

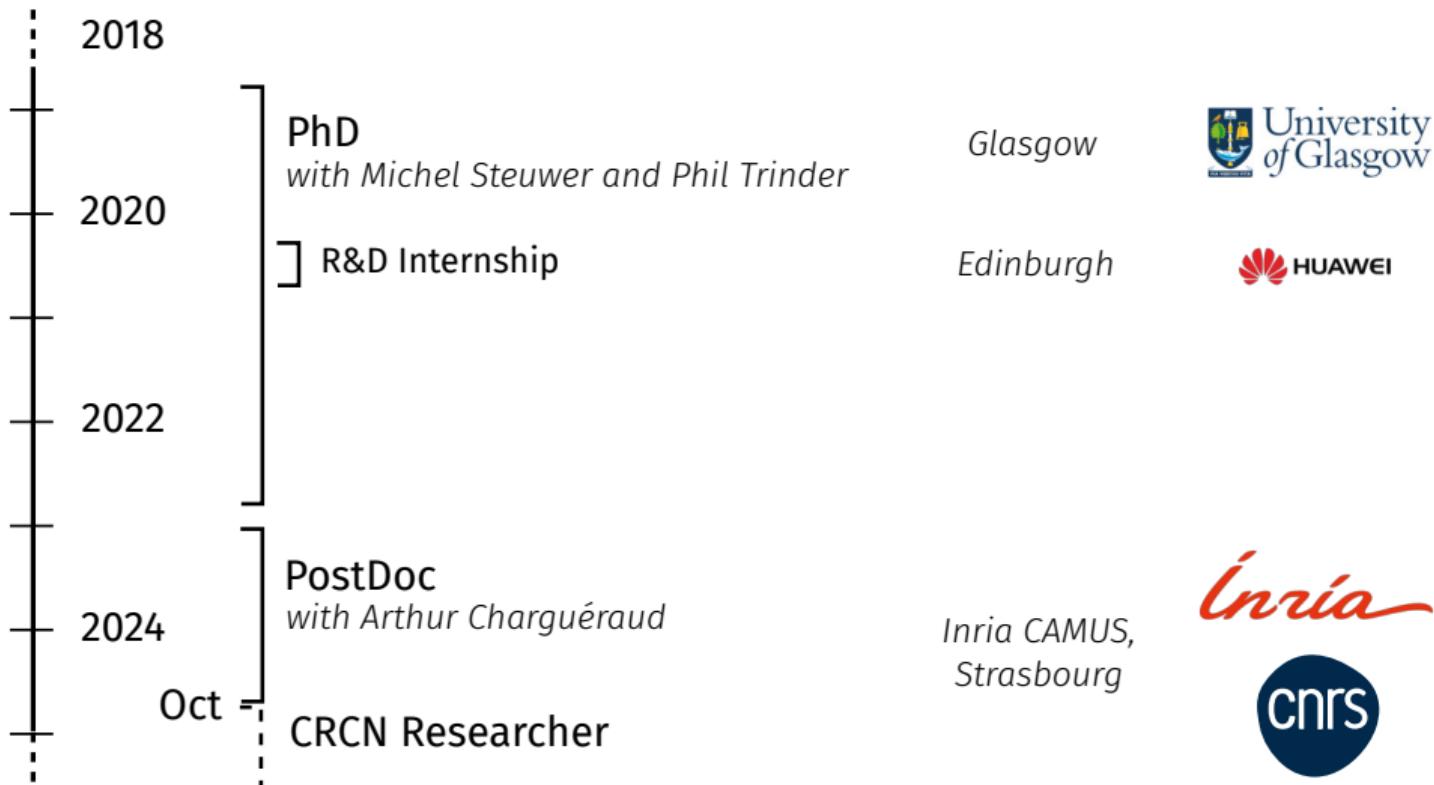
Towards Safe Interactive Optimization Across Layers

