A Term Rewriting Path to High-Performance

Thomas KŒHLER Phil TRINDER Michel STEUWER



INRIA CAMUS, Strasbourg – August 2022

Optimizing Low-Level Code is Hard

- hand optimization is time-consuming and error-prone e.g. in C, OpenCL, CUDA
- critical in performance-demanding domains
 e.g. image processing, numeric simulation, machine learning
- ► typically leads to orders of magnitude performance improvements

Optimizing Matrix Multiplication in C

```
for (int im = 0; im < m; im++) {
  for (int in = 0; in < n; in++) {
    float acc = 0.0f;
    for (int ik = 0; ik < k; ik++) {
        acc += a[ik + (k * im)] * b[in + (n * ik)];
    }
    output[in + (n * im)] = acc;
  }
}</pre>
```

Optimized program on the right:

+ $110 \times$ faster runtime

Intel i5-4670K CPU

 6× more lines, more complex code threads, SIMD, indexing

```
float aT[n + k].
#pragma onp parallel for
for (int in = 0; in \leq (n / 32); in = 1 + in) {
  for (int ik = 0: ik < k: ik = 1 + ik) {
    #pragma omp simd
    for (int in = 0; in < 32; in = 1 + in) {
      aT[(ik + ((32 + in) + k)) + (in + k)] = a[(in + (32 + in)) + (ik + n)];
#pragma onp parallel for
for (int im = 0: im < (m / 32): im = 1 + im) {
  for (int in = 0; in < (n / 32); in = 1 + in) {
    float tmp1[1024]:
    for (int jm = 0; jm < 32; jm = 1 + jm)
      for (int jn = 0; jn < 32; jn = 1 + jn) {
         tmp1[jn + (32 * jm)] = 0.0f;
    for (int ik = 0; ik < (k / 4); ik = 1 + ik) {
      for (int im = 0: im < 32: im = 1 + im) {
         float tmp2[32];
         for (int in = 0: in < 32: in = 1 + in) {
           tmp2[jn] = tmp1[jn + (32 * jm)];
         Apragma omp sind
         for (int in = 0; in < 22; in = 1 + in) (
           tmp2[in] += (a[((4 * ik) + ((32 * im) * k)) + (im * k)] * aT[((4 * ik) + ((32 * in) * k)) + (in * k)]);
         Noragma ono sind
         for (int jn = 0; jn < 32; jn = 1 + jn) {
   tmp2[in] += (a[((1 + (4 * ik)) + ((32 * im) * k)) + (jm * k)] *</pre>
             aT[((1 + (4 + ik)) + ((32 + in) + k)) + (jn + k)]);
         #pragma omp sind
         for (int in = 0; in < 32; in = 1 + in) {
            \begin{array}{l} tmp2[jn] += (a[(2 + (4 + ik)) + ((32 + im) + k)) + (jm + k)] \\  aT[(2 + (4 + ik)) + ((32 + in) + k)) + (jm + k)] \\ \end{array} 
         #pragma omp simd
         for (int jn = 0; jn < 32; jn = 1 + jn) {
    tmp2[jn] += (a[((3 + (4 + ik)) + ((32 + im) + k)) + (jm + k)] +</pre>
             aT[((3 + (4 + ik)) + ((32 + in) + k)) + (in + k)]);
         for (int in = 0; in < 32; in = 1 + in) {
           tmp1[in + (32 * im)] = tmp2[in];
    for (int im = 0; im < 32; im = 1 + im)
      for (int j_n = 0; j_n < 32; j_n = 1 + j_n) {
output[((j_n + ((32 * im) * n)) + (j_2 * in)) + (j_m * n)] = tmp1[j_n + (32 * j_m)];
```

Automating Optimization via Term Rewriting



- + convenient, hardware agnostic programming
- + high-performance code generation
- + extensible set of abstractions and optimizations

