

e-graphs, four ways

Max Willsey

2023

e-graphs, four ways

- e-graphs, as data structures
- e-graphs, as datalog
- ... but fast
- ... as datalog again

$$(a * 2) / 2$$



a

$(a * 2) / 2$

$=$

a

$$(a * 2) / 2 \rightarrow a$$

rewrite it!

useful

$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

not so useful

$$x * 2 = x \ll 1$$

$$x * y = y * x$$

$$x = x * 1$$

$$(a * 2) / 2 \rightarrow a * (2 / 2) \rightarrow a * 1 \rightarrow a$$

happy path

$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$



$(a * 2) / 2 \Rightarrow (a \ll 1) / 2$ ❌ wrong turn

$(a * 2) / 2 \Rightarrow (2 * a) / 2 \Rightarrow (a * 2) / 2$ loop

$a \Rightarrow a * 1 \Rightarrow a * 1 * 1 \Rightarrow \dots$ infinite size

pitfalls

$$x * 2 = x \ll 1$$

$$x * y = y * x$$

$$x = x * 1$$

but critical for
other inputs

$$(a * 2) / 2 \rightarrow a$$

which rewrite? when?

useful

$$(x * y) / z = x * (y / z)$$

$$x / x = 1$$

$$x * 1 = x$$

not so useful

$$x * 2 = x \ll 1$$

$$x * y = y * x$$

$$x = x * 1$$

$$x * 2 = x \ll 1$$

$$x = x * 1$$

$$x = 1 * x$$

$$x / x = 1$$

which rewrite? when?

all of them! all the time!

$$(z | h) * x = z | (h * x)$$

$$x * y = y * x$$

e-graphs
+
equality saturation

$x/x=1$

$x * 2 = x \ll 1$

$x = x * 1$

$x = 1 * x$

a

e!

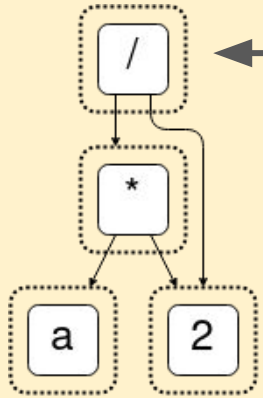
$(z | h) * x = z | (h * x)$

$x * y = y * x$

e-graphs, four ways

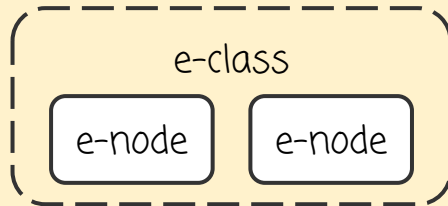
- e-graphs, as data structures
- e-graphs, as datalog
- ... but fast
- ... as datalog again

e-graphs?

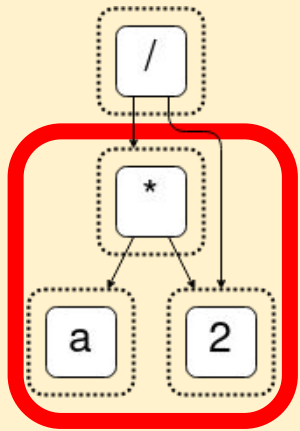


← this e-class represents

$$(a * 2) / 2$$

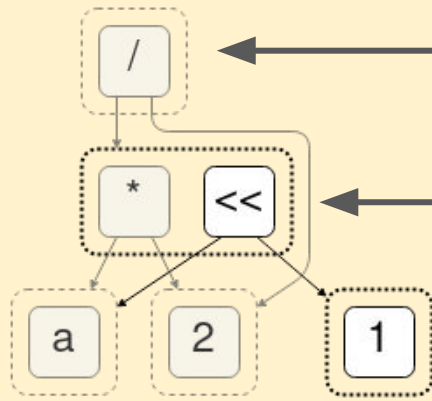
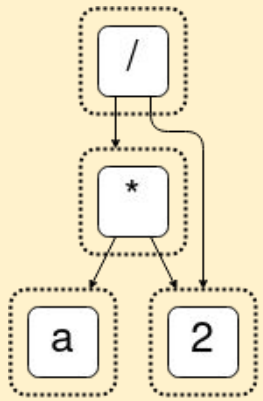


