

Guided Equality Saturation

Thomas K  HLER

Andr  s GOENS

Siddharth BHAT

Tobias GROSSER

Phil TRINDER

Michel STEUWER

Inria



University
of Glasgow



UNIVERSITEIT
VAN AMSTERDAM



TECHNISCHE
UNIVERSIT  T
BERLIN



THE UNIVERSITY
of EDINBURGH



UNIVERSITY OF
CAMBRIDGE

POPL Conference — London, January 2024

The Limits of Greedy Term Rewriting

Example: Eliminating Intermediate Memory with Fusion

Rewrite Rules:

(1) $\text{map } a \circ \text{map } b \mapsto \text{map } (a \circ b)$

uses less memory

Program to Optimize:

$(\text{map } (\text{map } f)) \circ (\text{transpose} \circ (\text{map } (\text{map } g)))$



local optimum,
cannot use less memory?

The Limits of Greedy Term Rewriting

Example: Eliminating Intermediate Memory with Fusion

Rewrite Rules:

(1) $\text{map } a \circ \text{map } b \mapsto \underline{\text{map } (a \circ b)}$

uses less memory

(2) $\text{transpose} \circ \text{map } (\text{map } a) \mapsto \text{map } (\text{map } a) \circ \text{transpose}$

(3) $a \circ (b \circ c) \mapsto (a \circ b) \circ c$

no effect on memory

Program to Optimize:

$(\text{map } (\text{map } f)) \circ (\text{transpose} \circ (\text{map } (\text{map } g)))$

not explored by greedy rewriting: \downarrow (2); (3)

$((\text{map } (\text{map } f)) \circ (\text{map } (\text{map } g))) \circ \text{transpose}$

\downarrow (1); (1)

$\underline{(\text{map } (\text{map } (f \circ g))) \circ \text{transpose}}$

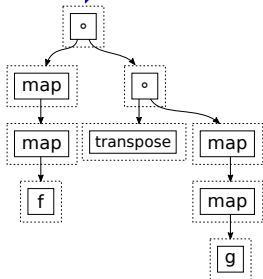
global optimum,
uses less memory!

Equality Saturation

Example: Eliminating Intermediate Memory with Fusion

$(\text{map } (\text{map } f)) \circ (\text{transpose} \circ (\text{map } (\text{map } g)))$

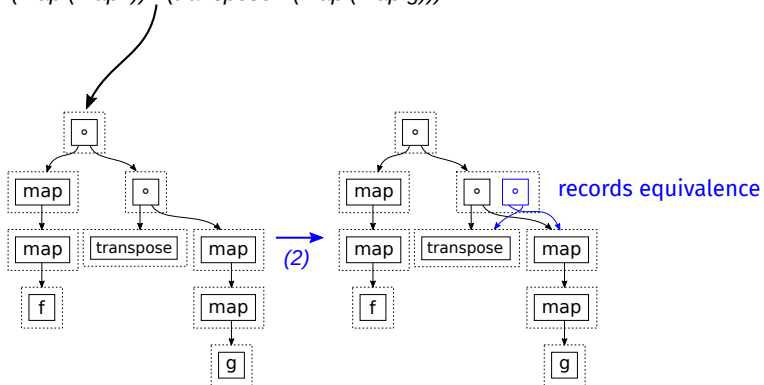
initialize equivalence graph
a.k.a. e-graph



Equality Saturation

Example: Eliminating Intermediate Memory with Fusion

$(\text{map } (\text{map } f)) \circ (\text{transpose} \circ (\text{map } (\text{map } g)))$

