

Interactive Source-to-Source Optimizations Validated using Static Resource Analysis (OptiTrust)

Guillaume BERTHOLON

Arthur CHARGUÉRAUD

Thomas KÉHLER

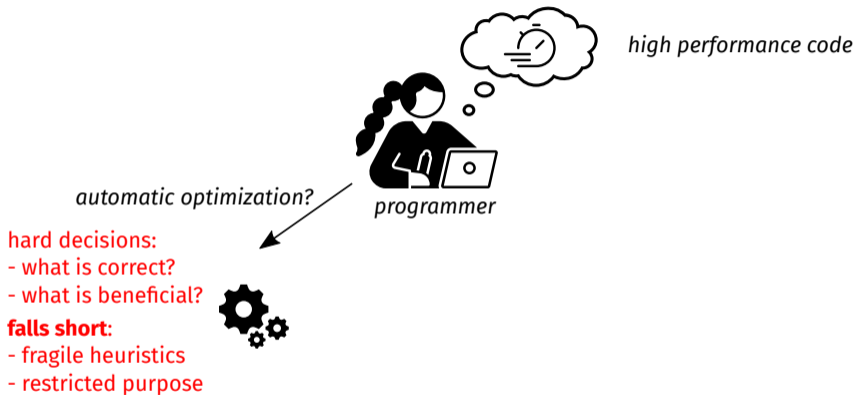
Begatim BYTYQI

Damien ROUHLING

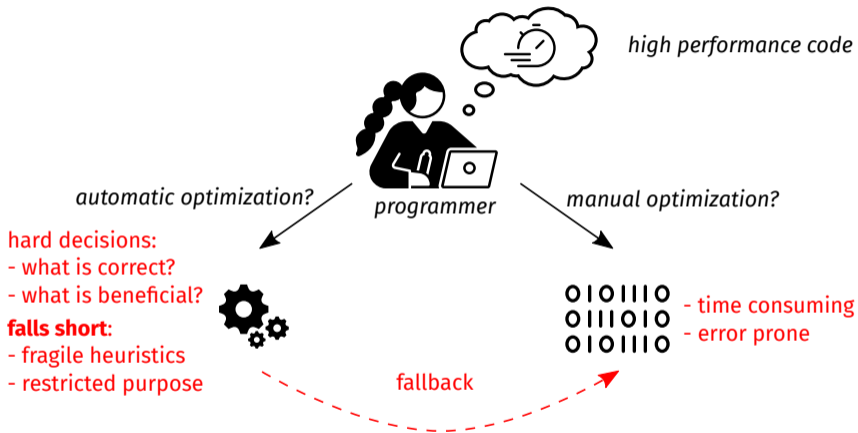
The logo for Inria, featuring the word "Inria" in a stylized, red, cursive script font.

SOAP Workshop – Copenhagen, June 2024

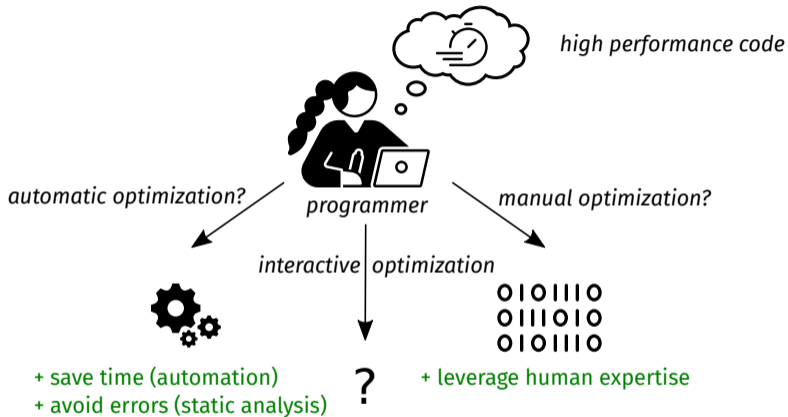
Producing High-Performance Code is Hard



Producing High-Performance Code is Hard



Producing High-Performance Code is Hard



The OptiTrust Interactive Optimization Framework

Users direct C code optimizations, OptiTrust checks the correctness of the optimizations

Case studies:

- ▶ **Matrix Multiplication reproducing TVM baseline** *linear algebra*
loop and data layout transformations, OpenMP parallelism
- ▶ **Box Filter (on rows) reproducing handwritten OpenCV code** *image processing*
multi-versioning, algorithmic sliding window optimization

Work in progress:

- ▶ **Harris Corner Detection reproducing Halide baseline** *image processing*
uneven tiling, operator fusion, circular buffers, CSE
- ▶ **Particle-In-Cell reproducing handwritten code** *physics simulation*
AoS \leftrightarrow SoA, transforming values in memory, atomic bags

Hand Optimizing Matrix Multiplication

Optimized C code:

Naive 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;
  }
}
```



```
float* pB = (float*) malloc(sizeof(float)[32][256][4][32]);
#pragma omp parallel for
for (int bj = 0; bj < 32; bj++)
  for (int bk = 0; bk < 256; bk++)
    for (int k = 0; k < 4; k++)
      for (int j = 0; j < 32; j++)
        pB[32768 * bj + 128 * bk + 32 * k + j] = B[1024 * (4 * bk + k) + 32 * bj + j];
#pragma omp parallel for
for (int bi = 0; bi < 32; bi++) {
  for (int bj = 0; bj < 32; bj++) {
    float* sum = (float*) malloc(sizeof(float)[32][32]);
    for (int i = 0; i < 32; i++)
      for (int j = 0; j < 32; j++)
        sum[32 * i + j] = 0.0f;
    for (int bk = 0; bk < 256; bk++)
      for (int i = 0; i < 32; i++)
        float s[32];
        memcpy(s, &sum[32 * i + bj], sizeof(s));
    #pragma omp simd
    for (int j = 0; j < 32; j++)
      s[j] += A[1024 * bi + 32 * bj + j];
    #pragma omp simd
    for (int j = 0; j < 32; j++)
      s[j] += A[1024 * bi + 32 * bj + j];
    #pragma omp simd
    for (int j = 0; j < 32; j++)
      s[j] += A[1024 * bi + 32 * bj + j];
    #pragma omp simd
    for (int j = 0; j < 32; j++)
      s[j] += A[1024 * bi + 32 * bj + j];
    memcpy(6sum[32 * i], s, sizeof(float)[32]);
  }
}
for (int i = 0; i < 32; i++)
  for (int j = 0; j < 32; j++)
    C[1024 * (32 * bi + i) + 32 * bj + j] = sum[32 * i + j];
free(sum);
}
free(pB);
```

150× faster

7× more characters
time consuming
error prone

Hand Optimizing Matrix Multiplication

Optimized C code:

Naive 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;
  }
}
```



```
float* pB = (float*) malloc(sizeof(float)[32][256][4][32]);
#pragma omp parallel for
for (int bj = 0; bj < 32; bj++)
  for (int bk = 0; bk < 256; bk++)
    for (int k = 0; k < 4; k++)
      for (int j = 0; j < 32; j++)
        pB[32768 * bj + 128 * bk + 32 * k + j] = B[1024 * (4 * bk + k) + 32 * bj + j];
#pragma omp parallel for
for (int bi = 0; bi < 32; bi++)
  for (int bj = 0; bj < 32; bj++)
    float* sum = (float*) malloc(sizeof(float)[32]);
    for (int i = 0; i < 1024; i++)
      for (int bk = 0; bk < 32; bk++)
        float s[32];
        memcpy(s, pB[32768 * bj + 128 * bk + 32 * bi + i], sizeof(s));
#pragma omp simd
for (int j = 0; j < 32; j++)
  s[j] += A[1024 * (32 * bi + i) + 4 * bk + 3] * pB[32768 * bj + 128 * bk + 32 * bi + i];
#pragma omp simd
for (int j = 0; j < 32; j++)
  s[j] += A[1024 * (32 * bi + i) + 4 * bk + 3] * pB[32768 * bj + 128 * bk + 32 * bi + i];
#pragma omp simd
for (int j = 0; j < 32; j++)
  s[j] += A[1024 * (32 * bi + i) + 4 * bk + 3] * pB[32768 * bj + 128 * bk + 32 * bi + i];
#pragma omp simd
for (int j = 0; j < 32; j++)
  s[j] += A[1024 * (32 * bi + i) + 4 * bk + 3] * pB[32768 * bj + 128 * bk + 32 * bi + i];
  memcpy(6sum[32 * i], s, sizeof(float)[32]);
}
}
for (int i = 0; i < 32; i++)
  for (int j = 0; j < 32; j++)
    C[1024 * (32 * bi + i) + 32 * bj + j] = sum[32 * i + j];
free(sum);
}
free(pB);
```

specialized input sizes
 $m = n = p = 1024$
improved data locality
transformed loops and data layout
parallel
vectorization, multi-threading