Thomas KŒHLER  thok.eu

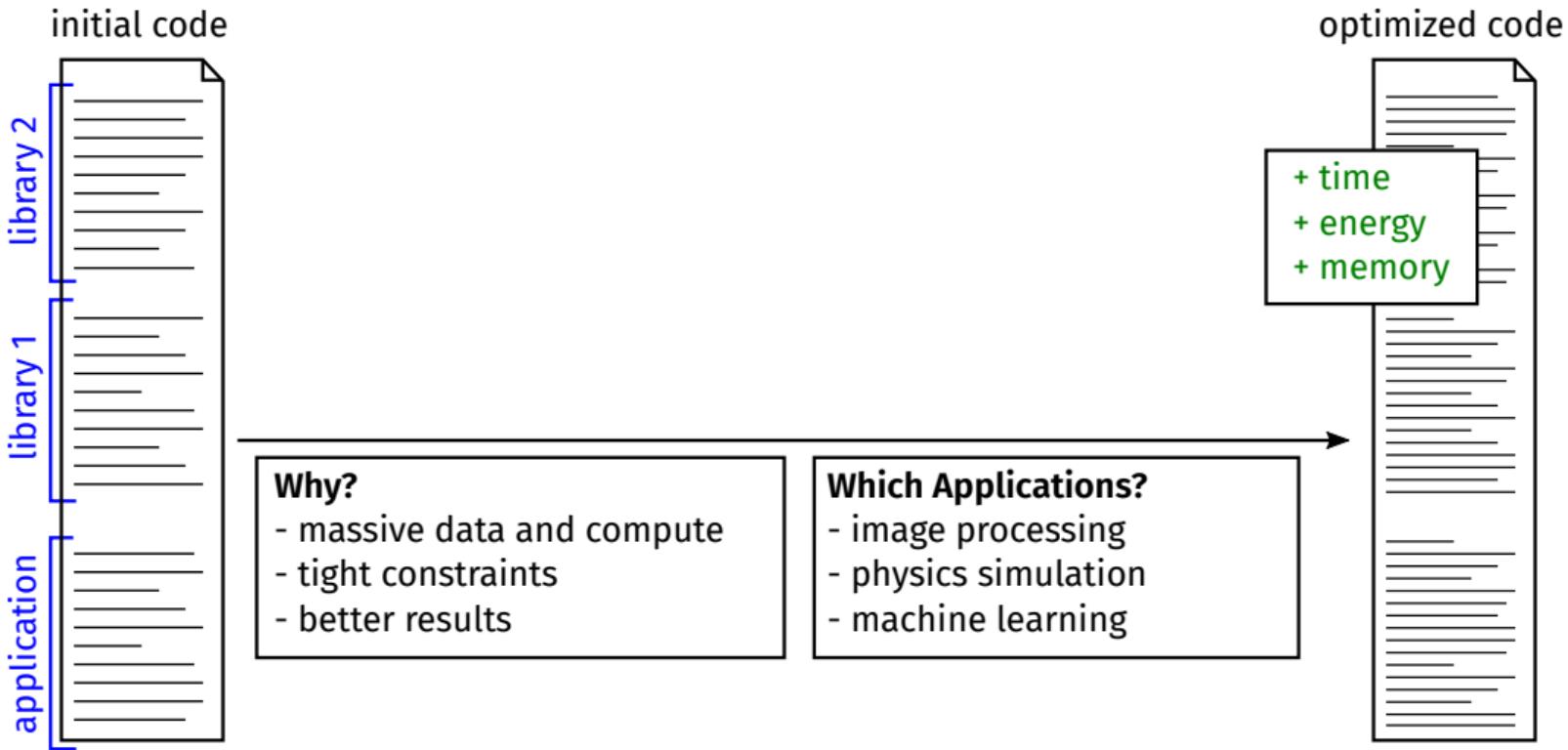


Séminaire CASH – July 2024, Lyon

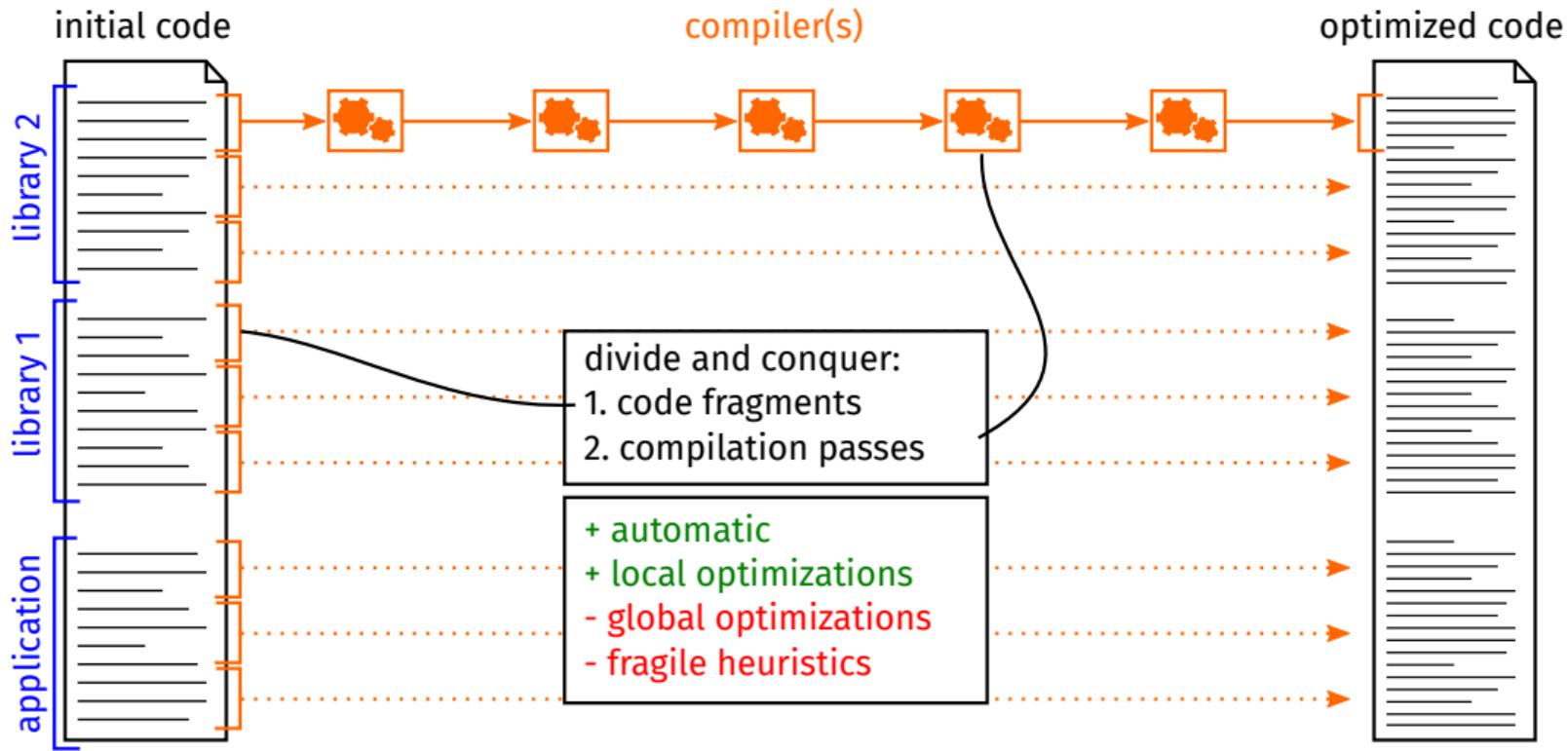
My Career



Optimizing Programs



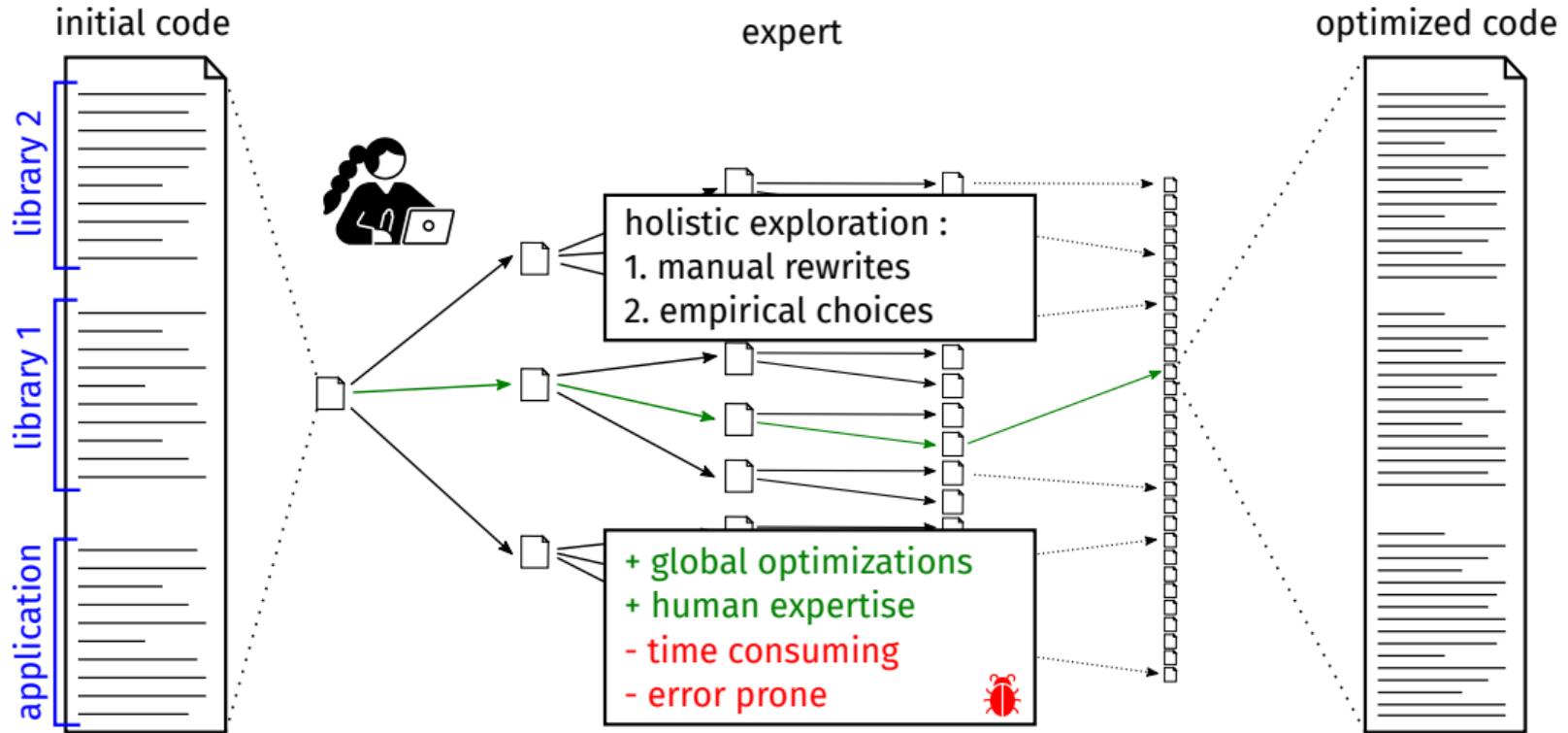
Automatic Program Optimization



Automatic Program Optimization



Manual Program Optimization



Manual Program Optimization

initial code

library 2
library 1
application

Matrix Multiplication, initial C code

```
for (int i = 0; i < m; i++) {  
    for (int j = 0; j < n; j++) {  
        float sum = 0.0f;  
        for (int k = 0; k < p; k++) {  
            sum += A[i][k] * B[k][j];  
        }  
        C[i][j] = sum;  
    }  
}
```

