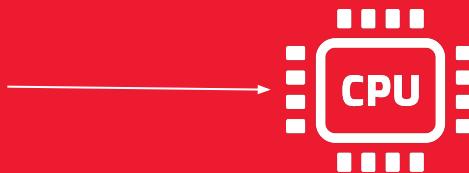


# Latent Idiom Recognition for a Minimalist Functional Array Language using Equality Saturation

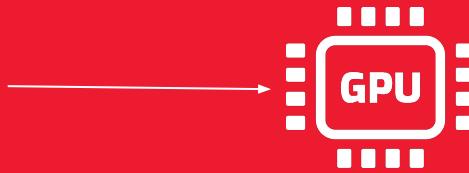
## No library

```
for (i = 0; i < N; i++) {  
    sum += xs[i];  
}
```



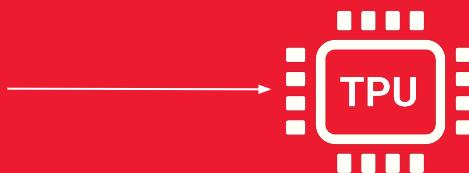
## BLAS

```
double ones[N] = { 1, 1, ... };  
ddot(N, xs, 1, ones, 1);
```



## PyTorch

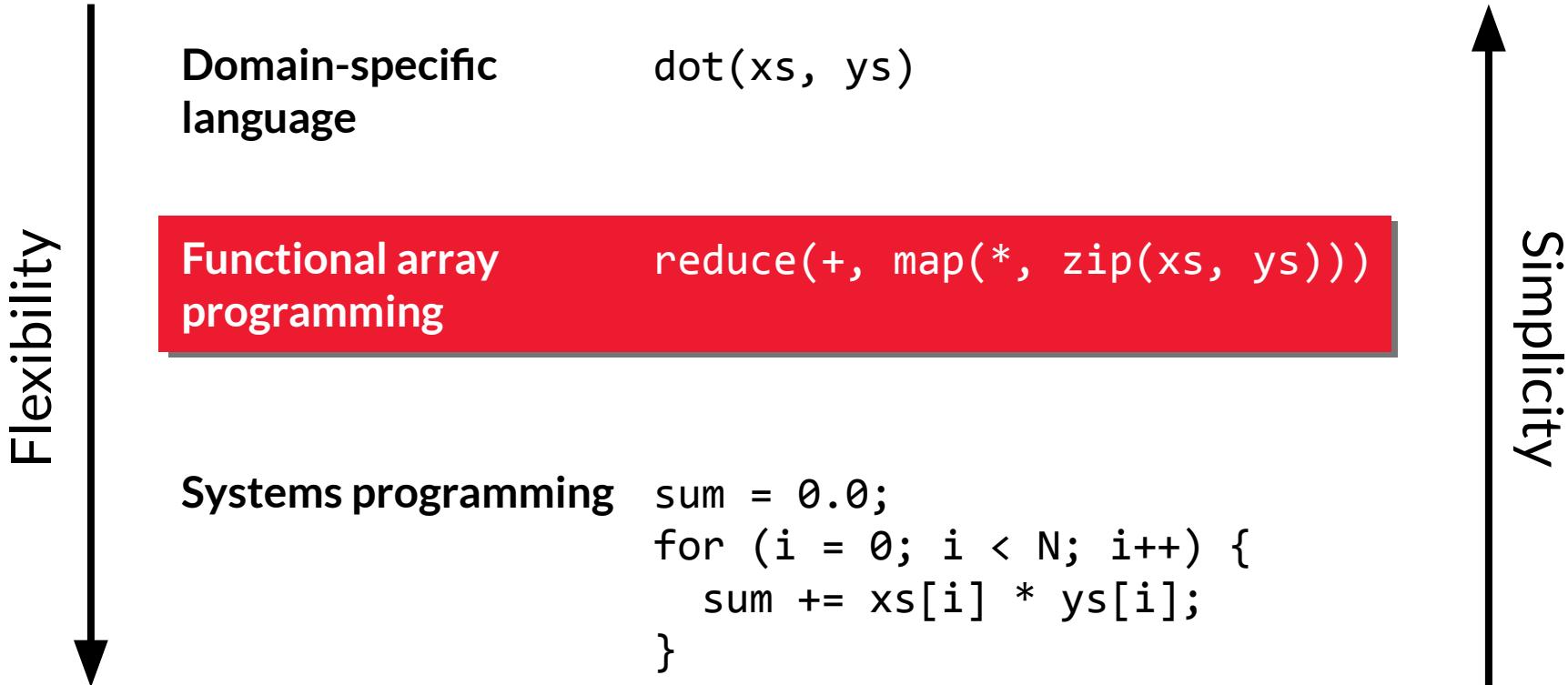
```
dot(xs, full(N, 1))  
-OR-  
sum(xs)
```



# McGill

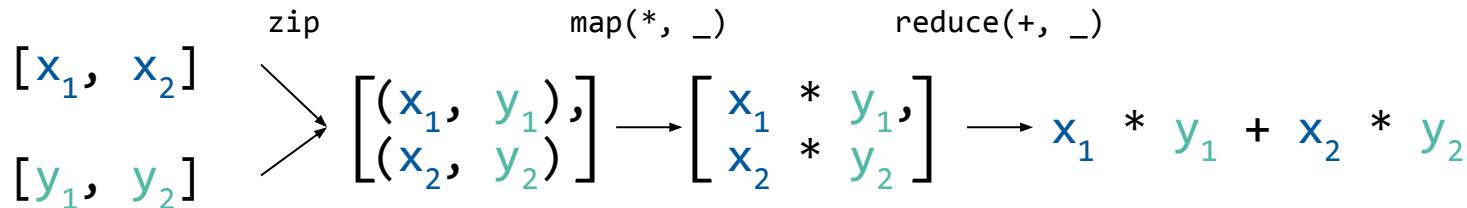
Jonathan Van der Cruysse and Christophe Dubach  
McGill University - CGO '24 - Tuesday March 5, 2024

# Levels of abstraction



# Functional array programming

```
dot(xs, ys) = reduce(+, map(*, zip(xs, ys)))
```



# Idiom recognition



```
dot(xs, ys) = reduce(+, map(*, zip(xs, ys)))
```



Dot product

```
reduce(+, map(*, zip(xs, ys)))
```

→ dot(xs, ys)



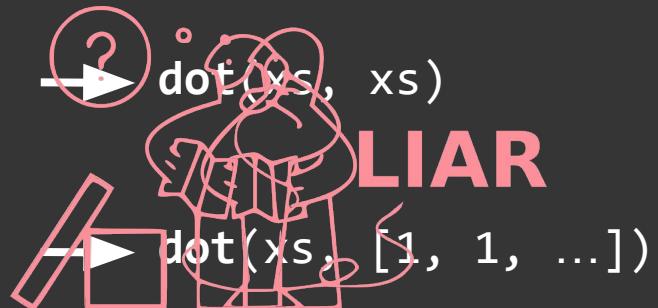
Vector norm

```
reduce(+, map(λx. x * x, xs))
```



Array sum

```
reduce(+, xs)
```



# Three ingredients

## Minimalist IR

map  
reduce  
zip  
join  
split  
concat  
...