Optimizing MatMul with Domain-Specific Compilers

How to compute: TVM Schedule

What to compute: TVM Algorithm

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

```
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 bi, k, i:  
    B  
C =  
l  
s  
    pB[j // 32, k, j % 32],  
    axis=k))
```

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]  
i2,  
(kax  
bk, l  
s[CC]  
s[CC]  
s[CC]  
s[C].parallel(bj1)  
bj3, _, j3 = s[pB].op.axis  
s[pB].vectorize(j3)  
s[pB].parallel(bj3)
```

concise

safe

built-in optimizations

black box code generation

implicit heuristics + unfamiliar IR

Optimizing MatMul 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"];
```

Optimizing MatMul 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.choist_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"];
```

Optimizing MatMul 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-1");
    for (int j = 0; j < n; j++)
      __xmodifies("&C[i][j] ~> float");
      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
  ("j", 32); ("k", 4)];
  ["i"; "k"; "j"] [cPlusEq ()];
  cFor "bi" "pB"
  yRead "B"];
  n" -copy_var:"s"
  [cPlusEq ()] "k"];
  sEq ()] "j"];
  "mm1024"; cStrict; cFor "";
  PlusEq ()] "k"];
}
```

source-to-source transformations

validated using static resource analysis

based on separation logic

no black box code generation

readable intermediate steps

easier to compose and extend

Resource Analysis in MatMul

user-provided function contracts

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

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

Resource Analysis in MatMul

user-provided partial loop contracts

```
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;
  }
}
```

Resource Analysis in MatMul

automatically inferred annotations

```
for (int i = 0; i < m; i++) {
  __xmodifies("for j in 0..n -> &C[i][j] ~ Cell");
  __sreads("A ~ Matrix2(m, p), B ~ Matrix2(p, n)");
  for (int j = 0; j < n; j++) {
    __xmodifies("&C[i][j] ~ Cell");
    __sreads("A ~ Matrix2(m, p), B ~ Matrix2(p, n)");
    float sum = 0.0f;
    for (int k = 0; k < p; k++) {
      __smodifies("&sum ~ Cell");
      __sreads("A ~ Matrix2(m, p), B ~ Matrix2(p, n)");
      __ghost(matrix2_ro_focus, "A, i, k");
      __ghost(matrix2_ro_focus, "B, k, j");
      sum += A[i][k] * B[k][j];
      __ghost(matrix2_ro_unfocus, "A");
      __ghost(matrix2_ro_unfocus, "B");
    }
    C[i][j] = sum;
  }
}
```

Resource Analysis in MatMul

automatically computed resources at every program point

```
 $R_1 : \alpha_1(A \rightsquigarrow \text{Matrix2}(m, p)) * R_2 : \alpha_2(B \rightsquigarrow \text{Matrix2}(p, n)) *$   
 $R_3 : \star i \in 0..m \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$   
for (int i = 0; i < m; i++) {  
   $R_4 : (\alpha_4/m)(A \rightsquigarrow \text{Matrix2}(m, p)) * R_5 : (\alpha_5/m)(B \rightsquigarrow \text{Matrix2}(p, n)) *$   
   $R_6 : \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$   
  for (int j = 0; j < n; j++) {  
     $R_7 : \dots * R_8 : \dots *$   
     $R_9 : \&C[i][j] \rightsquigarrow \text{Cell}$   
    float sum = 0.f;  
     $R_7 * R_8 * R_9 * R_{10} : \&\text{sum} \rightsquigarrow \text{Cell}$   
    for (int k = 0; k < p; k++) { /* [...] */ }  
     $R_7 * R_8 * R_9 * R_{11} : \&\text{sum} \rightsquigarrow \text{Cell}$   
    C[i][j] = sum;  
     $R_7 * R_8 * R_{11} * R_{12} : \&C[i][j] \rightsquigarrow \text{Cell}$   
  }  
   $R_4 * R_5 * R_{13} : \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$   
}  
 $R_1 * R_2 * R_{14} : \star i \in 0..m \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$ 
```