150× faster,
7× bigger

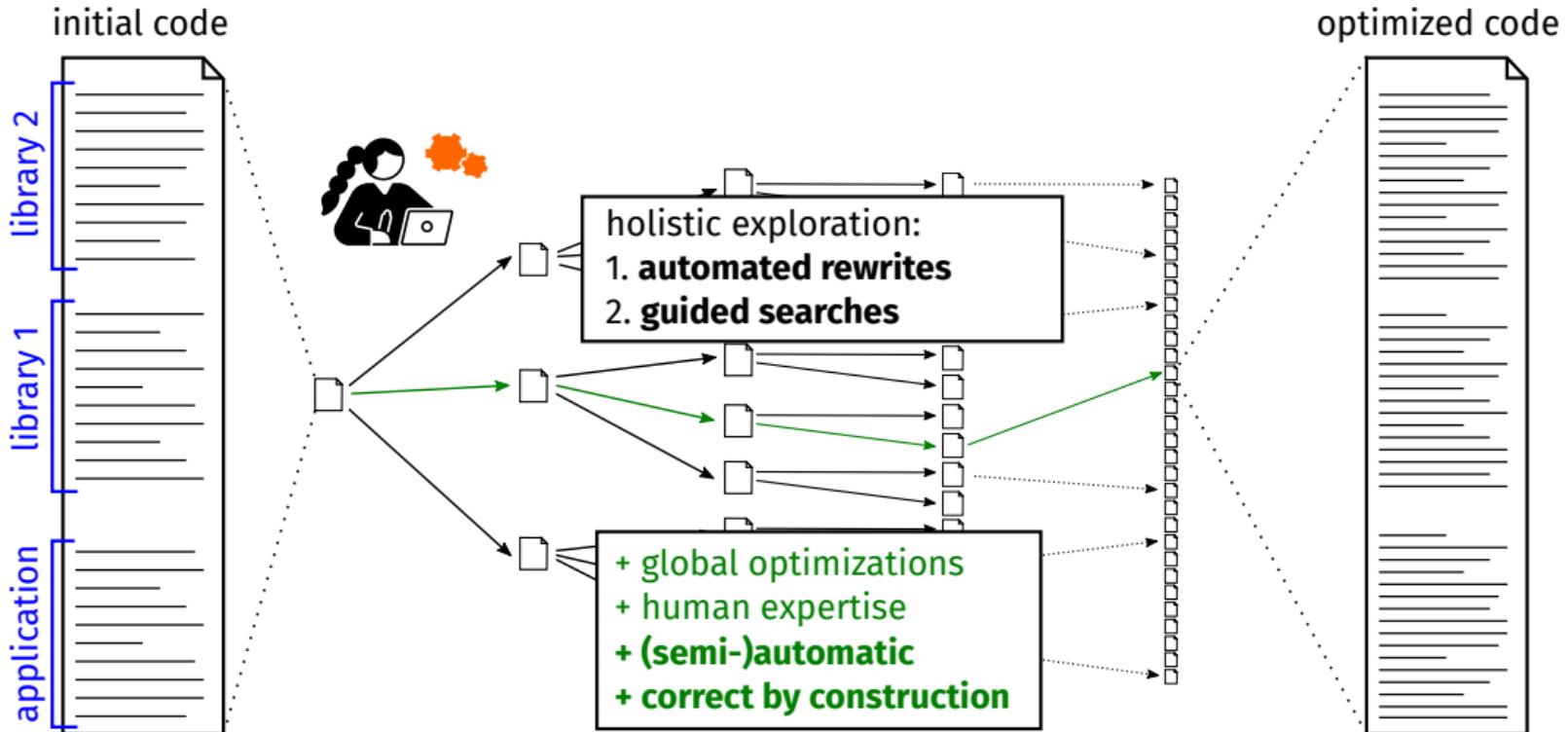
```
float* pB = (float*) malloc(sizeof(float)*32*32*4*32);  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum = 0.0f;  
                    for (int l = 0; l < 4; l++) {  
                        sum += A[i][b1 + 32 * b2 + 32 * k + l] * B[32 * b1 + (4 * b2 + k) + 32 * b3 + j];  
                    }  
                }  
            }  
        }  
    }  
} #pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        float* sum = (float*) malloc(sizeof(float)*32*32);  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                sum[32 * i + 32 * b1 + 32 * b2 + k] = 0.0f;  
            }  
        }  
    }  
}  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum[32 * i + 32 * b1 + 32 * b2 + 32 * k + j] = 0.0f;  
                }  
            }  
        }  
    }  
}  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum[32 * i + 32 * b1 + 32 * b2 + 32 * k + j] += A[32 * b1 + (32 * b2 + k) + 32 * b3 + j] * pB[32768 * b1 + 328 * b2 + 32 * k + j];  
                }  
            }  
        }  
    }  
}  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum[32 * i + 32 * b1 + 32 * b2 + 32 * k + j] += A[32 * b1 + (32 * b2 + k) + 32 * b3 + 1 + j] * pB[32768 * b1 + 328 * b2 + 32 * k + 1 + j];  
                }  
            }  
        }  
    }  
}  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum[32 * i + 32 * b1 + 32 * b2 + 32 * k + j] += A[32 * b1 + (32 * b2 + k) + 32 * b3 + 2 + j] * pB[32768 * b1 + 328 * b2 + 32 * k + 2 + j];  
                }  
            }  
        }  
    }  
}  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum[32 * i + 32 * b1 + 32 * b2 + 32 * k + j] += A[32 * b1 + (32 * b2 + k) + 32 * b3 + 3 + j] * pB[32768 * b1 + 328 * b2 + 32 * k + 3 + j];  
                }  
            }  
        }  
    }  
}  
#pragma omp parallel for  
for (int i = 0; i < 32; i++) {  
    for (int b1 = 0; b1 < 32; b1++) {  
        for (int b2 = 0; b2 < 32; b2++) {  
            for (int k = 0; k < 4; k++) {  
                for (int j = 0; j < 32; j++) {  
                    sum[32 * i + 32 * b1 + 32 * b2 + 32 * k + j] += sum[32 * i + j];  
                }  
            }  
        }  
    }  
}  
free(sum);  
free(pB);
```

What Future for Program Optimization?

- ▶ compilers answer neither current nor future optimization needs
- ▶ algorithms and hardware architectures evolve faster than compilers
- ▶ falling back to manual optimization **slows down progress**