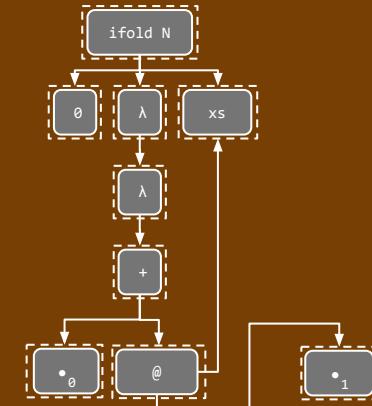
build  
ifold

## Language rules + library idioms

`build(N, f)@i = f(i)`  
(8 rules in total)

`fill(c) = build(N, λi. c)`  
(10 for BLAS, 11 for PyTorch)

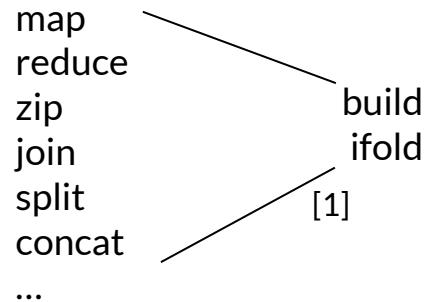
## Equality saturation



# Ingredient #1: Minimalist IR

```
build(3, f) = [f(0), f(1), f(2)]
```

```
ifold(3, f, init) = f(f(f(init, 0), 1), 2)
```



# Idiom recognition on minimalist IR



```
dot(xs, ys) = reduce(+, map(*, zip(xs, ys@i))) ys@i)
```



Dot product

```
reduce(+, map(*, zip(xs, ys@i))) → dot(xs, ys)  
reduce(+, map(*, zip(xs, ys)))
```



Vector norm

```
reduce(+, map(λx. x * x, xs)) → dot(xs, xs)  
reduce(+, map(λx. x * x, xs))
```



Array sum

```
reduce(+, 0, λas)λi. a + xs@i)  
reduce(+, 0, xs)
```



# Ingredient #2: Rules and idioms

Language rules

```
x = x * 1
```

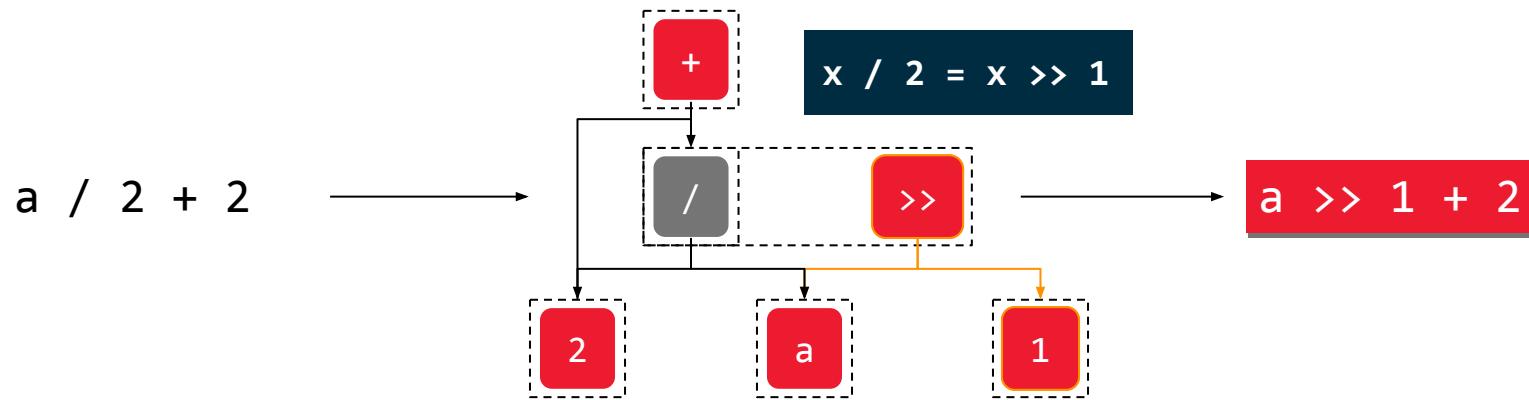
```
build(N, f)@i = f(i)
```

Library idioms (e.g., PyTorch, BLAS)

```
dot(A, B) =  
ifold(N, λa. λi. A@i * B@i + a)
```

```
fill(c) = build(N, λi. c)
```

# Ingredient #3: Equality saturation



# Combining the three ingredients

Expressions in e-graph:

$\text{ifold}(N, 0, \lambda a. \lambda i. xs @ i + a)$

$\text{ifold}(N, 0, \lambda a. \lambda i. xs @ i * 1 + a)$

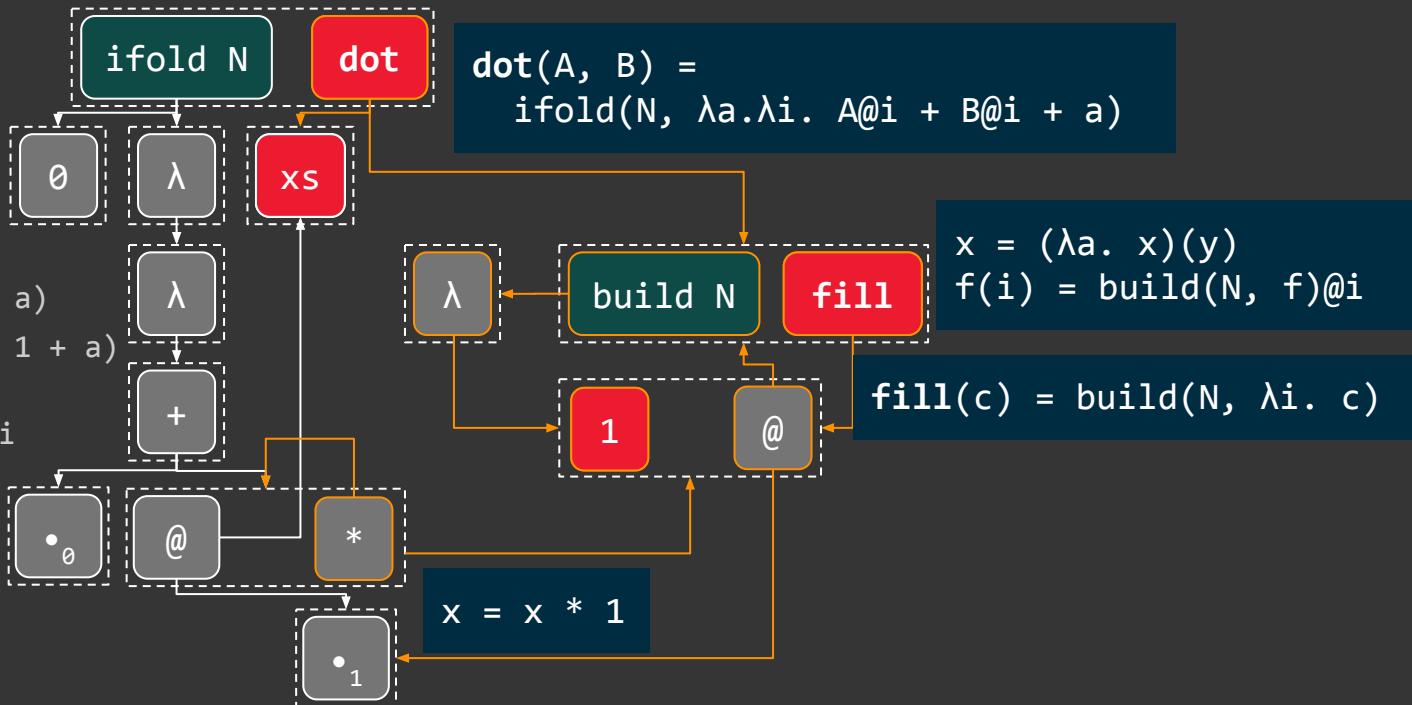
$\text{ifold}(N, 0, \lambda a. \lambda i.$

