Sketch-Guided Equality Saturation

Thomas KŒHLER Phil TRINDER Michel STEUWER



Shameless PLUG Seminar — February 2022

The Optimizing Compiler Challenge



- + convenient, hardware agnostic programming
- + high-performance execution

Domain-Specific Compilers



- fixed set of abstractions and optimizations
- need to design and maintain multiple compilers

Optimization decisions are encoded into a schedule

Schedules for Optimization

Halide algorithm: what to compute

blur_x(x, y) = (input(x-1, y) + input(x, y) + input(x+1, y))/3; blur_y(x, y) = (blur_x(x, y-1) + blur_x(x, y) + blur_x(x, y+1))/3;

Halide schedule: *how to optimize*

blur_y.tile(x, y, xi, yi, 256, 32)
 .vectorize(xi, 8).parallel(y);
blur_x.compute_at(blur_y, x).vectorize(x, 8);

Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. "Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines". In: *Acm Sigplan Notices* (2013)

A Domain-Extensible Compiler



+ extensible set of abstractions and optimizations

Optimization decisions are encoded via rewriting

Rewriting Strategies for Optimization



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)

Rewriting Strategies vs Schedules

Advantages of rewriting strategies:

- ► 5 principles: separate concerns, facilitate reuse, enable composition, allow reasoning, explicit by default
- in practice, rewriting strategies are a competitive alternative to schedules matrix multiplication and corner detection case studies

Rewriting Strategies vs Schedules

Advantages of rewriting strategies:

- ► 5 principles: separate concerns, facilitate reuse, enable composition, allow reasoning, explicit by default
- in practice, rewriting strategies are a competitive alternative to schedules matrix multiplication and corner detection case studies

But ... both are difficult to write! optimization strategy = schedule or rewriting strategy

Optimization Strategies are Difficult to Write

100	If while enhanced down who downloaded even down who endewides	1.4.4	// We can do 4-wide vectorized loads from shared memory if
100	77 This schedule fuses the deprivate convinto the pointwise	1.45	// we can do whate vectorized today from shared memory if
187	77 conv. The results of the depthwise conv are computed inside	140	// we direct the reduction toop by a factor of four above
105	// the outer of the two pointwise reduction loops.	140	// and stage the toads from the depthwise_convolved
169		147	// output.
110	Var xi, yi, di, dii, xii, yii;	148	
111	RVar ro, ri;	149	depthwise_convolved.in()
112		150	.in()
113	// The pointwise convolution kernel. Produces a 4x4 tile of output.	151	.compute_at(pointwise_compled, x)
114	Func(output)	152	.bound_extent(d, 4)
115	.tile({d, x, y}, {di, xi, yi}, {16, 4, 4})	153	.vectorize(d)
116	.tile({di, xi, yi}, {dii, xii, yii}, {1, 2, 2})	154	unrother
117	.gpu_threads(di, xi, yi)	155	* TONING COL
118	.fuse(y, b, b)	156	2115 +215
119	.gpu_blocks(d, x, b)	157	77 The dependence on the convolution kernel. Produces a 4x4 tile
120	.unroll(xii)	150	// Refrectiate state, storing the result in shared.
121	.unroll(yii)	ne ri	(Wetwise_convolved.in()
122	.unroll(dii);	100 hai	.compute_at(output, d)
123	. Je W	" 9/s	.tile({d, x, y}, {di, xi, yi}, {32, 4, 4}, TailStrategy::RoundUp)
124	pointwise_convolved.compute_at(output, di)	15.	.tile({di, xi, yi}, {dii, xii, yii}, {2, 2, 2})
125	reorder(x, y, d)	10 40	.gpu_threads(di, xi, yi)
126	unroll(x)	allos.	.unroll(xii)
127	unroll(y)	Oros	.unroll(yii)
128	unroll(d) ine? http://	166	.unroll(dii);
129	update()	167	
130		168	depthwise_convolved
131	unroll(y) 215	169	.compute_at(depthwise_convolved.in(), di)
132	.unroll(d)	170	.unroll(x)
133	<pre>split(rc, ro, ri, 4)</pre>	171	.unroll(y)
134	.reorder(ri, x, y, d, ro)	172	.unroll(d)
135	unroll(ri);	173	.update()
136		174	.reorder(d, x, y, rx, ry, rd)
137	// We're going to call in() on depthwise_convolved twice.	175	.unroll(x)
138	// The first will be to give it a wrapper to do the	176	.unroll(y)
139	// accumulation in registers before writing the result to	177	.unroll(d);
140	// shared. The second will be staging the loads from		
141	// shared into registers. We write them in reverse order		
142	// heless		