Resource Analysis in MatMul

automatically computed resources at every program point

```
 $R_1 : \alpha_1 (A \rightsquigarrow \text{Matrix2}(m, p)) * R_2 : \alpha_2 (B \rightsquigarrow \text{Matrix2}(p, n)) *$   
 $R_3 : \star i \in 0..m \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$   
for (int i = 0; i < m; i++) {  
   $R_4 : (\alpha_4/m) (A \rightsquigarrow \text{Matrix2}(m, p)) * R_5 : (\alpha_5/m) (B \rightsquigarrow \text{Matrix2}(p, n)) *$   
   $R_6 : \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$   
  for (int j = 0; j < n; j++) {  
     $R_7 : \dots * R_8 : \dots *$   
     $R_9 : \&C[i][j] \rightsquigarrow \text{Cell}$   
    float sum = 0.f;  
     $R_7 * R_8 * R_9 * R_{10} : \&\text{sum} \rightsquigarrow \text{Cell}$   
    for (int k = 0; k < p; k++) { /* [...] */ }  
     $R_7 * R_8 * R_9 * R_{11} : \&\text{sum} \rightsquigarrow \text{Cell}$   
    C[i][j] = sum;  
     $R_7 * R_8 * R_{11} * R_{12} : \&C[i][j] \rightsquigarrow \text{Cell}$   
  }  
   $R_4 * R_5 * R_{13} : \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$   
}  
 $R_1 * R_2 * R_{14} : \star i \in 0..m \star j \in 0..n \ \&C[i][j] \rightsquigarrow \text{Cell}$ 
```

Combined Transformations on MatMul

Combined transformations:

- ▶ compose basic transformations
- ▶ are valid by composition

MatMul: 8 script steps result in 55 basic transformations (+61 ghost transformations).

Zoom on the 3rd transformation step:

```
!! Loop.reorder_at ~order:["bi"; "bj"; "bk"; "i"; "k"; "j"] [cPlusEq ()];
```

- ▶ 18 basic transformations
 - ▶ 4 `Loop.swap`
 - ▶ 6 `Loop.fission`
 - ▶ 2 `Loop.hoist`
 - ▶ .. 6 more
- ▶ 32 ghost transformations

- Trace for matmul_check ✓
- Preprocessing loop contracts ✓
- Function.inline_def [cFunDef "mm"]: ✓
- List.iter tile [("T", 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"] ()] "j"]: ✓
- Omp.parallel_for [nbMulti; cFunBody ""; cStrict; cFor ""]; ✓

```

1  @@ -1,20 +1,38 @@
2  #include <optitrust.h>
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  for (int bk = 0; bk < 256; bk++) {
10 for (int k = 0; k < 4; k++) {
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 }

```

```

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  }
7  for (int bj = 0; bj < 32; bj++) {
8  float* const sum = (float* const)MALLOC2(32, 32, sizeof(float));
9  for (int i = 0; i < 32; i++) {
10 for (int j = 0; j < 32; j++) {
11 sum[MINDEX2(32, 32, i, j)] = 0.f;
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 }
18 for (int k = 0; k < 4; k++) {
19 for (int j = 0; j < 32; j++) {
20 sum[MINDEX2(32, 32, i, j)] +=
21 A[MINDEX2(1024, 1024, bi * 32 + i, bk * 4 + k)] *
22 B[MINDEX2(1024, 1024, bk * 4 + k, bj * 32 + j)];
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 }

```

Basic Transformations

Basic transformations:

- ▶ transform the AST directly
 - internally, we transform an imperative lambda calculus AST with bijective encoding from C*
- ▶ check for sufficient correctness conditions

- ▶ **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_malloco,
- ▶ **Omp**: parallel_for, simd
- ▶ **Arith**: simpl
- ▶ ...

Correctness Criteria for `Instr.move`

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

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

Correctness Criteria for Loop.fusion

```

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



```

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

correct if:

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

with:

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

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



Conclusion

- ▶ OptiTrust is an interactive framework for optimizing general-purpose C code
- ▶ Transformations are validated using a static resource analysis
- ▶ 2 completed case studies, we are working on more: suggest your own!
- ▶ Future Work:
 - ▶ support more C features and other languages
function pointers, union types, switch, variable length arrays, abrupt termination and jumps
 - ▶ increase expressiveness towards full separation logic
 - ▶ make it easy for users to implement new transformations
 - ▶ extensions to guide users towards beneficial transformations

Conclusion

- ▶ OptiTrust is an interactive framework for optimizing general-purpose C code
- ▶ Transformations are validated using a static resource analysis
- ▶ 2 completed case studies, we are working on more: suggest your own!
- ▶ Future Work:
 - ▶ support more C features and other languages
function pointers, union types, switch, variable length arrays, abrupt termination and jumps
 - ▶ increase expressiveness towards full separation logic
 - ▶ make it easy for users to implement new transformations
 - ▶ extensions to guide users towards beneficial transformations

Thanks,
we are open to collaboration!

 github.com/charguer/optitrust
 chargeraud.org/softs/optitrust/traces/

Backup Slides

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