RISE is a Functional Array Language

High-level matrix multiplication in RISE:

```
def mm a b =
    map (λaRow.
        def mm a b =
    map (λaRow.
        dot aRow bCol)
        (transpose b)) a

def dot xs ys =
    reduce + 0
        (ap (λ(x, y), x × y))
        (zip xs ys))

for aRow in a:
    for aRow in a:
    for bCol in transpose(b):
        ... = dot(aRow, bCol)
        ... = dot(aRow, bCol)
        (aranspose b)) a
```

Rewrite Rules Encode Valid RISE Transformations

	map f x	for m: = f()
split-join:	→ join (map (map f) (split n x))	<pre>for m / n: for n: for = f()</pre>

transpose-around	<i>l-map-map</i> :
------------------	--------------------

map (map f) x	for m: for n: = f()	
\mapsto transpose		
(map	for n:	
(map f)	for m:	
(transpose x))	I = f()	

Complex Optimizations Emerge from Simple Rules

Matrix multiplication blocking in RISE:

 \mapsto^*

map (λ aRow.	for m:
map (λbCol.	for n:
dot aRow bCol)	for k:
(transpose b)) a	1

```
join (map (map join) (map transpose
  map
                                    for m / 32:
                                       for n / 32:
    (map \lambda x_2.
       reduceSeg (\lambda x_3, \lambda x_4.
                                        for k / 4:
         reduceSeq \lambda x_5, \lambda x_6.
                                         for 4:
                                          for 32:
           map
              (map (\lambda x_7.
                                           for 32:
               (fst x7) + (fst (snd x7)) \times
                   (snd (snd x7)))
                (map (\lambda x_7, zip (fst x_7) (snd x_7))
                  (zip x5 x6)))
          (transpose (map transpose
          (snd (unzip (map unzip map (\lambdax5.
            zip (fst x5) (snd x5))
            (zip x3 x4)))))))
         (generate (\lambdax3. generate (\lambdax4. 0)))
         transpose (map transpose x2))
    (map (map (map (map (split 4))))
       (map transpose
         (map (map (\lambdax2. map (map (zip x2))
           (split 32 (transpose b)))))
              split 32 a))))))
```

Deciding which Rewrites to Apply is Hard



ELEVATE Rewriting Strategies

- ▶ programmers describe optimizations as compositions of rewrite rules
- ► MM blocking:

- 1 def blocking = (baseline ';'
 2 tile(32,32) '@' outermost(mapNest(2)) ';;'
 3 fissionReduceMap '@' outermost(appliedReduce)';;'
 4 split(4) '@' innermost(appliedReduce)';;'
 5 reorder(List(1,2,5,6,3,4)))
- + empowers programmers to manually control the rewrite process
- + to define their own abstractions: tile, split, reorder are not built-in

Bastian Hagedorn, Johannes Lenfers, Thomas Koehler, Xueying Qin, Sergei Gorlatch, and Michel Steuwer. "Achieving high-performance the functional way: a functional pearl on expressing high-performance optimizations as rewrite strategies". In: ICFP (2020)

ELEVATE Rewriting Strategies

- ▶ programmers describe optimizations as compositions of rewrite rules
 - 1 def blocking = (baseline ';'
 2 tile(32,32) '@' outermost(mapNest(2)) ';;'
 3 fissionReduceMap '@' outermost(appliedReduce)';;'
 4 split(4) '@' innermost(appliedReduce)';;'
 5 reorder(List(1,2,5,6,3,4)))
- requires programmers to order all rewrite steps
- strategies are often program-specific and tedious to implement

► MM blocking:

Achieving High-Performance with ELEVATE

Case Study 1) Matrix Multiplication Optimizations for Intel CPU

- ► Transform loops blocking, permutation, unrolling
- Change data layout array packing
- ► Add parallelism vectorization, multi-threading
- Performance is on par with reference schedules from TVM



Achieving High-Performance with **ELEVATE**

Case Study 2) Harris Corner Detection Optimizations for ARM CPU

► 4 optimizations from reference Halide schedule

 $circular\ buffering,\ operator\ fusion,\ multi-threading,\ vectorization$

► 2 optimizations not supported by Halide convolution separation, register rotation