Optimization Strategies are Difficult to Write

```
def normForReorder(implicit ev: Traversable[Rise]): Strategy[Rise] =
         (splitBeforeMap '8' topDown[Rise]) ';;'
         (fuseReduceMap '8' topDown[Rise]) ';;'
         (fuseReduceMap '@' topDown[Rise]) ';;' RNF()
       #strategy def reorder(1: List[Int])(implicit ey: Traversable[Rise]): Strategy[Rise] = normForReorder ':' (reorderRec(1) '@' topDown[Rise])
       Bstrategy def reorderRec(l: List[int])(implicit ev: Traversable[Rise]): Strategy[Rise] = e => {
         def freduce(s: Strategy[Rise]): Strategy[Rise] =
           function(function(argumentOf(reduceSeg.primitive, body(body(s)))))
         def freduceX(s: Strategy[Rise]): Strategy[Rise] =
           argument(function(function(argumentOf(reduceSeg.primitive, body(body(s))))))
         def stepDown(s: Strategy[Rise]): Strategy[Rise] = freduceX(s) <+ freduce(s) <+ fnap(s)
         val isFullyAppliedReduceSeq: Strategy[Rise] = isApplied(isApplied(isApplied(isReduceSeq))) <+</pre>
           argument(isApplied(isApplied(isApplied(isReduceSeg)))) // weird reduce special case
         val isEullyAppliedMap: Strateov(Risel = isApplied(isApplied(isMap)))
         def moveReductionUp(pos: Int): Strategy[Rise] = {
           if (pos <= 1) id
              applyNTimes(pos-2)(stepDown)(function(liftReduce)) ':' DENE() ':' RNE() ':' moveReductionUp(pos-1)
                                                              -40 lines of ELEVATE stratesy
to define loop reordering
for matrix multiplication
         1 match (
           // nothing to reorder, so further down
           case x :: xs if x == 1 => (stepDown(reorderRec(xs.map(y => if (y>x) y-1 else y ))))(e)
           // work to do
104
           case pos :: xs => (
             (applyNTimes(pos-1)(stepDown)(isFullyAppliedReduceSeg) ::
              // move reduction and normalize the AST
               moveReductionUp(pos) ';'
               stepDown(reorderRec(xs.map(y => if (y>pos) y-1 else y )))) <+</pre>
               // other case: is it a man?
               annlyMTimes(nos.1)(stenDown)(isEullyAnnliedMan) '.'
               basic fail
              )(e)
           case Nil => id(e)
           case _ => Failure(reorderRec(1))
117
110
```

Optimization Strategies are Difficult to Write

► to write: which sequence of transformations lead to an optimized program?



▶ to read: what is the resulting program? what are the intermediate programs?



Transformed program is hidden state that needs to be reasoned about

Reasoning about Strategies with Program Shapes



(a) Definitions and Example GPU Schedule of Group KWZ: Computation of K has been moved at the block (innermost inter-tile level) of the output Z and W has been interleaved inside the thread level that computes W.



optimized

elided

details

program shape

Savvas Sioutas, Sander Stuijk, Twan Basten, Henk Corporaal, and Lou Somers. "Schedule synthesis for halide pipelines on gpus". In: TACO (2020)

Reasoning about Strategies with Program Shapes

Not only found in papers ... but also in talks:



Reasoning about Strategies with Program Shapes



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

Sketches to Formalize Program Shapes

sketches = program patterns that leave details unspecified

S ::= ? | F(S, .., S) | contains(S)

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

Basic sketch grammar (top) and the terms it represents (bottom)

Sketches to Formalize Program Shapes

sketches = program patterns that leave details unspecified

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

map(grayLine) ▷ slide(3, 1) ▷
map(sobelLine) ▷ slide(3, 1) ▷ map(coarsityLine)

Example informal program shape (middle) and possible sketch (bottom)

Sketches for Optimization

Why not use sketches to specify optimization goals?

▶ instead of writing optimization strategy, write desired program shape:



► instead of reading optimization strategy, read program shape:

$$\frac{\text{optimized}}{\text{program}} \in \mathsf{R}(\frac{\text{desired}}{\text{sketch}})$$

Automated Search?





- ► An e-graph efficiently represents a large set of equivalent programs.
- ► The e-graph is grown by applying all possible rewrite rules in a purely additive way.
- ► After growing the e-graph, the best program found is extracted.

Ross Tate, Michael Stepp, Zachary Tatlock, and Sorin Lerner. "Equality saturation: a new approach to optimization". In: *POPL*. 2009

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

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



(a * 2)/2

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



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

Sketch-Guided Equality Saturation

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



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

Sketch-Guided Equality Saturation

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



cost = term size

16

Equality Saturation with Sketches



Questions:

- 1. How does it work for functional programs like RISE?
 - ▶ no efficient support for name bindings, rewritten languages are usually first order
- 2. Does it scale to complex optimizations of realistic programs?
 - ► the search could be too costly

Equality Saturation for Functional Programs

Consider 2 standard λ -calculus rules + 2 rules that introduce names on the right:

 $\begin{array}{ccc} (\lambda x. \ b)e \longrightarrow b[e/x] & (\beta \text{-reduction}) \\ \lambda x. \ fx \longrightarrow f & \text{if } x \text{ not free in } f & (\eta \text{-reduction}) \\ map \ f(map \ g \ arg) \longrightarrow map \ (\lambda x. \ f(g \ x)) \ arg & (map-fusion) \\ map \ (\lambda x. \ fgx) \longrightarrow \lambda y. \ map \ f(map \ (\lambda x. \ gx) \ y) & \text{if } x \text{ not free in } f & (map-fission) \end{array}$

How can we implement substitution, predicates and name bindings?