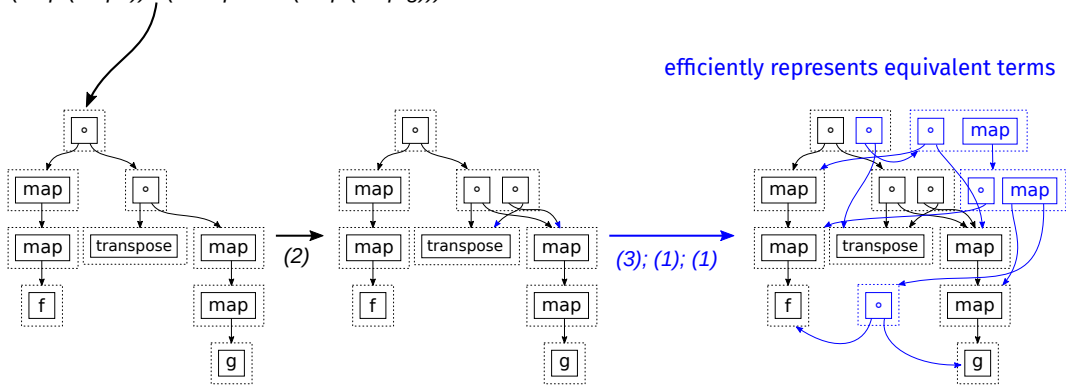


$$\begin{aligned} \text{transpose} \circ \text{map } (\text{map } g) \\ = \\ \text{map } (\text{map } g) \circ \text{transpose} \end{aligned}$$

Equality Saturation

Example: Eliminating Intermediate Memory with Fusion

$(\text{map } (\text{map } f)) \circ (\text{transpose} \circ (\text{map } (\text{map } g)))$



Equality Saturation

Example: Eliminating Intermediate Memory with Fusion

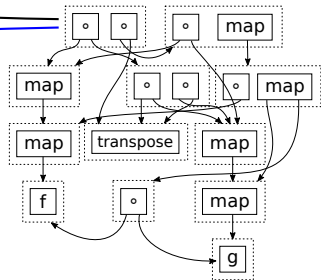
$(\text{map } (\text{map } f)) \circ (\text{transpose} \circ (\text{map } (\text{map } g)))$

extract term minimizing
memory usage

$(\text{map } (\text{map } (f \circ g))) \circ \text{transpose}$

global optimum,

found by equality saturation [Tate et al. 2009; Willsey et al. 2021]



The Limits of Equality Saturation

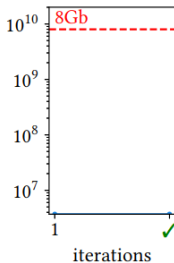
Example: Improving Memory Access Patterns with Tiling

(1),(2),(3) + (4) $\text{map } n_1 (\text{map } n_2 f) \mapsto \text{transpose} \circ (\text{map } n_2 (\text{map } n_1 f)) \circ \text{transpose}$
(5) $\text{map } (n_1 \times n_2) f \mapsto \text{join} \circ (\text{map } n_1 (\text{map } n_2 f)) \circ (\text{split } n_2)$

1D Tiling:

$\text{map } (n_1 \times 32) f$
↓ easy!
 $\text{join} \circ (\text{map } n_1 (\text{map } 32 f)) \circ (\text{split } 32)$

Memory Footprint (bytes):



1 iteration =
apply rules breadth-first

The Limits of Equality Saturation

Example: Improving Memory Access Patterns with Tiling

(1),(2),(3) + (4) $\text{map } n_1 (\text{map } n_2 f) \mapsto \text{transpose} \circ (\text{map } n_2 (\text{map } n_1 f)) \circ \text{transpose}$
(5) $\text{map } (n_1 \times n_2) f \mapsto \text{join} \circ (\text{map } n_1 (\text{map } n_2 f)) \circ (\text{split } n_2)$

2D Tiling:

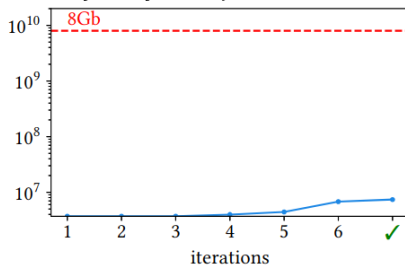
$\text{map } (n_1 \times 32) (\text{map } (n_2 \times 32) f)$



more challenging ...

$\text{join} \circ (\text{map } n_1 (\text{map } 32 \text{ join})) \circ (\text{map } n_1 \text{ transpose}) \circ$
 $(\text{map } n_1 (\text{map } n_2 (\text{map } 32 (\text{map } 32 f)))) \circ$
 $(\text{map } n_1 \text{ transpose}) \circ (\text{map } n_1 (\text{map } 32 (\text{split } 32))) \circ (\text{split } 32)$

Memory Footprint (bytes):