$xs @ i * \text{build}(N, \lambda j. 1) @ i$   
+ a)

$\text{ifold}(N, 0, \lambda a. \lambda i.$

$xs @ i * \text{fill}(1) @ i + a)$

**dot(xs, fill(1))**



# Evaluation methodology

PolyBench C

```
sum = 0.0;
for (i = 0; i < N; i++) {
    sum += xs[i] * ys[i];
}
```



Minimalist IR

```
ifold N 0 (λa. λi. a + xs@i)
```



LIAR



IR with  
library calls

**BLAS**  
`ddot(xs, [1, 1, ...])`

**PyTorch**  
`sum(xs)`



C code

**BLAS**  
`double ones[N] = { 1, 1, ... };`  
`ddot(N, xs, 1, ones, 1);`

# PolyBench: doitgen example

```
for (int r = 0; r < R; r++) {
    for (int q = 0; q < Q; q++) {
        for (int p = 0; p < P; p++) {
            sum[p] = 0.0;
            for (int s = 0; s < P; s++) {
                sum[p] += A[r][q][s] * B[s][p];
            }
        }
        for (int p = 0; p < P; p++) {
            A[r][q][p] = sum[p];
        }
    }
}
```

# PolyBench: doitgen example

```
build(N, λr. build(N, λq. build(N, λp.  
    ifold(N, 0, λa. λs. A@r@q@s * B@s@p + a))))
```

## Targeting PyTorch

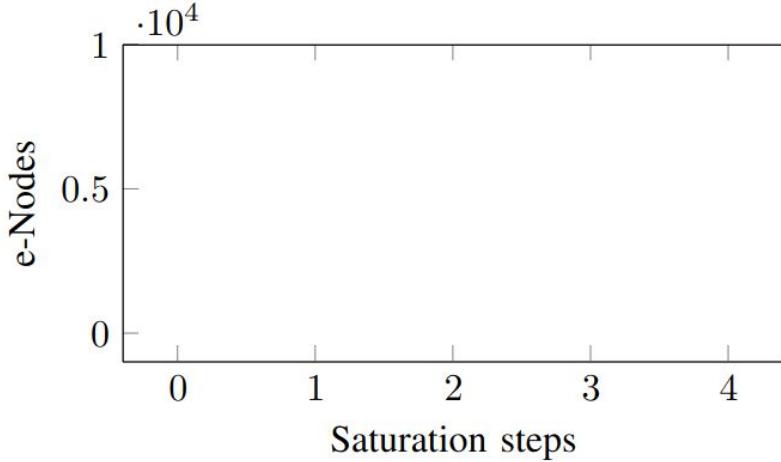
```
build(N, λr. mm(A@r, transpose(B)))
```

## Targeting BLAS

```
build(N, λr. dgemmF,T(1, A@r, B, 0, build(N, λq. memset(0))))
```

# Solutions over time

doitgen benchmark

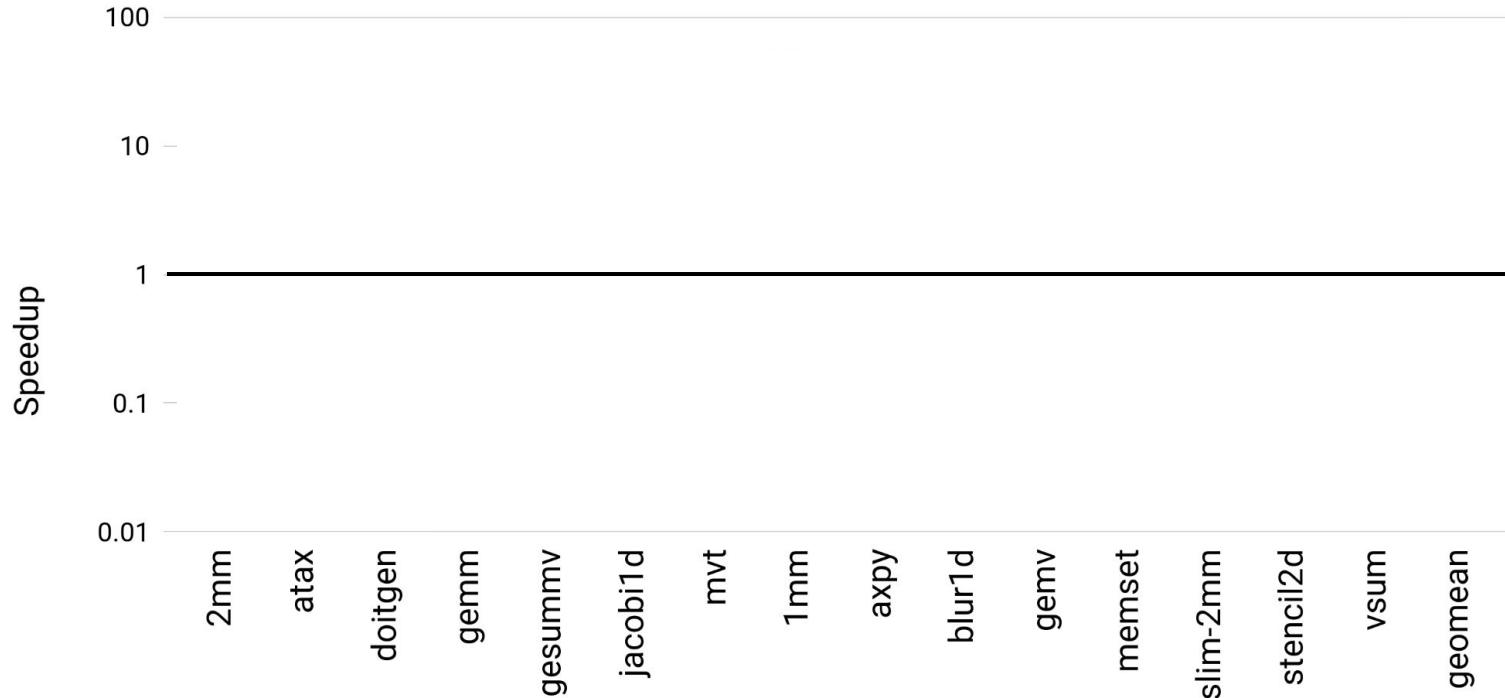


**PyTorch**

**BLAS**

# Run time speedup

Compared to PolyBench C kernels, BLAS parallel CPU library





Jonathan's  
home page

## Conclusion

- **Minimalist IR** simplifies semantics
- **Language rules + library idioms** scheme adapts to different APIs
- **Equality saturation** transforms programs to find idioms
- **More details in paper**
  - De Bruijn indices, full rules + idioms, etc.
- **Future work:** Improve scalability of equality saturation