Thomas Koehler and Michel Steuwer. "Towards a Domain-Extensible Compiler: Optimizing an Image Processing Pipeline on Mobile CPUs". In: *CGO*. 2021

Rewriting strategies achieve high-performance, but require **tedious manual rewrite ordering**

Equality Saturation



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

Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. "egg: fast and extensible equality saturation". In: POPL (2021)

Equality Saturation



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

No manual rewrite ordering, but does not scale to the RISE case studies

Sketch-Guided Equality Saturation

Observation:

► The *shape* of the optimized program is often used to explain optimizations:

for m:	for m	/ 32:
for n:	, * for	n / 32:
for k:	→ for	k / 4:
	fo	r 4:
	f	or 32:
		for 32:
		••



Sketch-Guided Equality Saturation

Observation:

► The *shape* of the optimized 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

https://arxiv.org/abs/2111.13040

► *Sketches* are program patterns that leave details unspecified

baseline sketch:

containsMap(m,	for m:
<pre>containsMap(n,</pre>	for n:
containsReduceSeq(k,	for k:
<pre>containsAddMul)))</pre>	+ ×

• Abstractions defined in terms of smaller building blocks:

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

► *Sketches* are program patterns that leave details unspecified

baseline sketch:

<pre>containsMap(m,</pre>	for m:
<pre>containsMap(n,</pre>	for n:
containsReduceSeq(k,	for k:
<pre>containsAddMul)))</pre>	+ ×

► A sketch s is satisfied by a set of terms R(s):

► *Sketches* are program patterns that leave details unspecified

baseline sketch:

<pre>containsMap(m,</pre>	for m:
<pre>containsMap(n,</pre>	for n:
containsReduceSeq(k,	for k:
<pre>containsAddMul)))</pre>	+ ×

<i>blocking</i> sketch: <i>blocking</i> sketch: <i>blocking</i> sketch: <i>containsReduceSeq(k / 4, containsReduceSeq(k / 4, containsReduceSeq(4, containsMap(32, containsMap(32, containsAddMul))))))</i>	<pre>for m / 32: for n / 32: for k / 4: for 4: for 32: for 32: +×</pre>
---	---

► *Sketches* are program patterns that leave details unspecified

<i>baseline</i> sketch:	<pre>containsMap(m, containsMap(n, containsReduceSeq(k, containsAddMul)))</pre>	for m: for n: for k: + ×
sketch guide: how to split the loops before reordering them?	<pre>containsMap(m / 32, containsMap(32, containsMap(1 / 32, containsMap(32, containsReduceSeq(k / 4, containsReduceSeq(4, containsAddMul))))))</pre>	<pre>for m / 32: for 32: for n / 32: for 32: for k / 4: for 4: + ×</pre>
blocking sketch:	<pre>containsMap(m / 32, containsMap(n / 32, containsReduceSeq(k / 4, containsReduceSeq(4, containsMap(32, containsMap(32, containsAddMul))))))</pre>	<pre>for m / 32: for n / 32: for k / 4: for 4: for 32: for 32: + ×</pre>

• Equality Saturation without Sketch Guides²:

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

► Sketch-Guided Equality Saturation³:

goal	sketch guides	found?	runtime	RAM
baseline	0	✓	0.5s	0.02 GB
blocking	1	✓	7s	0.3 GB
+ 5 others	2-3	✓	$\leq 7s$	$\leq 0.5 \text{ GB}$

²Intel Xeon E5-2640 v2
³AMD Ryzen 5 PRO 2500U

A Term Rewriting Path to High-Performance

• Equality Saturation without Sketch Guides²:

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

► Sketch-Guided Equality Saturation³:

goal	sketch guides	found?	runtime	RAM
baseline	0	✓	0.5s	0.02 GB
blocking	1	✓	7s	0.3 GB
+ 5 others	2-3	1	≤7s	$\leq 0.5 \text{ GB}$