The Limits of Equality Saturation

Example: Improving Memory Access Patterns with Tiling

(1),(2),(3) + (4) $\text{map } n_1 (\text{map } n_2 f) \mapsto \text{transpose} \circ (\text{map } n_2 (\text{map } n_1 f)) \circ \text{transpose}$
(5) $\text{map } (n_1 \times n_2) f \mapsto \text{join} \circ (\text{map } n_1 (\text{map } n_2 f)) \circ (\text{split } n_2)$

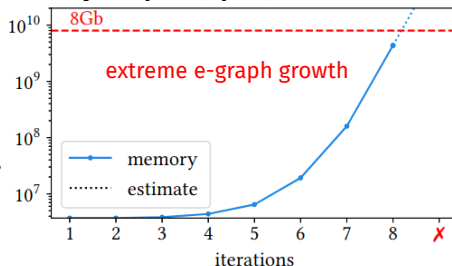
3D Tiling:

$\text{map } (n_1 \times 32) (\text{map } (n_2 \times 32) (\text{map } (n_3 \times 32) f))$

↓ unreachable with 8 Gb!

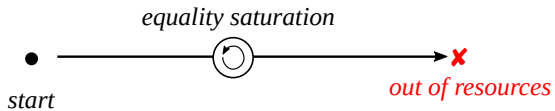
$(\text{map } (n_1 \times 32) (\text{map } (n_2 \times 32) \text{join}) \circ \text{join}) \circ \text{join} \circ$
 $(\text{map } n_1 \text{transpose} \circ (\text{map } n_2 (\text{map } 32 \text{transpose}) \circ \text{transpose})) \circ$
 $(\text{map } n_1 (\text{map } n_2 (\text{map } n_3 (\text{map } 32 (\text{map } 32 (\text{map } 32 f)))))) \circ$
 $(\text{map } n_1 (\text{map } n_2 \text{transpose})) \circ (\text{map } n_1 (\text{map } n_2 (\text{map } 32 \text{transpose})) \circ \text{transpose}) \circ$
 $(\text{split } 32) \circ (\text{map } (n_1 \times 32) (\text{map } n_2 (\text{map } 32 (\text{split } 32)))) \circ (\text{split } 32)$

Memory Footprint (bytes):



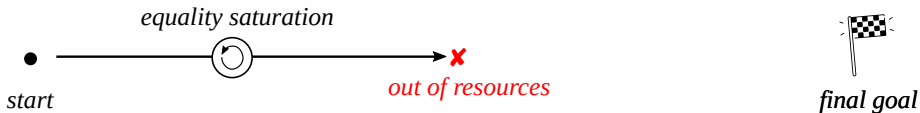
Guided Equality Saturation

unguided:

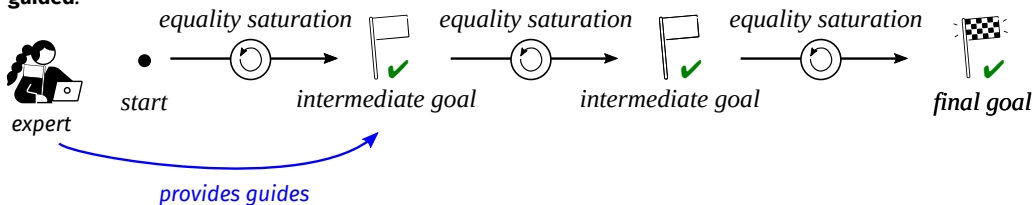


Guided Equality Saturation

unguided:



guided:



Guided Equality Saturation

Example: Improving Memory Access Patterns with Tiling (3D)

$map\ (n_1 \times 32)\ (map\ (n_2 \times 32)\ (map\ (n_3 \times 32)\ f))$

1. split the loops

$(map\ (n_1 \times 32)\ (map\ (n_2 \times 32)\ join) \circ join) \circ join \circ$
 $(map\ n_1\ (map\ 32\ (map\ n_2\ (map\ 32\ (map\ n_3\ (map\ 32\ f)))))) \circ$
 $(split\ 32) \circ (map\ (n_1 \times 32)\ (map\ n_2\ (map\ 32\ (split\ 32)))) \circ (split\ 32))$

guide provides insight

2. reorder them

$(map\ (n_1 \times 32)\ (map\ (n_2 \times 32)\ join) \circ join) \circ join \circ$
 $(map\ n_1\ transpose \circ (map\ n_2\ (map\ 32\ transpose) \circ transpose)) \circ$
 $(map\ n_1\ (map\ n_2\ (map\ n_3\ (map\ 32\ (map\ 32\ f)))))) \circ$
 $(map\ n_1\ (map\ n_2\ transpose)) \circ (map\ n_1\ (map\ n_2\ (map\ 32\ transpose)) \circ transpose) \circ$
 $(split\ 32) \circ (map\ (n_1 \times 32)\ (map\ n_2\ (map\ 32\ (split\ 32)))) \circ (split\ 32))$