Latent Idiom Recognition for a Minimalist Functional Array Language using Equality Saturation  
Christophe Dubach  
McGill University & Mila  
Montreal, Quebec, Canada  
christophe.dubach@cs.mcgill.ca

Jonathan Van der Cruyssen  
*McGill University*  
Montreal, Quebec, Canada  
jathan.vandercruyssen@mail.mcgill.ca

pling programs is typically done by writing high-performance libraries or hardware drivers, recognizing such idioms is difficult to work out. Recognition memory may be difficult to work out, especially if the user has little knowledge. In addition, slight variations in the way a program is used may hide the need for a minimalist function. For example, a user may be accustomed to using a small, but expensive, set of operations on a tiny set of variables. The user may also use language constructs of the same variety, which are not well suited for the type of computation he is trying to perform. This removes the need for a minimalist function, as the user can simply change the language to define the idioms.

## 1 INTRODUCTION

**L. INTRODUCTION**

Generating high-performance code for today's heterogeneous specialized hardware is challenging. A promising approach is to automatically rewrite specific idioms found in programs as highly optimized memory [3, 5, 13], which decouples compiler's pattern recognition from the hardware-specific knowledge embedded in the library implementation.

However, library is only useful when its idioms are found. Most prior work [3, 5, 13] focuses on idioms in the compiler at a low level of abstraction. Crafting low-level idioms is tedious, requiring expert computer knowledge fail to respond to analysis results. Furthermore, recognition might fail to respond to analysis results. For example, program *S*, which idioms recognition fails to find, contains a high-level recognition rule:

To solve the first challenge, this paper proposes to both programs and idioms using a high-level functional array language. Functional array programming for high-performance computing has become increasingly popular in recent years [6, 16, 19] and the concise, high-level nature of functional languages facilitates pattern detection and rewriting [7]. The second challenge of robust pattern matching considers the second challenge of robust pattern matching. Consider the following two examples:

In a functional language, the second pattern detection amounts to finding hidden idioms. Consider the following vector sum program:  $\text{sum}(\mathbf{v}) = \text{fold } (+) \ 0 \ \mathbf{v}$ .

If we have at our disposal a library function that performs such a sum, then we present the problem as a call to that function. However, if the library supports more general primitives, the rewriting problem becomes more complicated.

Suppose the library has a function *fill* that creates an array of *n* identical elements. A human could both implement the vector  $\text{vector}[n] = \text{det}(\text{v}, \text{fill}(1))$ . A pattern generator would have to use its intuition and would be hard-pressed to find this solution as neither the *sum* corresponds to *dot*, nor that to *fill* appear in the original program.

This paper proposes to find an alternative solution, using Latent Action Rewriting (LAR), a trustworthy technique that finds all explored latent forms using equality saturation. Equalized saturation [20] discovers all possible program variants under a fixed limit. LAR relies on a minimalist Intermediate Representation (IR) based on a few simple functional programming primitives, resulting in a compact set of rewrite rules. This makes it easier to capture the essential structure of a program without getting bogged down in language-specific details. Using equalized saturation to analyze LAR and library-specific identities, LAR efficiently programs that hide hidden subjects and improve program efficiency.

Section 3 shows how this technique can be applied on a concrete example of a library SubPrograms (BLAS) [1] that contains subroutines for matrix multiplication, effectiveness,

Basic Linear Algebra [13] libraries. To evaluate LIAR's efficiency, we apply it to custom kernels, and to linear algebra simulation kernels from the PolyBench suite. We show that LIAR leads to significant improvements in efficiency. Overall, this work demonstrates that LIAR is a powerful tool for idiom recognition. LIAR can be easily adapted to different libraries by appropriate idiom descriptions for those libraries. In summary, this paper makes the following contributions:

- appropriate domain. To summarize, this paper makes the following contributions: a minimalist IR and its tiny subset of rewrites; a new language for capturing the language semantics; how this minimalist IR can be used to express programs found in BLAS and PyTorch; demonstrates the effectiveness of using equality saturation with a minimalist IR on a set of computational kernels.

$\beta_2$  N. From

ally looks based on illustrate function to r division. at selects

Examiner (jV-C) 

Applies a set of rules  $\mathcal{R}_1$  to e-graph  $\mathcal{E}_1$  from e-graph  $\mathcal{E}_0$ , of which are highlighted in blue.

organization  
participate  
are  
soil  
the  
occurred  
and  
application  
was

... zero or more  
is an anonymous  
problem domain.  
S but not s...

*ijin index.*

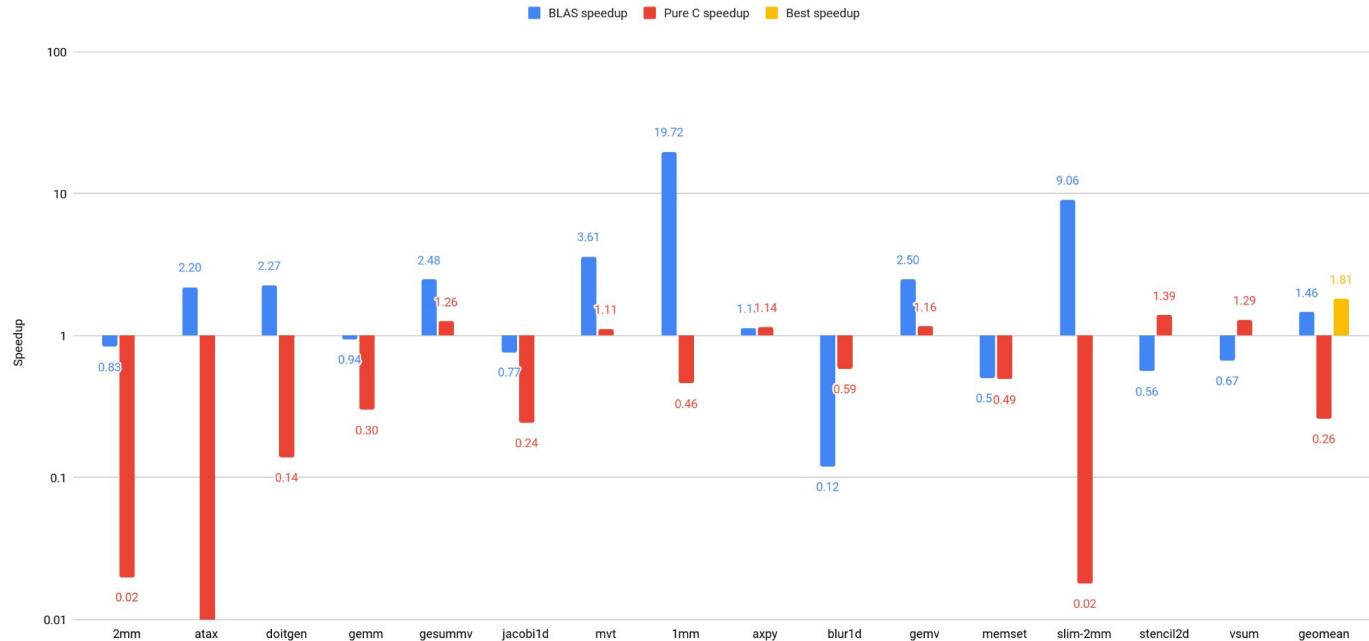
*with*

1

1

# Run time speedup

## Compared to reference kernels

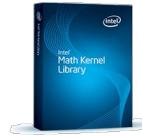


# BLAS and PyTorch

BLAS



cuBLAS



- Library specification
- Linear algebra
- Complex API

PyTorch



- Library
- Machine learning
- Streamlined API

```
dgemm(  
    transa, transb,  
    m, n, k,  
    alpha, a, lda,  
    b, ldb,  
    beta, c, ldc)
```

```
mm(a, b)  
add(a, b)  
mul(a, b)  
transpose(a)
```