PLDI 2024 Keynote: The Future of Fast Code: Giving Hardware What It Wants

Towards Safe Interactive Optimization Across Layers



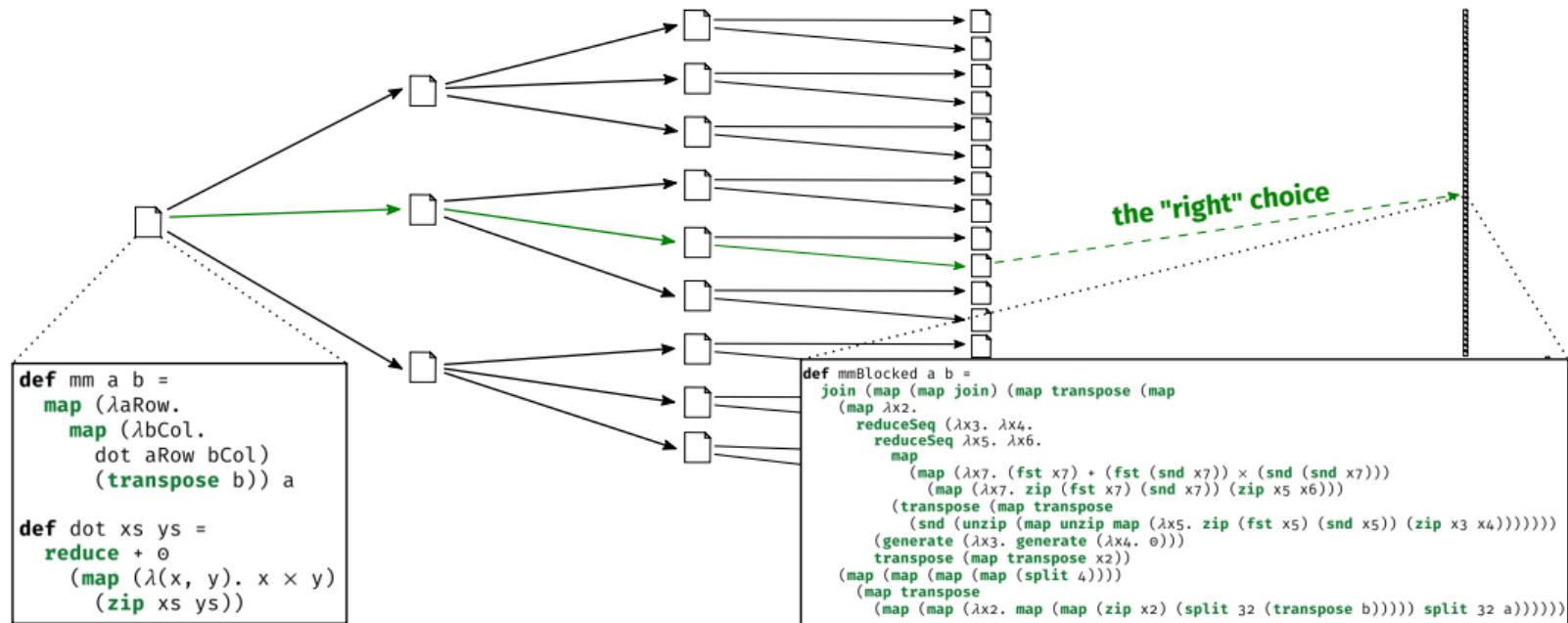
PhD: Functional Program Rewriting

rewrite rules



$$(\text{map } f) . (\text{map } g) = \text{map } (f . g)$$

combinatorial rewriting space, correct and extensible

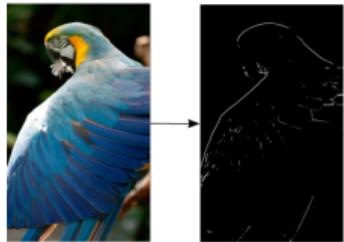


Achieving Expert Optimizations by Composition

6 expert optimizations

→ decomposed into 74 rules

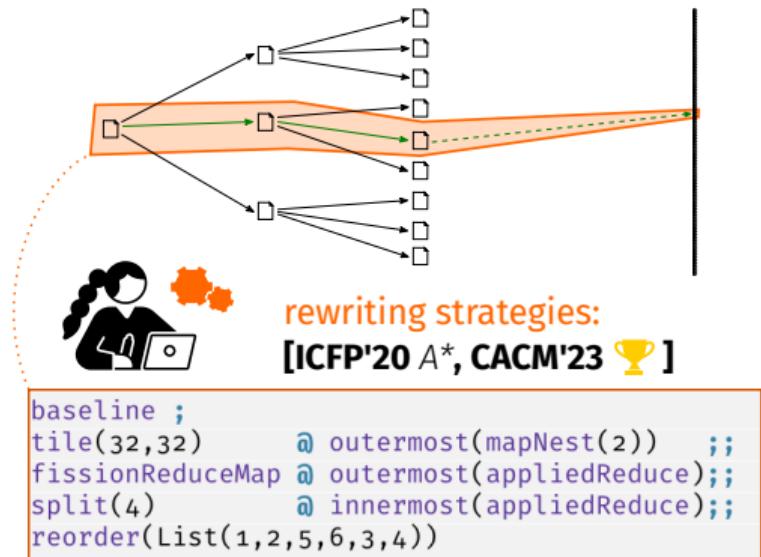
Harris:
corner
and edge
detection



16x faster than OpenCV,
1.4x faster than Halide! [CGO'21 A]

↑
specialized compiler
of industrial strength

enormous rewrite space,
10x bigger programs than before,
thousands of rewrite steps



Combining Automatic Search and Human Expertise

Guided Equality Saturation [POPL'24 A*]

=

automatic rewrite search,
sharing equivalent subterms

+

specifying guides as program sketches

MIT PL Review:



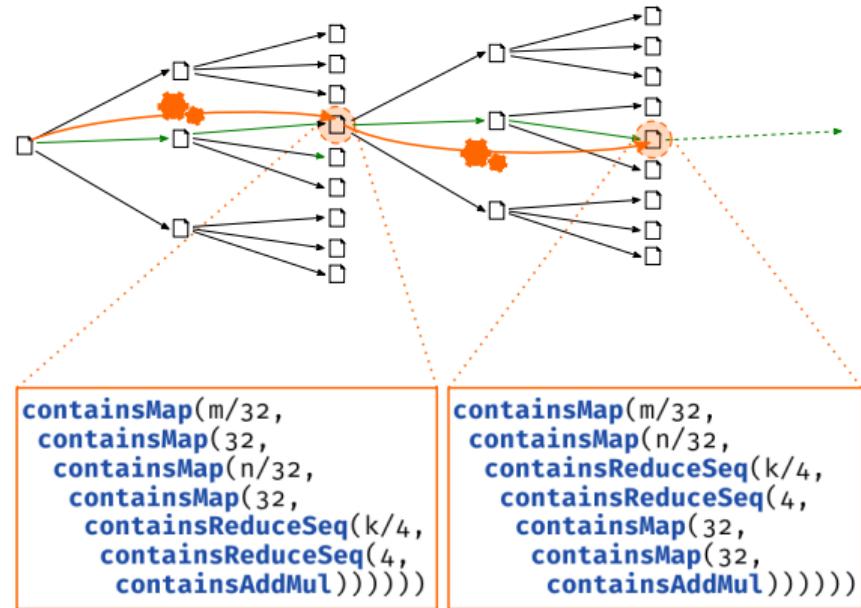
"significant potential to shape the future direction of PL research and/or industry practice"

theorem proving collaboration:

Andrés Goens, University of Amsterdam

guided search :

582x faster, 116x less memory



PostDoc: Transforming Imperative Programs

- ▶ add **imperative choices** to the optimization space *reusing memory*
- ▶ **general-purpose language** more expressive than functional array language $C, C++$

```
for (int i = 0; i < 1024; i++) {                                for (int bi = 0; bi < 32; bi++) {  
    for (int j = 0; j < 1024; j++) {                            for (int i = 0; i < 32; i++) {  
        float sum = 0.f;                                    for (int j = 0; j < 1024; j++) {  
            for (int k = 0; k < 1024; k++) {                float sum = 0.f;  
                sum += A[i][k] * B[k][j];                    for (int k = 0; k < 1024; k++) {  
                    sum += A[bi * 32 + i][k] * B[k][j];  
                }  
                C[bi * 32 + i][j] = sum;  
            }  
        }  
        C[i][j] = sum;  
    }  
}
```

The OptiTrust Project:

- ▶ user-guided source-to-source transformations
- ▶ validated through static resource analysis
based on separation logic
- ▶ internally: imperative λ -calculus

Expert Optimizations over General-Purpose Imperative Code

Realized Case Studies:

with Arthur Chaguéraud and Guillaume Bertholon

► Matrix Multiplication

linear algebra, baseline: TVM specialized compiler

loop and data layout transformations, OpenMP parallelism

► Box Filter (row-based): blurring and denoising

image processing, baseline: OpenCV optimized library

multi-versioning, algorithmic sliding window optimization

In Progress:

► Harris: corner and edge detection

image processing, baseline: specialized compiler Halide

uneven tiling, operator fusion, circular buffers, CSE

► Particle-In-Cell: plasma simulation

physics simulation, baseline: manual expert optimizations

AoS \leftrightarrow SoA, transforming values in memory, atomic bags

Optimizing MatMul with Domain-Specific Compilers

What to compute: TVM Algorithm

```
pB = tvm.compute((N / 32, P, 32),
  lambda bj, k, j:
    B[k, bj * 32 + j])

C = tvm.compute((M, N),
  lambda i, j:
    sum(A[i, k] *
      pB[j // 32, k, j % 32],
      axis=k))
```

How to compute: TVM Schedule

```
CC = s.cache_write(C, "global")
bi, bj, i, j = s[C].tile(
  C.op.axis[0], C.op.axis[1], 32, 32)
s[CC].compute_at(s[C], bj)
i2, j2 = s[CC].op.axis
(kaxis,) = s[CC].op.reduce_axis
bk, k = s[CC].split(kaxis, factor=4)
s[CC].reorder(bk, i2, k, j2)
s[CC].vectorize(j2)
s[CC].unroll(k)
s[C].parallel(bi)
bj3, _, j3 = s[pB].op.axis
s[pB].vectorize(j3)
s[pB].parallel(bj3)
```

Optimizing MatMul with Domain-Specific Compilers

What to compute: TVM Algorithm

```
pB = tvm.compute((N / 32, P, 32),  
    lambda bj, k, j:  
        B[  
            high-level DSL  
            restricted domain  
            rewritten for schedule  
            pB[j // 32, k, j % 32],  
            axis=k))
```

How to compute: TVM Schedule

```
CC = s.cache_write(C, "global")  
bi, bj, i, j = s[C].tile(  
    C.op.axis[0], C.op.axis[1], 32, 32)  
s[CC].compute_at(s[bi].op, bi)  
    concise  
    safe  
    built-in optimizations  
    black box code generation  
    implicit heuristics + unfamiliar IR  
    parallel  
bj3, _, j3 = s[pB].op.axis  
s[pB].vectorize(j3)  
s[pB].parallel(bj3)
```

MatMul Optimization with OptiTrust

What to compute: general-purpose code

```
void mm(float* C, float* A, float* B,
       int m, int n, int p) {
    for (int i = 0; i < m; i++) {
        for (int j = 0; j < n; j++) {
            float sum = 0.0f;
            for (int k = 0; k < p; k++)
                sum += A[i][k] * B[k][j];
            C[i][j] = sum;
        }
    }
}
```