Guided Equality Saturation

Example: Improving Memory Access Patterns with Tiling (3D)

$\text{map } (n_1 \times 32) \text{ (map } (n_2 \times 32) \text{ (map } (n_3 \times 32) \text{ f))}$

1. split the loops

```
(contains  
(map n1 (map 32 (map n2 (map 32 (map n3 (map 32 f))))))  
)
```

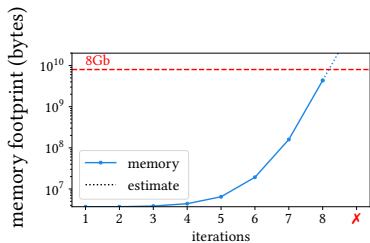
guide is more or less
precise sketch

2. reorder them

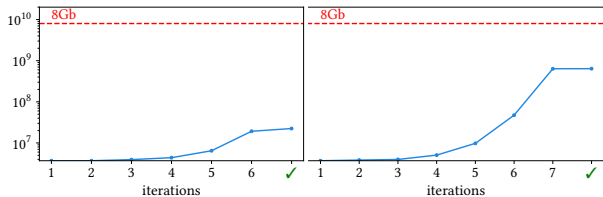
```
(map (n1 × 32) (map (n2 × 32) join) ◦ join) ◦ join ◦  
(map n1 transpose ◦ (map n2 (map 32 transpose) ◦ transpose)) ◦  
(map n1 (map n2 (map n3 (map 32 (map 32 (map 32 f)))))) ◦  
(map n1 (map n2 transpose)) ◦ (map n1 (map n2 (map 32 transpose)) ◦ transpose) ◦  
(split 32) ◦ (map (n1 × 32) (map n2 (map 32 (split 32))) ◦ (split 32))
```

Guided Equality Saturation

Example: Improving Memory Access Patterns with Tiling (3D)



(a) Equality saturation (found: X)



(b) **Guided** equality saturation (found: ✓)

- A single guide makes 3D Tiling reachable with 8Gb!

Case Study: Program Optimization

- ▶ We reproduced Matrix Multiplication optimizations from TVM:

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

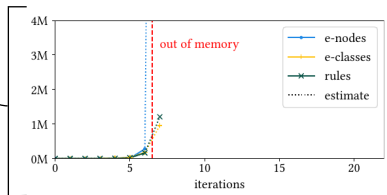
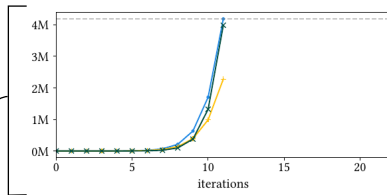
- ▶ transform loops *blocking, permutation, unrolling*
 - ▶ change data layout
 - ▶ add parallelism *vectorization, multi-threading*
- ▶ Prior work performs them by manually composing rewrite rules [*ICFP 2020; CACM 2023*]

Case Study: Program Optimization

Unguided Runtime and Memory Consumption

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

5 goals are too hard to find with
unguided equality saturation



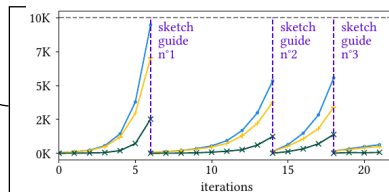
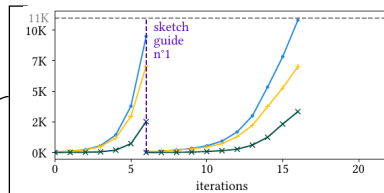
Case Study: Program Optimization

Guided Runtime and Memory Consumption

goal	sketch guides	found?	runtime	RAM
baseline	0	✓	0.5 s	0.02 GB
blocking	1	✓	7 s	0.3 GB
vectorization	2	✓	7 s	0.4 GB
loop-perm	2	✓	4 s	0.3 GB
array-packing	3	✓	5 s	0.4 GB
cache-blocks	3	✓	5 s	0.5 GB
parallel	3	✓	5 s	0.4 GB