# LIAR in action



```
dot(xs, ys) = ifold N 0 (λa. λi. a + xsi * ysi)
```



**Dot product**

```
ifold N 0 (λa. λi. a + xsi * ysi) → dot(xs, ys)  
reduce (+) (map (*) (zip xs ys))
```



**Vector norm**

```
ifold N 0 (λa. λi. a + xsi * xsi) → dot(xs, xs)  
reduce (+) (map (λx. x * x) xs)
```



**Array sum**

```
ifold N 0 (λa. λi. a + xsi) → dot(xs, fill(1))  
reduce (+) 0 xs
```

# Equality saturation pipeline

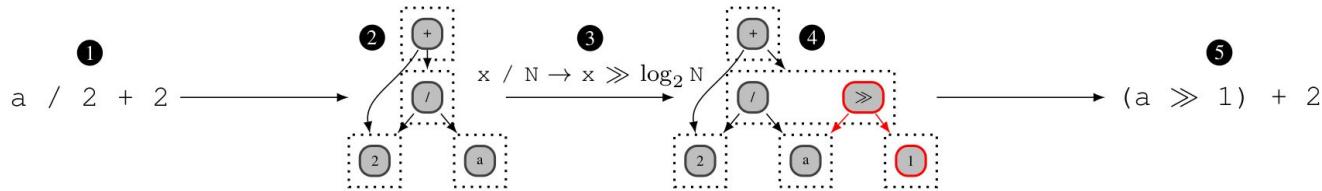


Fig. 1: Expression ① is converted to e-graph ②, which is subsequently saturated. In this example, only rule ③ is applied:  $x / N \rightarrow x \gg \log_2 N$ . From saturated e-graph ④, expression ⑤ is selected by an extractor that prefers bitwise shift.

# LIAR pipeline overview

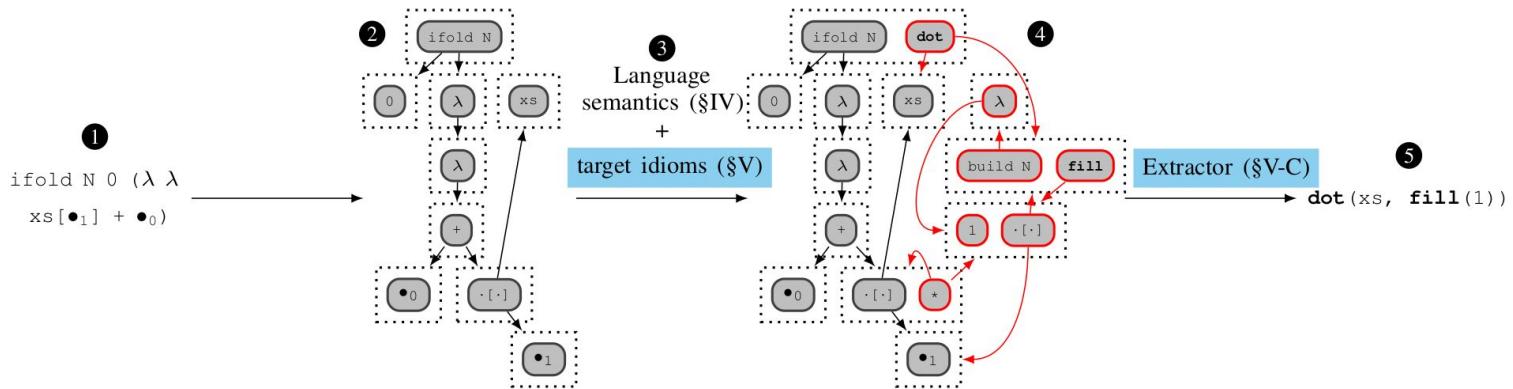


Fig. 2: Overview of LIAR, the proposed technique. Vector sum expression ① is converted to e-graph ②. Equality saturation applies a set of rules ③ to e-graph ②. These rules consist of target-independent language semantics and target-specific idioms. Rule application yields updated e-graph ④. From e-graph ④, a target-specific extractor chooses expression ⑤. The only target-specific components are the target idioms and extractor, both of which are highlighted in blue.

# IR grammar

$e ::= \lambda e$	<i>lambda abstraction</i>
$e e$	<i>lambda application</i>
$\bullet_i$	<i>parameter use</i>
$\text{build } N f$	<i>array construction</i>
$e[e]$	<i>array indexing</i>
$\text{ifold } N e e$	<i>iteration with accumulator</i>
$\text{tuple } e e$	<i>tuple creation</i>
$\text{fst } e   \text{snd } e$	<i>tuple unpacking</i>
$\mathbf{f}(\bar{e})$	<i>named function application</i>

Fig. 3: The grammar describing the minimalist IR.  $\bar{e}$  indicates zero or more instances of  $e$ .  $N$  is a compile-time integer constant and  $\mathbf{f}$  is an anonymous function. The set of available named functions depends on the problem domain. For example, **gemm** is a named function when targeting BLAS but not when targeting PyTorch.

# IR semantics

```
(λ e) y = subst(e, y)          (E-BETAREDUCE)
(build f N) [i] = f i           (E-INDEXBUILD)
fst (tuple a b) = a            (E-FSTTUPLE)
snd (tuple a b) = b            (E-SNDTUPLE)
ifold 0 init f = init          (E-FOLDINIT)
ifold (N + 1) init f = f N (ifold N init f) (E-FOLDSTEP)
```

Listing 1: Reduction semantics for the minimalist IR.

```
(λ e) y → subst(e, y)          (R-BETAREDUCE)
e → (λ e ↑) y                  (R-INTROLAMBDA)
(build f N) [i] → f i           (R-ELIMINDEXBUILD)
f i → (build f N) [i]           (R-INTROINDEXBUILD)
fst (tuple a b) → a             (R-ELIMFSTTUPLE)
a → fst (tuple a b)             (R-INTROFSTTUPLE)
snd (tuple a b) → b             (R-ELIMSNDTUPLE)
b → snd (tuple a b)             (R-INTROSNDTUPLE)
```

Listing 2: Eight rewrite rules that capture the relationships between `build`, array access, tuple construction, and tuple deconstruction,  $\lambda$ -abstraction and  $\beta$ -reduction.