How to optimize: transformation script

```
Function.inline_def [cFunDef "mm"];
let tile (id, tile_size) = Loop.tile (int tile_size)
  -index:"b" ^ id -bound:TileDivides [cFor id] in
List.iter tile [(("i", 32); ("j", 32); ("k", 4));
Loop.reorder_at
  -order:["bi"; "bj"; "bk"; "i"; "k"; "j"] [cPlusEq ()];
Loop.hoist_expr -dest:[tBefore; cFor "bi"] "pB"
  -indep:["bi"; "i"] [cArrayRead "B"];
Matrix.stack_copy -var:"sum" -copy_var:"s"
  -copy_dims:1 [cFor -body:[cPlusEq ()] "k"];
Omp.simd [cFor -body:[cPlusEq ()] "j"];
Omp.parallel_for [cFunBody "mm1024"; cStrict; cFor ""];
Loop.unroll [cFor -body:[cPlusEq ()] "k"];
```

MatMul Optimization with OptiTrust

What to compute: general-purpose code

```
void mm(float* C, float* A, float* B,
       int m, int n, int p) {
    __reads("A ~ Matrix2(m, p)");
    __reads("B ~ Matrix2(p, n)");
    __modifies("C ~ Matrix2(m, n)");
    for (int i = 0; i < m; i++) {
        __xmodifies("for j in 0..n ~ &C[i][j] ~ Cell")
        ;
        for (int j = 0; j < n; j++) {
            __xmodifies("&C[i][j] ~ Cell");
            float sum = 0.0f;
            for (int k = 0; k < p; k++)
                sum += A[i][k] * B[k][j];
            C[i][j] = sum;
        }
    }
}
```

How to optimize: transformation script

```
Function.inline_def [cFunDef "mm"];
let tile (id, tile_size) = Loop.tile (int tile_size)
  -index:"b" ^ id -bound:TileDivides [cFor id] in
List.iter tile [( "i", 32); ( "j", 32); ( "k", 4)];
Loop.reorder_at
  -order:[ "bi"; "bj"; "bk"; "i"; "k"; "j" ] [cPlusEq ()];
Loop.hoist_expr -dest:[tBefore; cFor "bi"] "pB"
  -indep:[ "bi"; "i" ] [cArrayRead "B"];
Matrix.stack_copy -var:"sum" -copy_var:"s"
  -copy_dims:1 [cFor -body:[cPlusEq ()] "k"];
Omp.simd [cFor -body:[cPlusEq ()] "j"];
Omp.parallel_for [cFunBody "mm1024"; cStrict; cFor ""];
Loop.unroll [cFor -body:[cPlusEq ()] "k"];
```

MatMul Resource Contract for OptiTrust

```
void mm(float* C, float* A, float* B, int m, int n, int p) {
    __reads("A ~ Matrix2(m, p), B ~ Matrix2(p, n)");
    __modifies("C ~ Matrix2(m, n)");
    // [...]
}
```

- ▶ clauses describe ownership of resources (permissions)
- ▶ a matrix is a conjunction of cells: $\star_{i \in 0..m} \star_{j \in 0..n} (\&C[i][j] \rightsquigarrow \text{Cell})$

Trace for matmul_check ✓

- Preprocessing loop contracts ✓
- Function.inline_def [cFunDef "mm"]; ✓
- List.iter tile [("i", 32); ("i", 32); ("k", 4)]; ✓
- Loop.reorder_at ~order:[("bi"; "bj"; "bk"; "i"; "k"; "j")]; ✓
- [cPlusEq ~lhs:[cVar "sum"] ()]; ✓

- Loop.reorder_at ✓
 - bring down j ✓
 - Loop.hoist_alloc_loop_list ✓
 - Loop.fission ✓
 - Loop.fission ✓
 - Loop_swap.swap ✓
 - bring down j ✓
 - Loop.hoist_alloc_loop_list ✓
 - Loop.fission ✓
 - Loop.fission ✓
 - Loop_swap.swap ✓
 - bring down i ✓
 - Loop.hoist_alloc_loop_list ✓
 - Loop.fission ✓
 - Loop.fission ✓
 - Loop_swap.swap ✓
 - bring down i ✓
 - Loop.hoist_alloc_loop_list ✓
 - Loop.fission ✓
 - Loop.fission ✓
 - Loop_swap.swap ✓

Loop.hoist_expr ~dest:[tBefore; cFor "bi" "pB"
~indep:["bi"; "i"] [cArrayRead "B"]; ✓

Matrix.stack_copy ~var:"sum" ~copy_var:"s"
~copy_dims:1 [cFor ~body:[cPlusEq ~lhs:[cVar "sum"] ()] "K"]; ✓

Omp.simd [nbMulti; cFor ~body:[cPlusEq ~lhs:[cVar "s"] ()] "I"]; ✓

Omp.parallel_for [nbMulti; cFunBody "", cStrict;
cFor ""]; ✓

```

@@ -1,20 +1,38 @@
1 #include <optitrust.h>
2
3 void mm1024(float* C, float* A, float* B) {
4     for (int bi = 0; bi < 32; bi++) {
5         for (int i = 0; i < 32; i++) {
6             for (int bj = 0; bj < 32; bj++) {
7                 for (int j = 0; j < 32; j++) {
8                     float sum = 0.f;
9
10                for (int bk = 0; bk < 256; bk++) {
11                    sum += A[MINDEX2(1024, 1024, bi * 32 + i, bk * 4 + k)] *
12                        B[MINDEX2(1024, 1024, bk * 4 + k, bj * 32 + j)];
13                }
14            }
15            C[MINDEX2(1024, 1024, bi * 32 + i, bj * 32 + j)] = sum;
16        }
17    }
18 }
19 }
20 }

#include <optitrust.h>
2
3 void mm1024(float* C, float* A, float* B) {
4     for (int bi = 0; bi < 32; bi++) {
5         for (int i = 0; i < 32; i++) {
6             for (int bj = 0; bj < 32; bj++) {
7                 float const sum = (float*)const MALLOC2(32, 32, sizeof(float));
8                 for (int i = 0; i < 32; i++) {
9                     for (int j = 0; j < 32; j++) {
10                        sum[MINDEX2(32, 32, i, j)] = 0.f;
11                    }
12                }
13            }
14        for (int bk = 0; bk < 256; bk++) {
15            for (int i = 0; i < 32; i++) {
16                for (int j = 0; j < 32; j++) {
17                    for (int k = 0; k < 4; k++) {
18                        for (int j = 0; j < 32; j++) {
19                            sum[MINDEX2(32, 32, i, j)] +=
20                                A[MINDEX2(1024, 1024, bi * 32 + i, bk * 4 + k)] *
21                                B[MINDEX2(1024, 1024, bk * 4 + k, bj * 32 + j)];
22                        }
23                    }
24                }
25            }
26        }
27        for (int i = 0; i < 32; i++) {
28            for (int j = 0; j < 32; j++) {
29                C[MINDEX2(1024, 1024, bi * 32 + i, bj * 32 + j)] =
30                    sum[MINDEX2(32, 32, i, j)];
31            }
32        }
33        MFREE2(32, 32, sum);
34    }
35    for (int i = 0; i < 32; i++) {
36    }
37 }
38 }
```

advanced arguments justification
 normal full

Diff Code before Code after Decode Hide res Res annot Res context Res usage Full res Compact

Decomposition of MatMul Steps by OptiTrust

Combined Transformations:

- ▶ compose basic transformations
- ▶ validated by transitivity

Basic Transformations:

- ▶ directly modify AST
imperative λ -calculus
- ▶ check sufficient correctness conditions

8 MatMul Steps:

- ▶ 55 basic transformations
- ▶ 61 “ghost” transformations

- ▶ **Instr**: move, delete, insert, duplicate
- ▶ **Function**: inline
- ▶ **Variable**: inline, bind, to_const, init_detach, local_name
- ▶ **Loop**: fusion, fission, swap, hoist, move_out, tile, collapse, unroll
- ▶ **Matrix**: intro_malloc0,
- ▶ **Omp**: parallel_for, simd
- ▶ **Arith**: simpl
- ▶ ...