582x faster 116x less

All found with up to 3 guides,
eliding 90% of the complete program



Case Study: Theorem Proving

- We implemented a **ges** tactic for the Lean theorem prover:

$$\begin{aligned} g^{-1-1} &= \boxed{g^{-1-1} * (g^{-1} * g)} \\ - &= g \end{aligned} \quad \begin{array}{l} \text{key} \\ \text{reasoning step} \end{array}$$

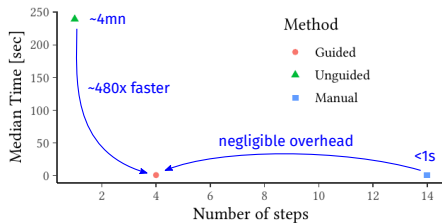
- Steps and details are omitted:

$$\begin{aligned} (g^{-1})^{-1} &\xrightarrow{\text{mul. one}} (g^{-1})^{-1} \cdot 1 \xrightarrow{\text{mul. inverse}} (g^{-1})^{-1} \cdot (g^{-1} \cdot g) \\ &\xrightarrow{\text{mul. assoc.}} ((g^{-1})^{-1} \cdot g^{-1}) \cdot g \xrightarrow{\text{mul. inverse}} 1 \cdot g \xrightarrow{\text{mul. one}} g \end{aligned}$$

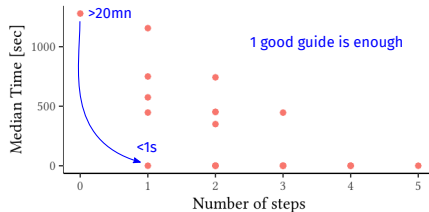
Case Study: Theorem Proving

Proving Theorems on Rings of Characteristic 2*

$$(x + y)^2 = x^2 + y^2$$



$$(x + y)^3 = x^3 + x \cdot y^2 + x^2 \cdot y + y^3$$



* $1 + 1 = 0, x + x = 0$

Conclusion

- ▶ **Guided Equality Saturation** offers an effective trade-off between manual and automated rewriting
- ▶ For program optimization, guides resemble explanatory code snippets
- ▶ For theorem proving, guides resemble key reasoning steps from textbooks
- ▶ More details in our paper, supplementary material and open-source code!

Conclusion

- ▶ **Guided Equality Saturation** offers an effective trade-off between manual and automated rewriting
- ▶ For program optimization, guides resemble explanatory code snippets
- ▶ For theorem proving, guides resemble key reasoning steps from textbooks
- ▶ More details in our paper, supplementary material and open-source code!

✉ thomas.koehler@thok.eu

🌐 thok.eu

Thanks!

Sketches

- *Sketches* are program patterns that leave details unspecified

baseline sketch:

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

- 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:

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

- 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:

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

sketch guide:

how to split the loops before reordering them?

<code>containsMap(m / 32, containsMap(32, containsMap(n / 32, containsMap(32, containsReduceSeq(k / 4, containsReduceSeq(4, containsAddMul))))))</code>	<code>for m / 32: for 32: for n / 32: for 32: for k / 4: for 4: .. + .. × ..</code>
---	---

blocking sketch:

<code>containsMap(m / 32, containsMap(n / 32, containsReduceSeq(k / 4, containsReduceSeq(4, containsMap(32, containsMap(32, containsAddMul))))))</code>	<code>for m / 32: for n / 32: for k / 4: for 4: for 32: for 32: .. + .. × ..</code>
---	---

Sketch Definition

$$S ::= ? \mid F(S, \dots, S) \mid \text{contains}(S)$$

$$R(?) = T = \{F(t_1, \dots, t_n)\}$$

$$R(F(s_1, \dots, s_n)) = \{F(t_1, \dots, t_n) \mid t_i \in R(s_i)\}$$

$$R(\text{contains}(s)) = R(s) \cup \{F(t_1, \dots, t_n) \mid \exists t_i \in R(\text{contains}(s))\}$$

```
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(×)))
```