Sketch-guidance enables to find all 7 optimization goals

²Intel Xeon E5-2640 v2 ³AMD Ryzen 5 PRO 2500U

• Equality Saturation without Sketch Guides²:

	go	bal	found?	runtin	ıe	RAM	ί .	
	ba	iseline	1	0.	5s	0.02 GE	5	
	bl	ocking	1		lh)	(35 GE	$\overline{\mathbf{D}}$	
	+	5 others	×	>35n	m	>60 GE	5	
Sketch-Guided Equality Saturation ³ : 582x						116x		
	goal	sketch	guides	found?	ru	ntime	RAM	
	baseline	C)	✓		0.5s	0.02 GB	J
	blocking	1		1		75	0.3 GB	\checkmark
	+ 5 others	2-	·3	1		$\leq 7s$	≤0.5 GB]

²Intel Xeon E5-2640 v2
³AMD Ryzen 5 PRO 2500U

A Term Rewriting Path to High-Performance

►

Sketches vs Full Program

all goals except baseline:sketch guidessketch goalsketch sizesprogram size1-317-1290-124

- ► sketches elide around 90% of the program
- sketches elide intricate details such as array reshaping patterns (e.g. split, join, transpose)

Conclusion

We talked about:

- ► The RISE language & SHINE compiler automating optimization via term rewriting
- ► ELEVATE *rewriting strategies* achieving high-performance by controlling rewriting
- ► *Sketch-guided equality saturation*, a novel, semi-automatic optimization technique

Conclusion

We talked about:

- ► The RISE language & SHINE compiler automating optimization via term rewriting
- ► ELEVATE *rewriting strategies* achieving high-performance by controlling rewriting
- ► *Sketch-guided equality saturation*, a novel, semi-automatic optimization technique

✓ thomas.koehler@thok.eu♦ thok.eu

Thanks!

rise-lang.orgelevate-lang.org

Future Work

- Combine sketch-guided with rewriting strategies
- Apply sketch-guiding to more diverse applications
- Improve automated search (rewrite rule scheduling, heuristics, pruning, leverage hardware knowledge)
- ► Can we synthesize sketch guides from a sketch goal?
- ► Use in an interactive optimization assistant

Sketch Definition

 $S ::= ? \mid F(S, ..., S) \mid contains(S)$

 $R(?) = T = \{F(t_1, ..., t_n)\}$ $R(F(s_1, ..., s_n)) = \{F(t_1, ..., t_n) \mid t_i \in R(s_i)\}$ $R(contains(s)) = R(s) \cup \{F(t_1, ..., t_n) \mid \exists t_i \in R(contains(s))\}$

```
def containsMap(n: NatSketch, f: Sketch): Sketch =
    contains(app(map :: ?t → n.?dt → ?y, f))
def containsReduceSeq(n: NatSketch, f: Sketch): Sketch =
    contains(app(reduceSeq :: ?t → ?t → n.?dt → ?t, f))
def containsAddMul: Sketch =
    contains(app(app(+, ?), contains(×)))
```

Sketch-Satisfying Equality Saturation



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

E-Graph Evolution



E-Graph Evolution



E-Graph Evolution



Sketches vs Full Program

goal	sketch guides	sketch goal	sketch sizes	program size
blocking	split	reorder ₁	7	90
vectorization	split + reorder ₁	lower ₁	7	124
loop-perm	split + reorder ₂	lower ₂	7	104
array-packing	split + reorder ₂ + store	lower ₃	7-12	121
cache-blocks	$split + reorder_2 + store$	lower ₄	7-12	121
parallel	$split + reorder_2 + store$	lower ₅	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)

$$(a*2)/2 \longrightarrow^* a$$



(a * 2)/2

$$(a*2)/2 \longrightarrow^* a$$



(a*2)/2 $x*2 \longrightarrow x \ll 1$

 $(a*2)/2 \longrightarrow^* a$



(a*2)/2 $x*2 \longrightarrow x \ll 1$ $(x*y)/z \longrightarrow x*(y/z)$

A Term Rewriting Path to High-Performance

 $(a*2)/2 \longrightarrow^* a$



cost = term size