OptiTrust's Sufficient Condition for `Instr.move`

$$E \begin{bmatrix} T_1; \Delta_1 \\ T_2; \Delta_2 \end{bmatrix} \mapsto E \begin{bmatrix} T_2; \\ T_1; \end{bmatrix} \quad \text{correct if:}$$
$$\begin{cases} \Delta_1.\text{notRO} \cap \Delta_2 = \emptyset \\ \Delta_2.\text{notRO} \cap \Delta_1 = \emptyset \end{cases}$$

- ▶ leverages computed resource information (Δ_1)
- ▶ correct if resources are exclusively shared in read-only

OptiTrust's Sufficient Condition for `Loop.fusion`

```
for  $\chi_1$   $i \in r_i \{$ 
     $T_1;$ 
}  $\Delta_1$ 
for  $\chi_2$   $i \in r_i \{$ 
     $T_2;$ 
}  $\Delta_2$ 
```

 \mapsto

```
for  $\chi$   $i \in r_i \{$ 
     $T_1;$ 
     $T_2;$ 
}
```

correct if:

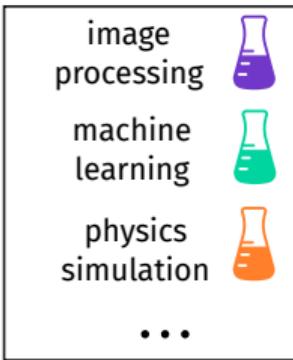
$$\begin{cases} i \text{ not free in } \chi_1.\text{shrd} \text{ and } \chi_2.\text{shrd} \\ \chi_1.\text{shrd} \setminus (\Delta_1.\text{notRO} \cap \Delta_2) = \emptyset \\ \chi_2.\text{shrd} \setminus (\Delta_2.\text{notRO} \cap \Delta_1) = \emptyset \end{cases}$$

with:

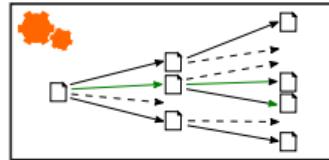
$$___, Q_1, P_2 \equiv \text{PartialSub}(\chi_1.\text{excl}.post, \chi_2.\text{excl}.pre)$$

$$\chi \equiv \begin{cases} \text{vars} \equiv \chi_1.\text{vars}, \chi_2.\text{vars} \\ \text{shrd} \equiv \chi_1.\text{shrd} * \chi_2.\text{shrd} \\ \text{excl}.pre \equiv \chi_1.\text{excl}.pre * P_2 \\ \text{excl}.post \equiv \chi_2.\text{excl}.post * Q_1 \end{cases}$$

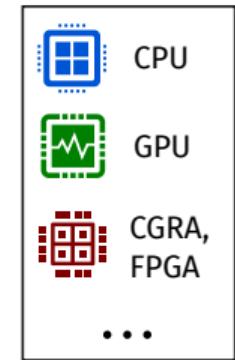
Future: Safe Interactive Optimization Across Layers



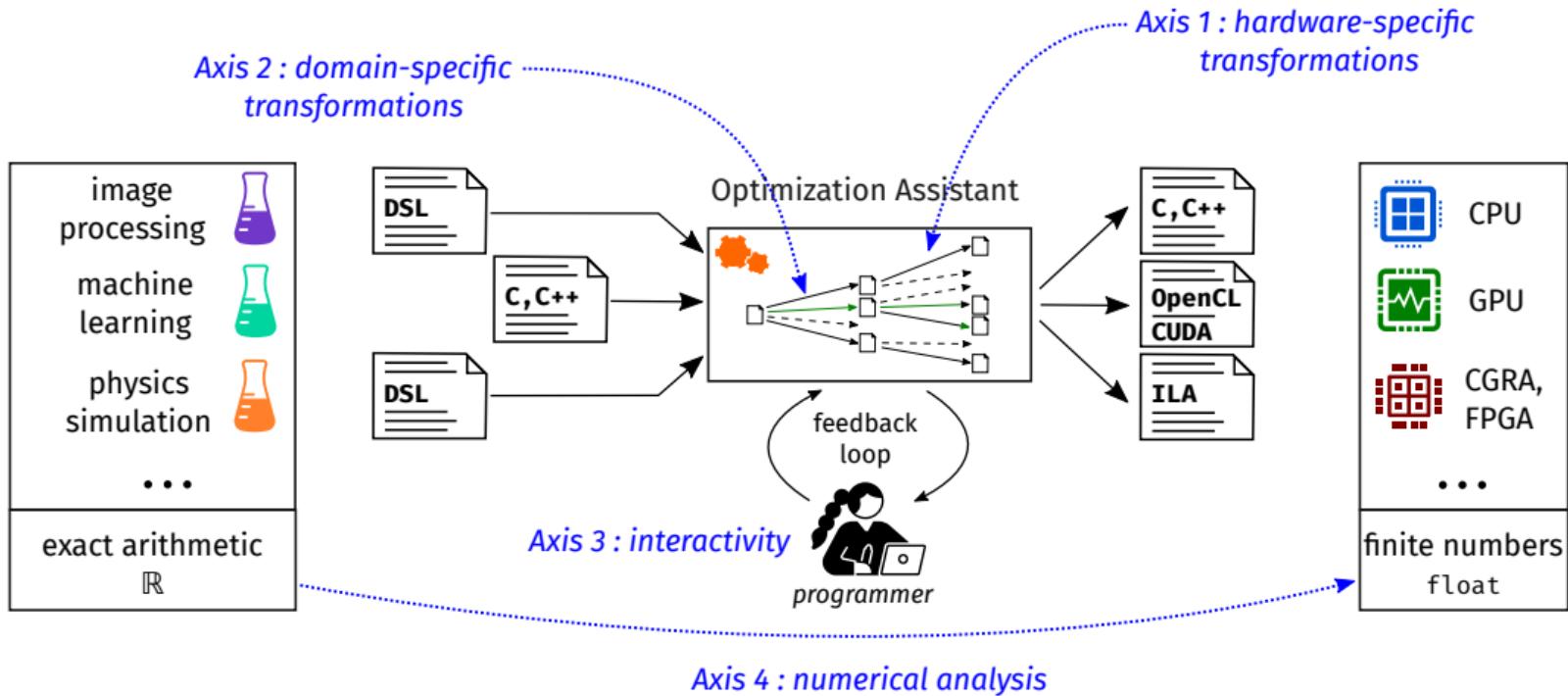
Goal:
build the first
Optimization Assistant



that adapts to domain and hardware evolution,
and avoids falling back to manual optimization



Future: Safe Interactive Optimization Across Layers



Future: Safe Interactive Optimization Across Layers

- ▶ Looking for Master students (interns)
- ▶ Preparing for PhD students (ANR JCJC)
- ▶ Starting fruitful collaborations
 - ▶ Combine different expertise
 - ▶ Leverage overlapping interests
 - ▶ Exchange benchmarks and applications
- ▶ Creating a French, European, and International network (workshops, ERC, ..)

Or, .. maximize fun and impact by not working alone

Axis 1: Correct and Composable Hardware-Specific Transformations

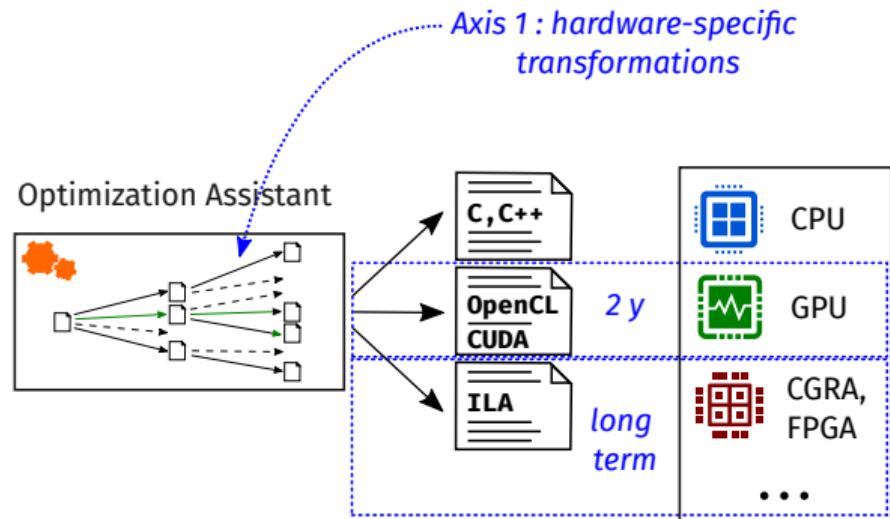
Goal: create an optimization space
that evolves with hardware
by composition of transformations

Approach:

- decompose expert optimizations
into simpler transformations

Difficulty:

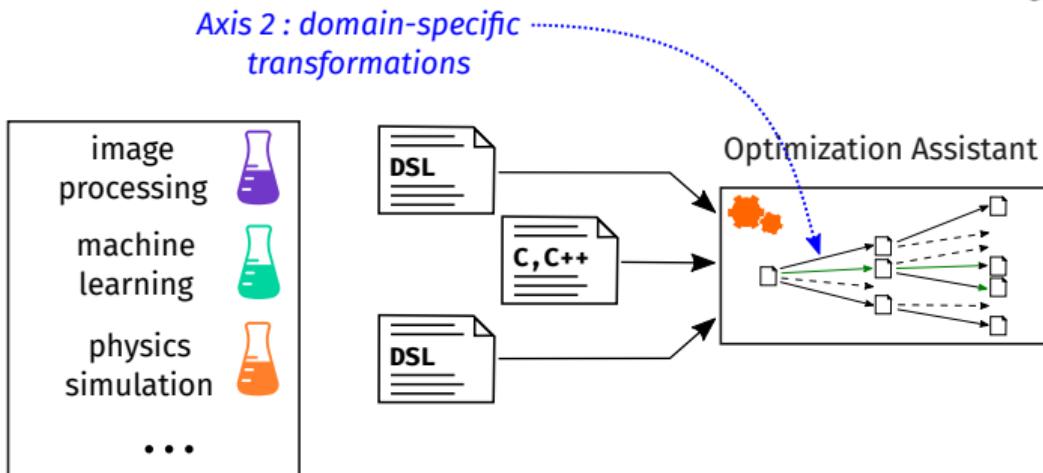
- identify and guarantee the correctness
of transformations, taking programming
models into account



International Collaborations:

- Bastian Köpcke, Universität Münster [[arXiv'22](#), GPGPU'22, PLDI'24]
- Jackson Woodruff, University of Edinburgh [PLDI'22, [arXiv'23](#), CC'23]

Axis 2: Libraries of Composable Domain-Specific Transformations



Goal: create an optimization space that evolves with domains and their specialized languages (DSLs)

Plan:

- 2 y: Halide-like DSL
- 4 y: combine with Axis 1 (GPUs)
- long term: mixing DSLs, libraries = functions
 - + transformations
 - + heuristics

Overlapping Interests?

- Gabriel Radanne
- AnyDSL team (Roland Leissa, Sebastian Hack)
- Inria EMERAUDE

Difficulties:

- synthesizing code annotations
- facilitate defining custom transformations

Axis 3: Interactivity between Programmer and Assistant

Goal: develop an interactive feedback loop allowing to productively explore the optimization space

Plan:

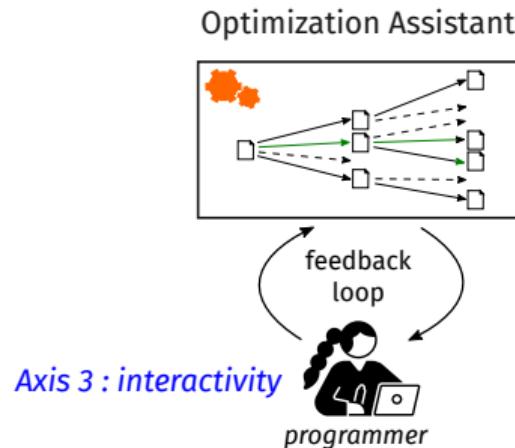
- 3 years: interactive equality saturation on imperative code
polyhedralic and/or beam search
- long term: optimization suggestions and visualizations
bottlenecks, hot code

Difficulty:

- exploring a complex and evolving optimization space

Interdisciplinary Collaboration:

- Géry Casiez (Inria LOKI, Human-Computer Interaction)



Axis 4: Guaranteeing Numerical Accuracy during Optimization

Goal: guarantee numerical accuracy during optimization,
avoid errors and allow optimizations

*error-tolerant trigonometric function = 4x faster
fixed point numbers for FPGA*

Approach:

- interval analysis on values and errors
error is below $e(x \circ y) + d$

Difficulty:

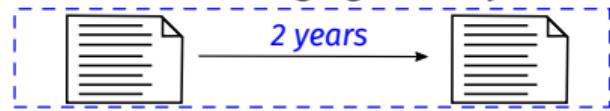
- most compilers do not perform numerical analyses
*'gcc': no optimisations over floats
'gcc -ffast-math': treats floats as reals*

International Collaboration:

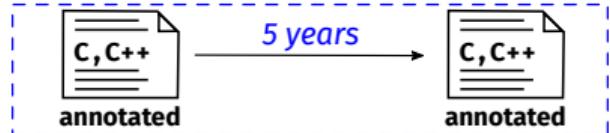
- Eva Darulova, Uppsala University [TOPLAS'17, SAS'23]