Rewritten Language

- 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
```

Sketches vs Full Program

goal	sketch guides	sketch goal	sketch sizes	program size
blocking	split	reorder₁	7	90
vectorization	split + reorder₁	lower₁	7	124
loop-perm	split + reorder₂	lower₂	7	104
array-packing	split + reorder₂ + store	lower₃	7-12	121
cache-blocks	split + reorder₂ + store	lower₄	7-12	121
parallel	split + reorder₂ + store	lower₅	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**)

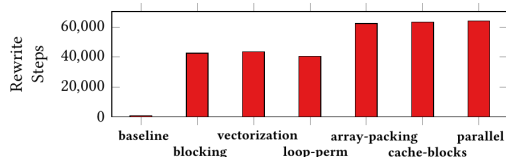
Difficulty 1. Long Rewrite Sequences

```

map (λaRow.      | for m:
  map (λbCol.    |   for n:
    dot aRow bCol) |   for k:
    (transpose b)) a | ...
  
```

↔ *

Prior work (*not shortest path*):



```

join (map (map join) (map transpose
  map (map λx2.      | for m / 32:
    reduceSeq (λx3. λx4. |   for n / 32:
      reduceSeq λx5. λx6. |   for k / 4:
        map      |   for 4:
          (map (λx7.   |   for 32:
            (fst x7) + (fst (snd x7)) ×
              (snd (snd x7)))
            (map (λx7. zip (fst x7) (snd x7))
              (zip x5 x6)))
            (transpose (map transpose
              (snd (unzip (map unzip map (λx5.
                zip (fst x5) (snd x5))
                (zip x3 x4)))))))
            (generate (λx3. generate (λx4. 0)))
            transpose (map transpose x2))
            (map (map (map (map (split 4))))
              (map transpose
                (map (map (λx2. map (map (zip x2)
                  (split 32 (transpose b))))
                    split 32 a))))))
  
```

Difficulty 2. Explosive Combinations of Rewrite Rules

Two example rules that quickly generate many possibilities:

split-join:

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

transpose-around-map-map:

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

Handwritten Matrix Multiplication

```
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;  
    }  
}
```

Optimised program on the right:

+ 110× faster runtime

Intel i5-4670K CPU

- 6× more lines of code where things
can go wrong

threads, SIMD, index computations

- hardware specific (not portable)

```
float aT[n * k];  
#pragma omp parallel for  
for (int in = 0; in < (n / 32); in = 1 + in) {  
    for (int ik = 0; ik < k; ik = 1 + ik) {  
        #pragma omp simd  
        for (int jn = 0; jn < 32; jn = 1 + jn) {  
            aT[(ik + ((32 * in) * k)) + (jn * k)] = a[(jn + (32 * in)) + (ik * n)];  
        }  
    }  
}  
#pragma omp 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 jm = 0; jm < 32; jm = 1 + jm) {  
                float tmp2[32];  
                for (int jn = 0; jn < 32; jn = 1 + jn) {  
                    tmp2[jn] = tmp1[jn + (32 * jm)];  
                }  
                #pragma omp simd  
                for (int jn = 0; jn < 32; jn = 1 + jn) {  
                    tmp2[jn] += (a[(((4 * ik) + ((32 * im) * k)) + (jm * k))] * aT[(((4 * ik) + ((32 * in) * k)) + (jn * k))]);  
                }  
                #pragma omp simd  
                for (int jn = 0; jn < 32; jn = 1 + jn) {  
                    tmp2[jn] += (a[(((2 * (4 * ik)) + ((32 * im) * k)) + (jm * k))] * aT[(((2 * (4 * ik)) + ((32 * in) * k)) + (jn * k))]);  
                }  
                #pragma omp simd  
                for (int jn = 0; jn < 32; jn = 1 + jn) {  
                    tmp2[jn] += (a[(((3 * (4 * ik)) + ((32 * im) * k)) + (jm * k))] * aT[(((3 * (4 * ik)) + ((32 * in) * k)) + (jn * k))]);  
                }  
                for (int jn = 0; jn < 32; jn = 1 + jn) {  
                    tmp1[jn + (32 * jm)] = tmp2[jn];  
                }  
            }  
        }  
        for (int jm = 0; jm < 32; jm = 1 + jm) {  
            for (int jn = 0; jn < 32; jn = 1 + jn) {  
                output[(((jn + ((32 * im) * n)) + (32 * in)) + (jm * n))] = tmp1[jn + (32 * jm)];  
            }  
        }  
    }  
}
```