# Arithmetic identities

$$x + 0 = x \quad (\text{E-ADDZERO})$$

$$1 * x = x \quad (\text{E-MULONEL})$$

$$x * 1 = x \quad (\text{E-MULONER})$$

$$x * y = y * x \quad (\text{E-COMMUTEMUL})$$

Listing 3: Scalar arithmetic identities. Each identity corresponds to two rewrite rules: a left-to-right rule and a right-to-left rule.  $x$  and  $y$  are numbers.

# BLAS and PyTorch idioms

```

axpy(alpha, A, B)
= build N (λ alpha↑ * A↑[•0] + B↑[•0])          (I-AXPY)
dot(A, B)
= ifold N 0 (λ λ A↑↑[•1] * B↑↑[•1] + •0)      (I-DOT)
gemvF(alpha, A, B, beta, C)
= build N (λ alpha↑ * dot(A↑[•0], B↑) + beta↑ * C↑[•0])
                                         (I-GEMV)
gemmF,T(alpha, A, B, beta, C)
= build N (λ gemvN(alpha↑, B↑, A↑[•0], beta↑, C↑[•0])) (I-GEMM)
transpose(A)
= build N (λ build M (λ A↑↑[•0][•1]))      (I-TRANSPOSE)
gemvX(alpha, transpose(A), B, beta, c)
= gemv¬X(alpha, A, B, beta, c)      (I-TRANSPOSEINGEMV)
gemmX,Y(alpha, transpose(A), B, beta, C)
= gemm¬X,Y(alpha, A, B, beta, C) (I-TRANSPOSEAINGEMM)
gemmX,Y(alpha, A, transpose(B), beta, C)
= gemmX,¬Y(alpha, A, B, beta, C) (I-TRANSPOSEBINGEMM)
dot(build N (λ alpha * A[•0]), B)
= alpha * dot(A, B)                      (I-HOISTMULFROMDOT)
memset(0)
= build N (λ 0)                           (I-MEMSETZERO)

```

Listing 4: BLAS idioms considered in this work. F in **gemv**<sup>F</sup> is short for *false* and indicates that matrix A is not transposed. F, T in **gemm**<sup>F,T</sup> indicate that A is not transposed and B is transposed.

```

dot(A, B)
= ifold N 0 (λ λ A↑↑[•1] * B↑↑[•1] + •0)      (I-DOT)
sum(A)
= ifold N 0 (λ λ A↑↑[•1] + •0)                  (I-VECSUM)
mv(A, B)
= build N (λ dot(A↑[•1], B↑))                    (I-MATVEC)
mm(A, B)
= build N (λ mv(B↑, A↑[•1]))                   (I-MATMAT)
transpose(A)
= build N (λ build M (λ A↑↑[•0][•1]))    (I-TRANSPOSE)
transpose(transpose(A))
= A                                         (I-TRANSPOSETWICE)
add(A, B)
= build N (λ A↑[•0] + B↑[•0])                (I-ADDVEC)
add(A, B)
= build N (λ add(A↑[•0], B↑[•0]))       (I-LIFTADD)
mul(alpha, A)
= build N (λ alpha * A↑[•0])                 (I-MULSCALARANDVEC)
mul(alpha, A)
= build N (λ mul(alpha, A↑[•0]))       (I-LIFTMUL)
full(c)
= build N (λ c↑)                                (I-FULLVEC)

```

Listing 5: PyTorch idioms considered in this work. I-TRANSPOSETWICE captures a property of the **transpose** function; all other rules recognize idioms.

# Cost model

$\text{cost}(\text{build } N \ f)$	$= N \cdot (\text{cost}(f) + 1) + 1$
$\text{cost}(A[i])$	$= \text{cost}(A) + \text{cost}(i) + 1$
$\text{cost}(\text{ifold } N \ \text{init } f)$	$= \text{cost}(\text{init}) + N \cdot \text{cost}(f) + 1$
$\text{cost}(\text{tuple } a \ b)$	$= \text{cost}(a) + \text{cost}(b) + 1$
$\text{cost}(\text{fst } t)$	$= \text{cost}(t) + 1$
$\text{cost}(\text{snd } t)$	$= \text{cost}(t) + 1$
$\text{cost}(\lambda e)$	$= \text{cost}(e) + 1$
$\text{cost}(f \ e)$	$= \text{cost}(f) + \text{cost}(e) + 1$
$\text{cost}(\bullet_k)$	$= 1 (\forall k \in \mathbb{N})$
$\text{cost}(a + b)$	$= \text{cost}(a) + \text{cost}(b) + 1$
$\text{cost}(a * b)$	$= \text{cost}(a) + \text{cost}(b) + 1$
$\text{cost}(c)$	$= 1 (\forall c \in \mathbb{R})$

Listing 6: Definition of the base cost function.

$\text{cost}(\text{memset}(c))$	$= \text{cost}(c) + .8N + 1$
$\text{cost}(\text{dot}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .8N$
$\text{cost}(\text{axpy}(a, A, B))$	$= \text{cost}(a) + \dots + \text{cost}(B) + .8N$
$\text{cost}(\text{gemv}(a, A, B, b, C))$	$= \text{cost}(a) + \dots + \text{cost}(C) + .7NM$
$\text{cost}(\text{gemm}(a, A, B, b, C))$	$= \text{cost}(a) + \dots + \text{cost}(C) + .6NMK$
$\text{cost}(\text{transpose}(A))$	$= \text{cost}(A) + .9NM$

Listing 7: BLAS-specific additions to  $\text{cost}$ . Calls to external functions are discounted to make them more attractive. Discounting factors are chosen semi-arbitrarily.  $N$ ,  $M$ , and  $K$  are array dimensions.

$\text{cost}(\text{full}(c))$	$= \text{cost}(c) + .8N + 1$
$\text{cost}(\text{add}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .4N + .4M$
$\text{cost}(\text{mul}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .4N + .4M$
$\text{cost}(\text{sum}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .8N$
$\text{cost}(\text{dot}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .8N$
$\text{cost}(\text{mv}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .7NM$
$\text{cost}(\text{mm}(A, B))$	$= \text{cost}(A) + \text{cost}(B) + .6NMK$
$\text{cost}(\text{transpose}(A))$	$= \text{cost}(A) + .9NM$

Listing 8: PyTorch-specific additions to  $\text{cost}$ .  $N$  and  $M$  are array dimensions. For polymorphic arrays,  $N$  and  $M$  represent the product of the arrays' dimensions.

# Overview of kernels

Kernel	Suite	Description
2mm	PolyBench	Two generalized matrix multiplications
atax	PolyBench	Matrix transpose and vector multiplication
doitgen	PolyBench	Multiresolution analysis kernel (MADNESS)
gemm	PolyBench	Generalized matrix product
gemver	PolyBench	Vector multiplication and matrix addition
gesummv	PolyBench	Scalar, vector and matrix multiplication
jacob1d	PolyBench	1D Jacobi stencil computation
mvt	PolyBench	Matrix-vector product and transpose
1mm	Custom	One matrix multiplication
axpy	Custom	Vector scaling and addition
blur1d	Custom	1D stencil
gemv	Custom	Generalized matrix-vector product
memset	Custom	Zero vector creation
slim-2mm	Custom	Two matrix multiplications
stencil2d	Custom	2D stencil
vsum	Custom	Vector reduction with sum

TABLE I: Overview of kernels examined in this work. Poly-Bench kernel descriptions adapted from benchmark suite [13].