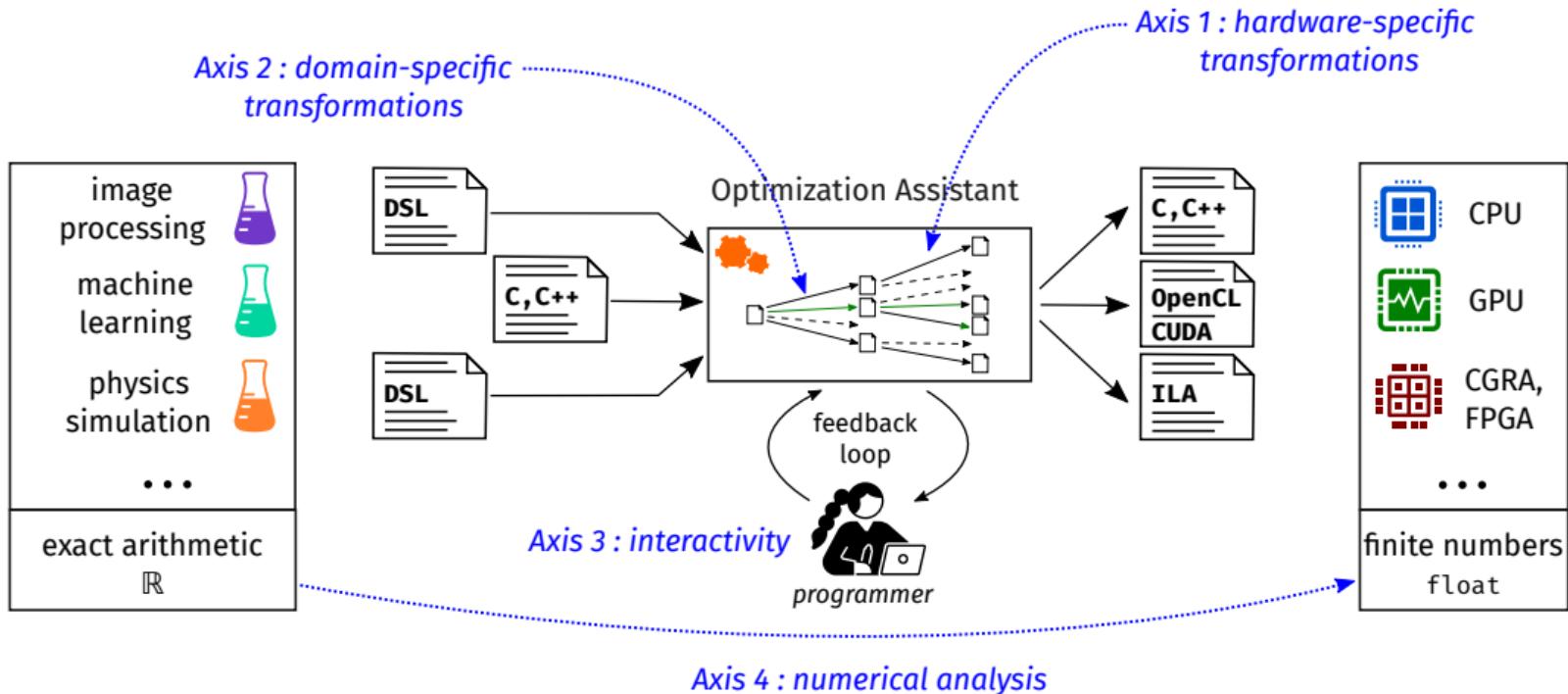
functional language on arrays



general-purpose imperative language



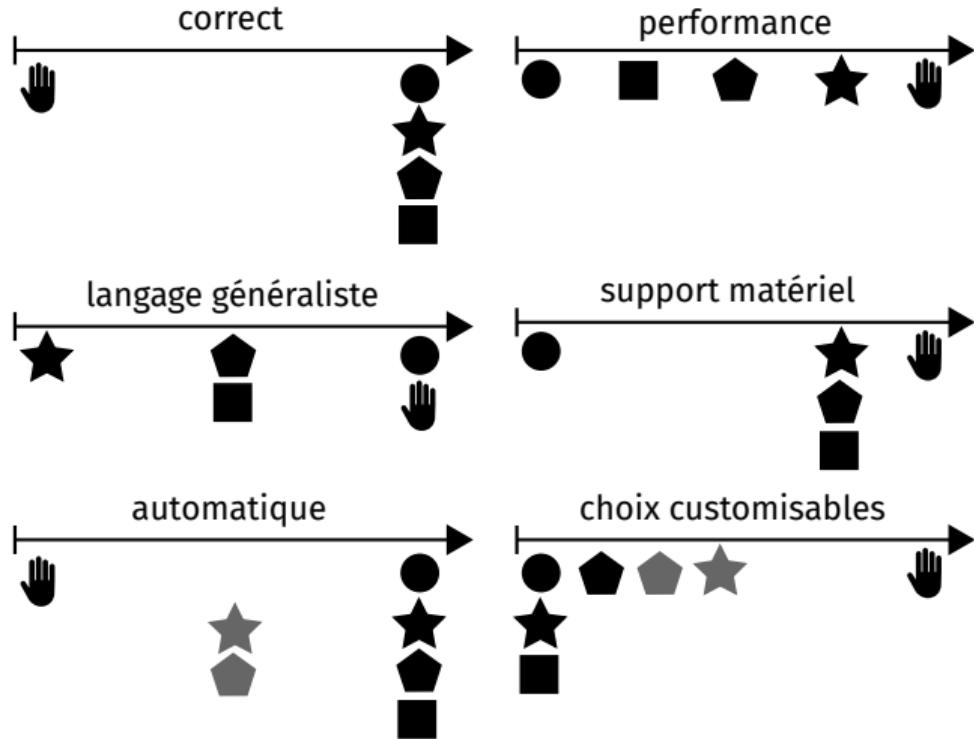
Towards Safe Interactive Optimization Across Layers



Backup Slides

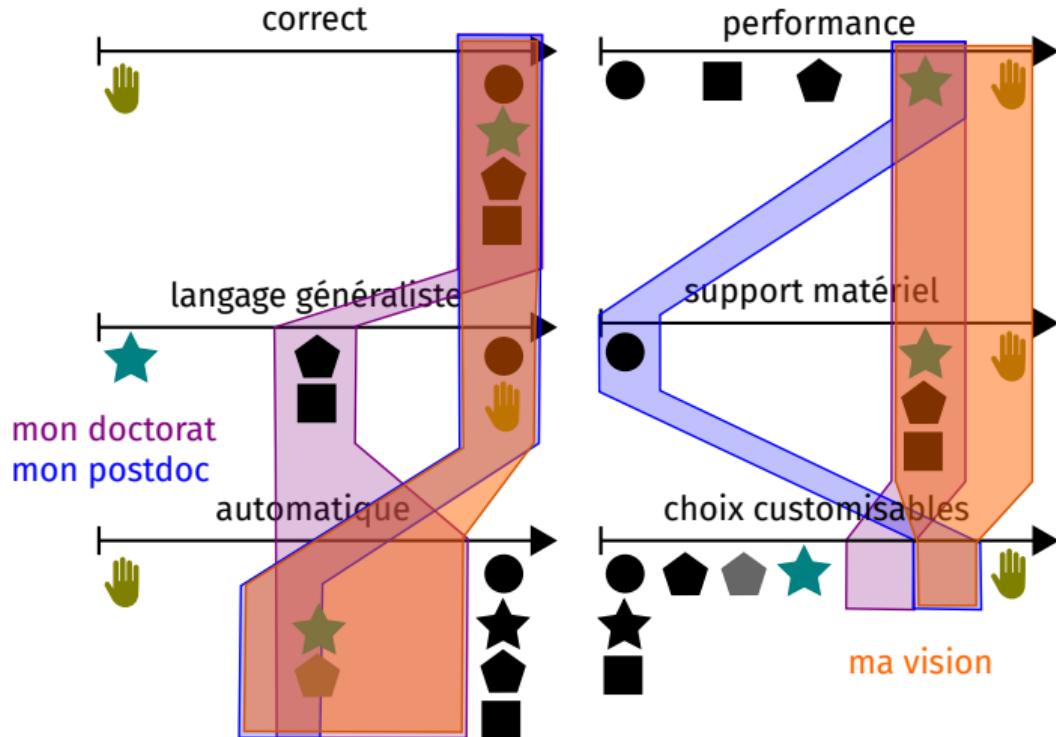
Les Compromis

- compilateurs généralistes
gcc, clang
- ✋ optimisation manuelle
- ★ compilateurs spécialisés
Halide, TVM
- ◤ compilateurs polyédriques
PLUTO, Polly
- compilateurs tableau
Accelerate, SaC, Futhark



Les Compromis

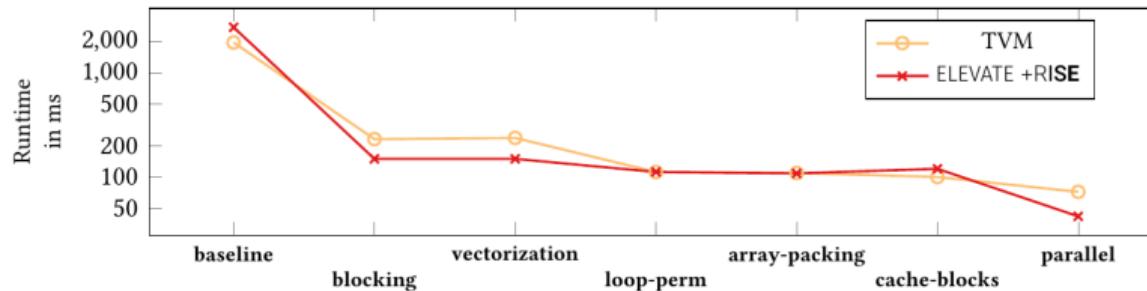
- compilateurs généralistes
gcc, clang
- ✋ optimisation manuelle
- ★ compilateurs spécialisés
Halide, TVM
- ◤ compilateurs polyédriques
PLUTO, Polly
- compilateurs tableau
Accelerate, SaC, Futhark



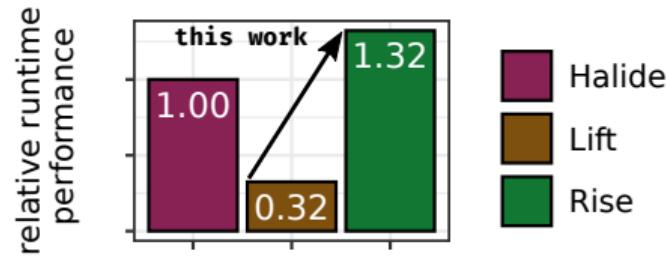
1st Benchmark: Matrix Multiplication on Intel CPU

- ▶ 6 optimizations
 - ▶ transform loops *blocking, permutation, unrolling*
 - ▶ change data layout *array packing*
 - ▶ add parallelism *vectorization, multi-threading*
- ▶ performance is on par with reference schedules from TVM.

https://tvm.apache.org/docs/how_to/optimize_operators/opt_gemm.html



2nd Benchmark: Corner Detection on ARM CPUs



- ▶ standard corner detection pipeline
- ▶ 6 well-known optimizations
 - circular buffering, operator fusion, multi-threading, vectorization, convolution separation, register rotation*
- ▶ extensibility + control
====>
faster code than Halide, with 2 additional optimizations

Matrix Multiplication Performance

- ▶ Intel(R) Core(TM) i7-8665U CPU, AVX2 (8 floats), 4 cores (8 hyperthreads)
- ▶ Relative speedup on 1024^3 input:

version	single-thread	multi-thread
unoptimized	1×	1×
optimized	46×	150×
TVM	46×	150×
numpy (Intel MKL) ¹	71×	183×

Both codes have 90th percentile runtime of 9.4ms over 200 benchmark runs, corresponding to a speedup of 150× compared to the 90th percentile of the naive code.

¹uses assembly code, explicit vectorization, custom thread library