Equality Saturation for Functional Programs

• We made substitution much more efficient, at the cost of ignoring possibilities. only substitute using one representative term by equivalence class

- ► State-of-the-art predicates are efficient but ignore possibilities. predicate has to hold ∀ terms in equivalence class
- We made name bindings much more efficient using DeBruijn indices. the goal is to avoid duplicating alpha-equivalent terms

Equality Saturation for Functional Programs

• We made substitution much more efficient, at the cost of ignoring possibilities. *discovering a rewrite goal requiring less than 10 rewrite rule applications:*

method	time	RAM	e-nodes	e-classes	found
before	40s	4GB	13M	3M	no
after	<1ms	MBs	364	277	yes

We made name bindings much more efficient using DeBruijn indices. discovering a rewrite goal requiring less than 50 rewrite rule applications:

method	time	RAM	e-nodes	e-classes	found
before	2mn	4GB	2.9M	1.4M	no
after	100ms	MBs	3K	1K	yes

Optimization Time and Memory Consumption

► ELEVATE strategies take ~1s to execute, reproducing 7 optimizations from TVM.

version	sketches	found	time	RAM	rules
baseline	1	yes	0.5s	20 MB	2
blocking	1	yes	1h+	35GB	5M
+ 5 others	1	no	~1h+	60GB+	???

• Equality Saturation with Sketches¹:

¹Intel Xeon E5-2640 v2

Optimization Time and Memory Consumption

► ELEVATE strategies take ~1s to execute, reproducing 7 optimizations from TVM.

►	Equality	Saturation	with	Sketches ¹ :	
---	----------	------------	------	-------------------------	--

version	sketches	found	time	RAM	rules
baseline	1	yes	0.5s	20 MB	2
blocking	1	yes	1h+	35GB	5M
+ 5 others	1	no	~1h+	60GB+	???

What else can we do to improve scaling?

¹Intel Xeon E5-2640 v2

Sketches



Intermediate Sketches

baseline version:

Intermediate sketch:

how to split the loops?

blocking version (3D):

<pre>containsMap(m,</pre>	m: r n: for k:
<pre>containsMap(m / 32, containsMap(n / 32, containsMap(n / 32, containsMap(32, containsReduceSeq(k / 4, containsReduceSeq(4, containsAddMul))))))</pre>	for m / 32: for 32: for n / 32: for 32: for k / 4: for 4: * *
<pre>containsMap(m / 32, containsMap(n / 32, containsReduceSeq(k / 4, containsReduceSeq(4, containsMap(32, con</pre>	for m / 32: for n / 32: for k / 4: for 4: for 32: for 32:

Sketch-Guided Equality Saturation



- decompose a complex search into a series of simpler searches
- additional sketches specify intermediate goals

Intermediate Sketches

version	sketches
blocking	split + reorder ₁
vectorization	split + reorder ₁ + lower ₁
loop-perm	split + reorder ₂ + lower ₂
array-packing	split + reorder ₂ + store + lower ₃
cache-blocks	split + reorder ₂ + store + lower ₄
parallel	split + reorder ₂ + store + lower ₅

Decomposition of each version into logical steps. A sketch is defined for each logical step.

Optimization Time and Memory Consumption

- ► ELEVATE strategies take ~1s to execute, reproducing 7 optimizations from TVM.
- Equality Saturation with Sketches²:

version	sketches	found	time	RAM	rules
baseline	1	yes	0.5s	20 MB	2
blocking	1	yes	1h+	35GB	5M
+ 5 others	1	no	~1h+	60GB+	???

► Sketch-Guided Equality Saturation³:

version	sketches	found	time	RAM	rules
baseline	1	yes	0.5s	20 MB	2
blocking	2	yes	7s	0.3 GB	11K
+ 5 others	3-4	yes	<7s	<0.5 GB	<11K

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

Conclusion

We propose:

- ► sketches for program optimization, alternative to optimization strategies
- ▶ practical techniques to support efficient equality saturation for lambda calculi
- ► *sketch-guided equality saturation*, a novel, semi-automatic optimization technique

Conclusion

We propose:

- ► sketches for program optimization, alternative to optimization strategies
- ▶ practical techniques to support efficient equality saturation for lambda calculi
- ► *sketch-guided equality saturation*, a novel, semi-automatic optimization technique

✓ thomas.koehler@thok.eu	Thanks!	🔇 rise-lang.org
🔇 thok.eu	We are open to collaboration!	Selevate-lang.org

paper: https://arxiv.org/abs/2111.13040

Deciding How to Apply Rewrite Rules

Fully automated search?

e.g. heuristic search, equality saturation, ...



Manually written recipe?

e.g. Halide/TVM schedules, Elevate strategies, ...



How can I generalize this recipe to other cases? 🖓