# Solutions found

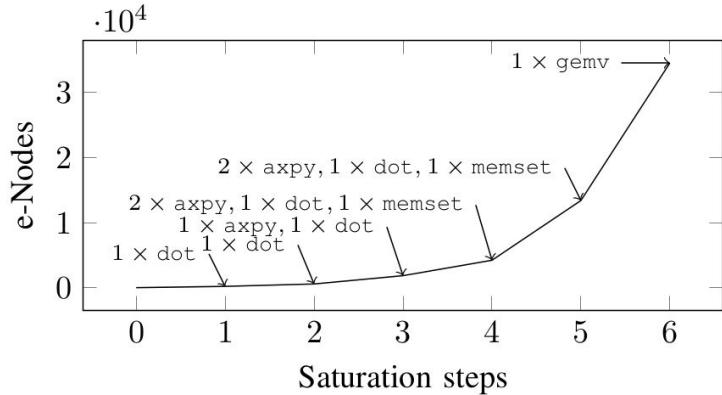
Kernel	Solution	Steps	e-Nodes
2mm	$3 \times \text{axpy} + 1 \times \text{dot}$ + $1 \times \text{gemv} + 3 \times \text{memset}$ + $1 \times \text{transpose}$	6	$3.44 \times 10^4$
atax	$2 \times \text{gemv} + 2 \times \text{memset}$	6	$1.49 \times 10^4$
doitgen	$1 \times \text{gemm} + 1 \times \text{memset}$	7	$2.14 \times 10^4$
gemm	$3 \times \text{axpy} + 1 \times \text{gemv}$ + $3 \times \text{memset}$	6	$1.93 \times 10^4$
gemver	$3 \times \text{axpy} + 2 \times \text{dot}$ + $1 \times \text{memset}$	5	$1.74 \times 10^4$
gesummv	$1 \times \text{dot} + 1 \times \text{gemv}$	6	$1.70 \times 10^4$
jacobi1d	$1 \times \text{gemv} + 1 \times \text{memset}$	5	$2.53 \times 10^4$
mvt	$2 \times \text{gemv} + 2 \times \text{memset}$	7	$2.69 \times 10^4$
lmm	$1 \times \text{gemm} + 1 \times \text{memset}$	7	$2.05 \times 10^4$
axpy	$1 \times \text{axpy}$	13	$2.57 \times 10^4$
blur1d	$1 \times \text{gemv} + 1 \times \text{memset}$	5	$1.71 \times 10^4$
gemv	$1 \times \text{gemv}$	7	$3.46 \times 10^4$
memset	$1 \times \text{memset}$	19	$2.86 \times 10^4$
slim-2mm	$1 \times \text{gemm} + 1 \times \text{gemv}$ + $2 \times \text{memset} + 2 \times \text{transpose}$	6	$2.04 \times 10^4$
stencil2d	$1 \times \text{gemv} + 1 \times \text{memset}$	5	$5.88 \times 10^4$
vsum	$1 \times \text{dot}$	10	$1.59 \times 10^4$

TABLE II: Solutions found for kernels when targeting BLAS. *Steps* describes the number of saturation steps applied to the kernel. *Solution* describes the external library calls found at the last step. *e-Nodes* counts the number of unique e-nodes in the e-graph, also at the last step.

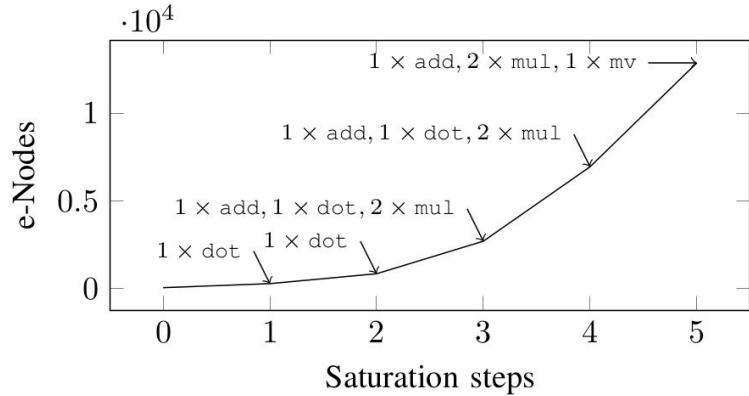
Kernel	Solution	Steps	e-Nodes
2mm	$1 \times \text{add} + 2 \times \text{mul} + 2 \times \text{mv}$ + $2 \times \text{transpose}$	5	$2.28 \times 10^4$
atax	$2 \times \text{mv} + 1 \times \text{transpose}$	7	$1.98 \times 10^4$
doitgen	$1 \times \text{mm} + 1 \times \text{transpose}$	6	$1.47 \times 10^4$
gemm	$1 \times \text{add} + 2 \times \text{mul} + 1 \times \text{mv}$ + $1 \times \text{transpose}$	5	$1.55 \times 10^4$
gemver	$2 \times \text{add} + 2 \times \text{dot} + 1 \times \text{mul}$	4	$9.06 \times 10^3$
gesummv	$1 \times \text{add} + 2 \times \text{mul} + 2 \times \text{mv}$	6	$1.54 \times 10^4$
jacobi1d	$1 \times \text{full} + 1 \times \text{mv}$	5	$3.13 \times 10^4$
mvt	$2 \times \text{mv} + 1 \times \text{transpose}$	7	$1.69 \times 10^4$
lmm	$1 \times \text{mm}$	6	$1.01 \times 10^4$
axpy	$1 \times \text{add} + 1 \times \text{mul}$	9	$1.52 \times 10^4$
blur1d	$1 \times \text{full} + 1 \times \text{mv}$	5	$2.13 \times 10^4$
gemv	$1 \times \text{add} + 2 \times \text{mul} + 1 \times \text{mv}$	6	$1.29 \times 10^4$
memset	$1 \times \text{full}$	15	$8.30 \times 10^3$
slim-2mm	$2 \times \text{mv} + 2 \times \text{transpose}$	5	$1.09 \times 10^4$
stencil2d	$1 \times \text{mv}$	4	$1.92 \times 10^4$
vsum	$1 \times \text{sum}$	9	$9.44 \times 10^3$

TABLE III: Solutions found for kernels when targeting PyTorch. Columns have the same meaning as in table II.

# Solutions over time



(a) Solutions over time for *gemv*, targeting BLAS.



(b) Solutions over time for *gemv*, targeting PyTorch.

Fig. 4: Solutions over time. The x-axes show equality saturation steps; the y-axes correspond to the number of e-nodes in the e-graph. Arrows indicate that a new best solution has been found, and are labeled with the external calls in that solution.

