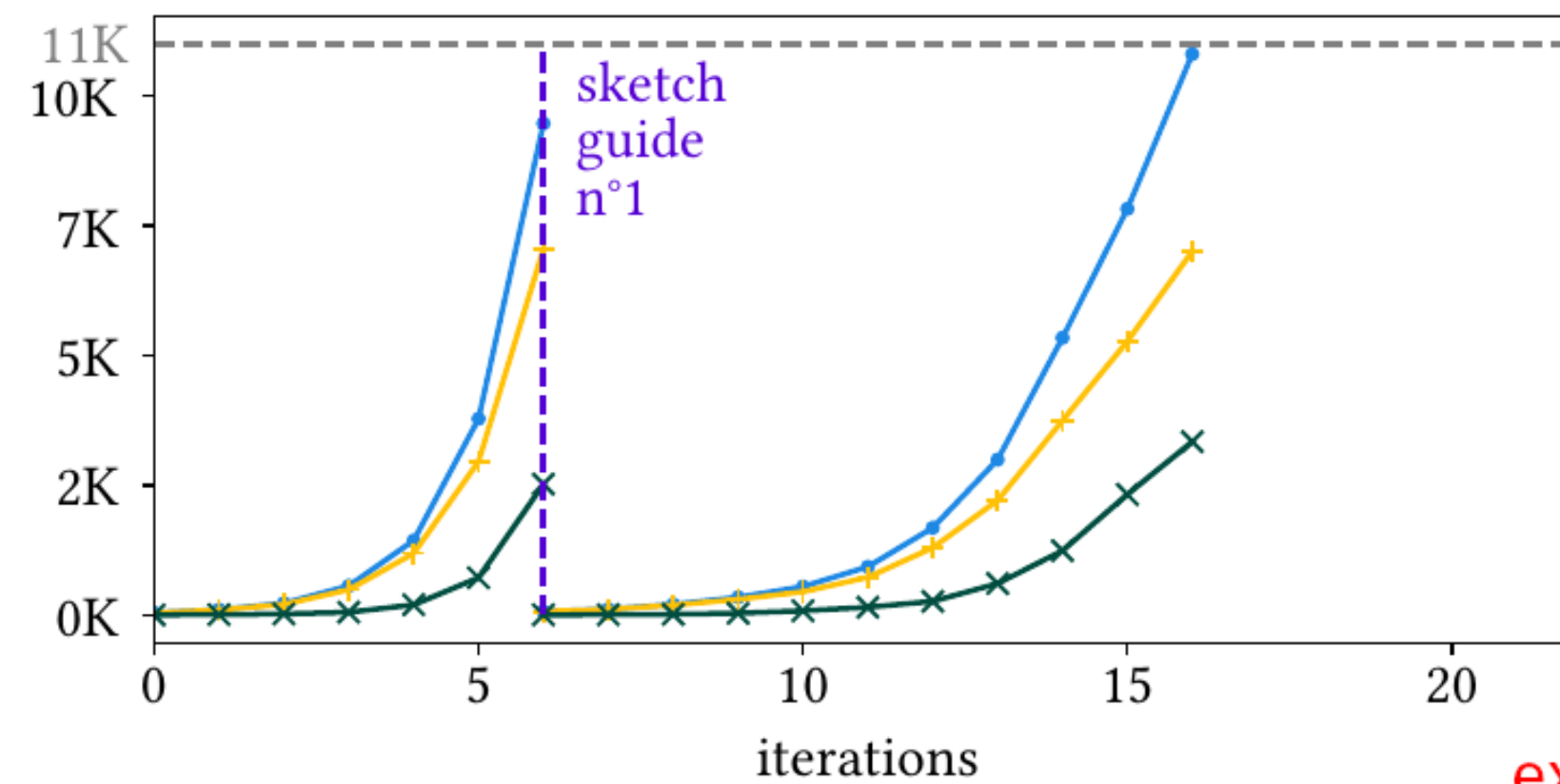
## E-Graph Evolution



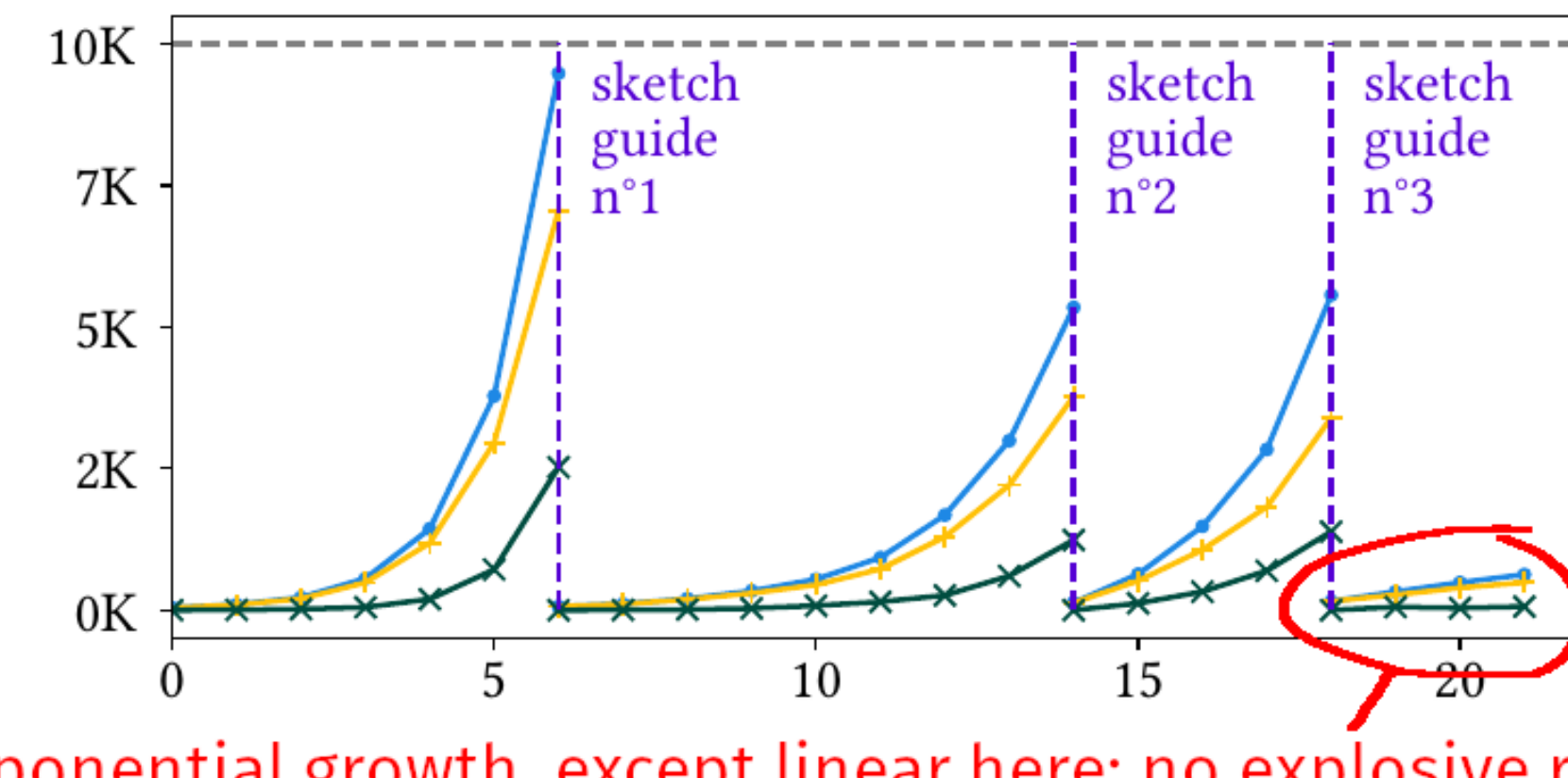
(a) unguided, *blocking*, found: ✓



(b) unguided, *parallel*, found: ✗



(c) sketch-guided, *blocking*, found: ✓



(d) sketch-guided, *parallel*, found: ✓

exponential growth, except linear here: no explosive rewrites



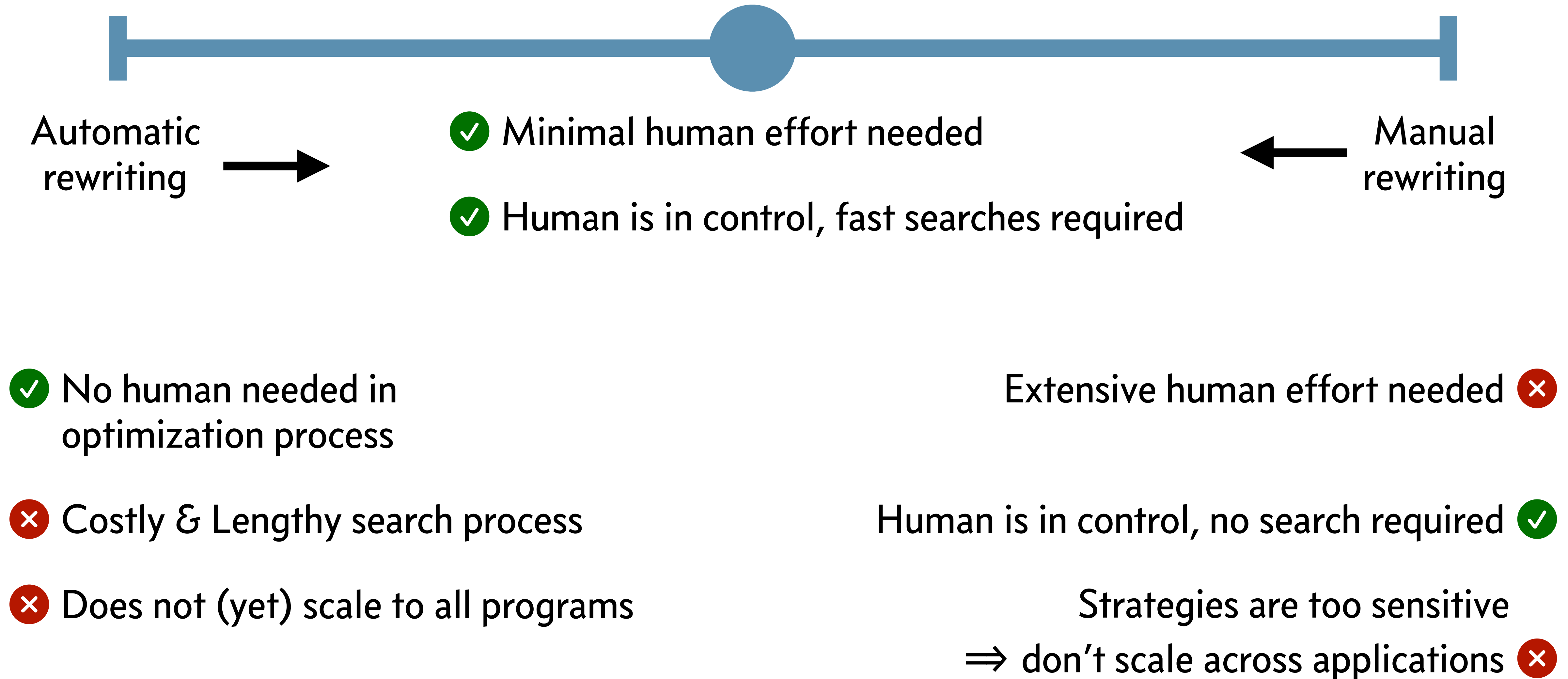
# Evaluation

## Sketches vs Full Program

goal	sketch guides	sketch goal	sketch sizes	program size
<i>blocking</i>	<i>split</i>	<i>reorder<sub>1</sub></i>	7	90
<i>vectorization</i>	<i>split + reorder<sub>1</sub></i>	<i>lower<sub>1</sub></i>	7	124
<i>loop-perm</i>	<i>split + reorder<sub>2</sub></i>	<i>lower<sub>2</sub></i>	7	104
<i>array-packing</i>	<i>split + reorder<sub>2</sub> + store</i>	<i>lower<sub>3</sub></i>	7-12	121
<i>cache-blocks</i>	<i>split + reorder<sub>2</sub> + store</i>	<i>lower<sub>4</sub></i>	7-12	121
<i>parallel</i>	<i>split + reorder<sub>2</sub> + store</i>	<i>lower<sub>5</sub></i>	7-12	121

- ▶ each sketch corresponds to a logical transformation step
- ▶ sketches elide around 90% of the program
- ▶ intricate details such as array reshaping patterns are not specified (e.g. **split**, **join**, **transpose**)

# Tradeoffs when optimizing with rewriting



# Thanks to all the PhD students



Johannes  
Lenfers



Martin  
Lücke



Federico  
Pizzuti



Xueying  
Qin



Bastian  
Hagedorn



Thomas  
Kœhler



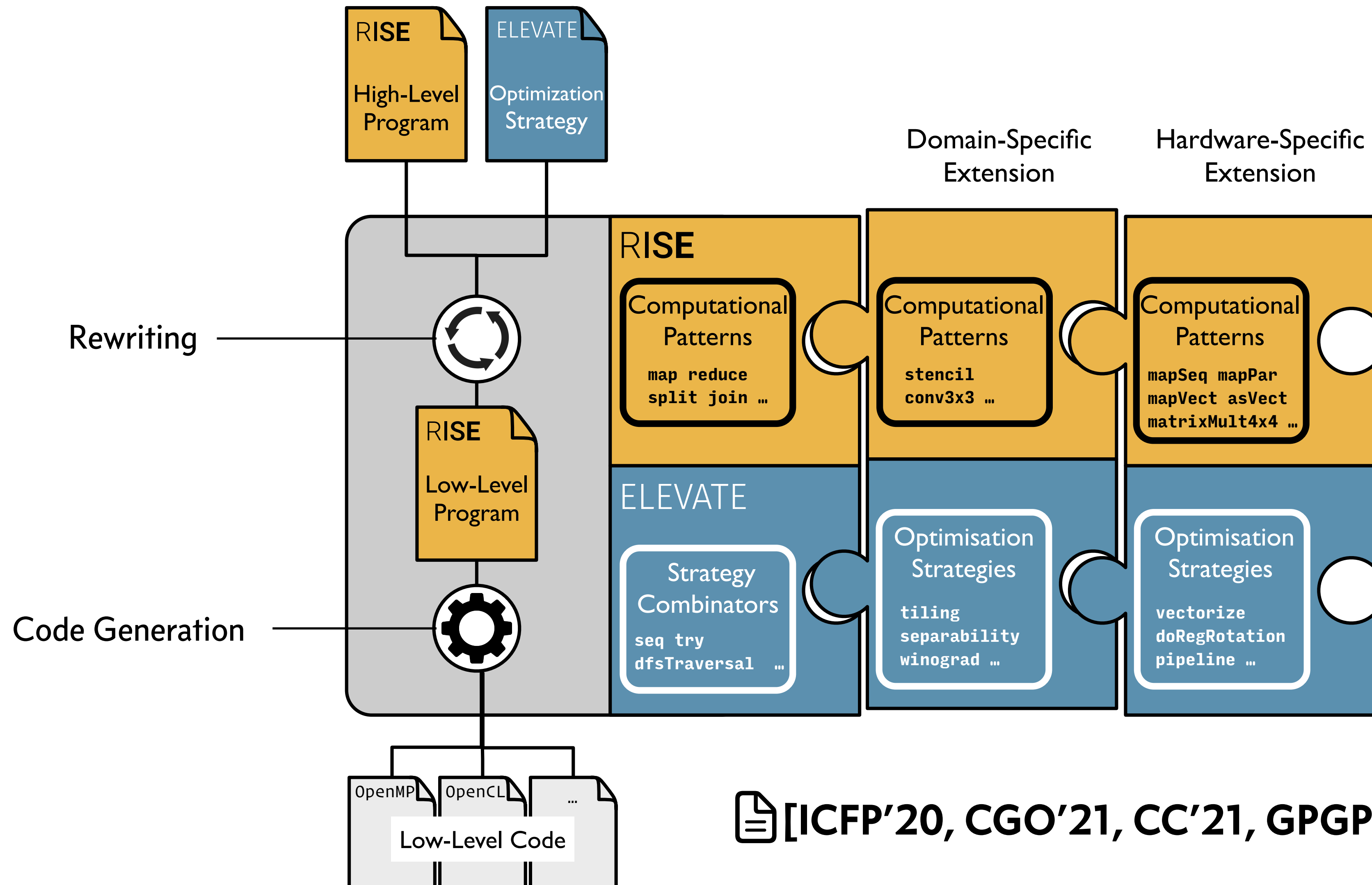
Rongxiao  
Fu



Bastian  
Köpcke

# How to design the next 700 optimizing compilers

Michel Steuwer



<https://michel.steuwer.info>

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

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

[ICFP'20, CGO'21, CC'21, GPGPU'22, arXiv'22]