# Coverage over time

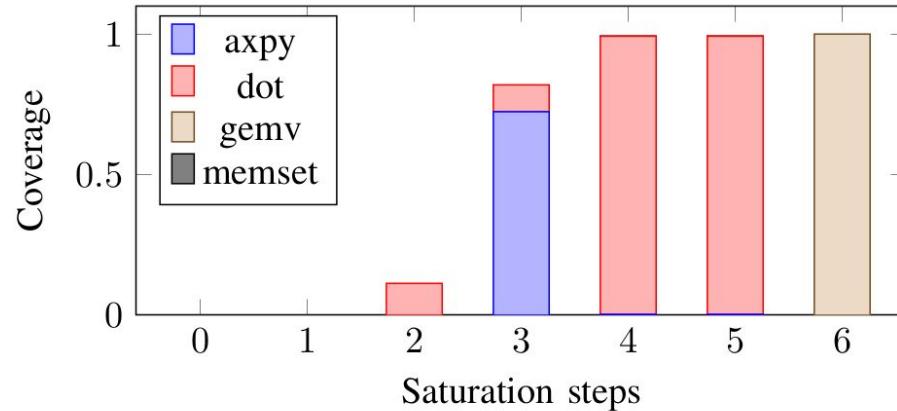


Fig. 5: Coverage over time for the *gemv* kernel, targeting BLAS. The x-axis shows saturation steps; the stacked bars depict the ratio of time spent in library functions. Higher bars indicate that more computation is offloaded to a library.

# Run time over time

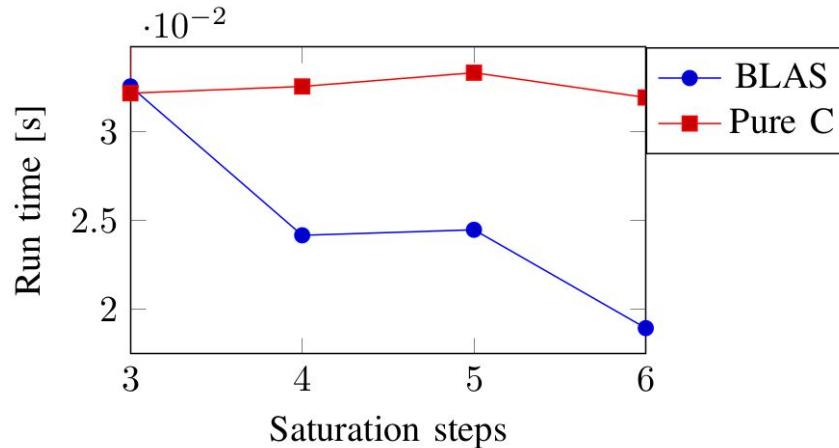


Fig. 6: Run time as a function of saturation steps. The x-axes show saturation steps; the y-axes depict the amount of time required to execute the solutions. Lower is better.

# Run time speedup

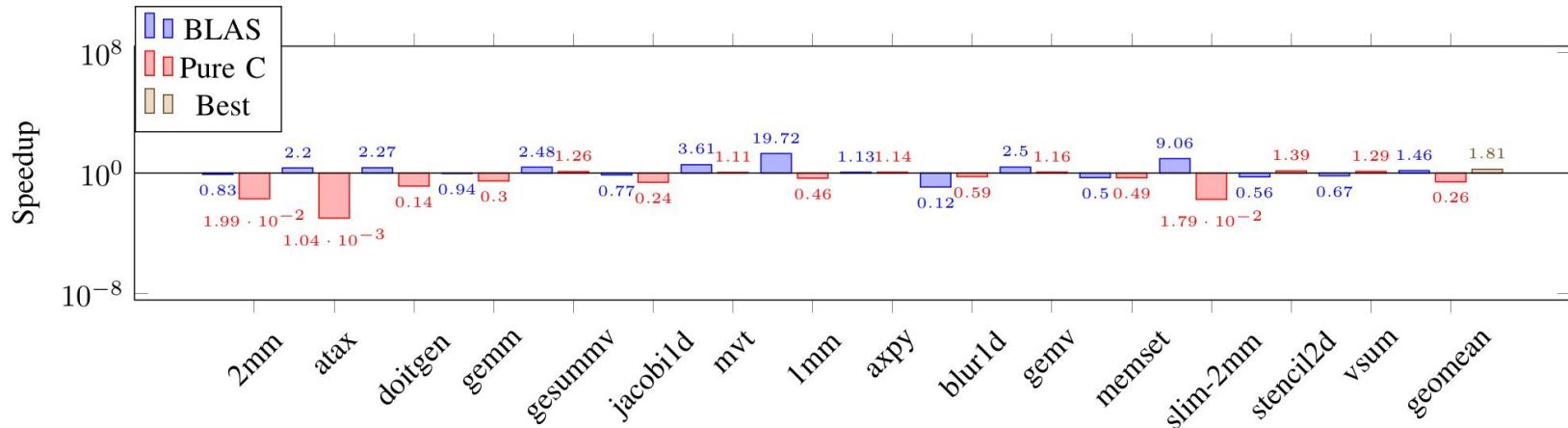


Fig. 7: Run time speedup of LIAR's solutions compared to reference implementations in C. For PolyBench kernels, the reference implementations are the original benchmarks. The reference for custom benchmarks is a hand-written C program coded in the style of a PolyBench kernel. Each bar represents the quotient of the reference solution run time and the LIAR solution run time. Higher is better.

# Functional array programming

**map**  $f$   $[x_1, x_2] = [f x_1, f x_2]$

**reduce**  $f$   $\text{init}$   $[x_1, x_2] = f (f \text{ init } x_1) x_2$

**zip**  $[x_1, x_2] [y_1, y_2] = [(x_1, y_1), (x_2, y_2)]$

Also: `join`, `split`, `concat`, `slide`, ...

# Rewrite rules: naive approach

$\text{map } g \ (\text{map } f \ \text{as}) \rightarrow \text{map } (g \circ f) \ \text{as}$

$\text{zip } (\text{map } f \ \text{as}) \ (\text{map } g \ \text{bs}) \rightarrow \text{map } (\lambda t. \ (f \ t_1, \ g \ t_2)) \ (\text{zip } \text{as} \ \text{bs})$

$\text{map } (\lambda a. \ f \ a \ a) \ \text{as} \rightarrow \text{map } (\lambda t. \ f \ t_1 \ t_2) \ (\text{zip } \text{as} \ \text{as})$

$\text{map } (\lambda a. \ f \ a \ c) \ \text{as} \rightarrow \text{map } (\lambda t. \ f \ t_1 \ t_2) \ (\text{zip } \text{as} \ [c, \ c, \ \dots])$

$\text{as} \rightarrow \text{map } (\lambda a. \ a) \ \text{as}$

**How many rules do we need?**

# De Bruijn indices

Named variables

$$\begin{aligned}\lambda x. & \quad x \\ \lambda x. & \quad \lambda y. \quad y \\ \lambda x. & \quad \lambda y. \quad x \\ \lambda x. & \quad \lambda y. \quad \lambda z. \quad z\end{aligned}$$

De Bruijn indices

$$\begin{aligned}\lambda & \quad \bullet_0 \\ \lambda & \quad \lambda \quad \bullet_0 \\ \lambda & \quad \lambda \quad \bullet_1 \\ \lambda & \quad \lambda \quad \lambda \quad \bullet_2\end{aligned}$$

# Latent idioms

Original computation

**2D convolution**

PyTorch

```
conv2d(  
    A,  
    weights  
)
```

BLAS

```
gemm(  
    im2col(...),  
    weights'  
)
```