growing an e-graph



$$x * 2 \rightarrow x \ll 1$$

growing an e-graph

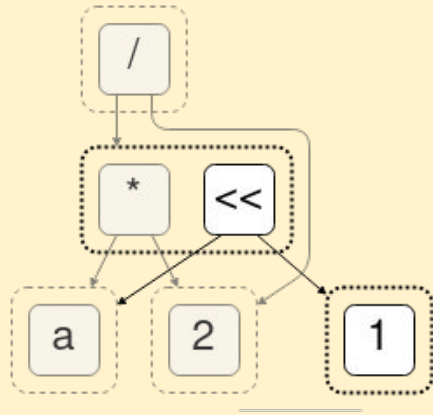
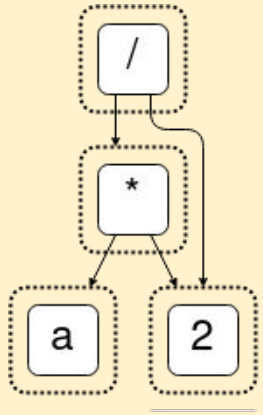


this e-class represents
 $(a * 2) / 2$ and $(a \ll 1) / 2$

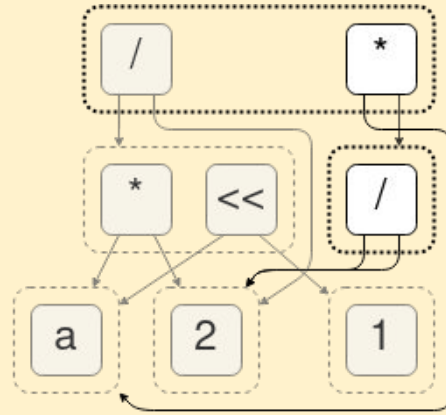
this e-class represents
 $(a * 2)$ and $(a \ll 1)$

$$x * 2 \rightarrow x \ll 1$$

growing an e-graph

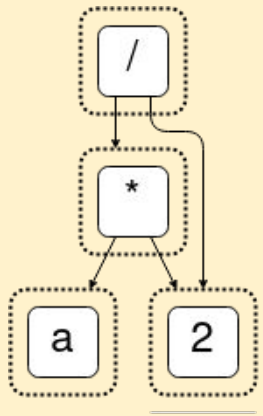


$$x * 2 \rightarrow x \ll 1$$

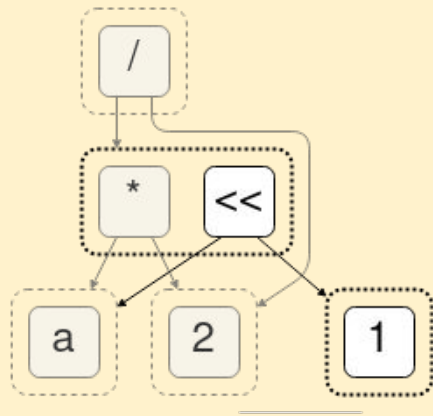


$$(x * y) / z \rightarrow x * (y / z)$$

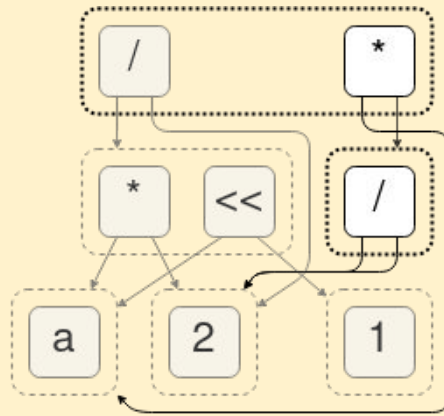
e-graphs are compact



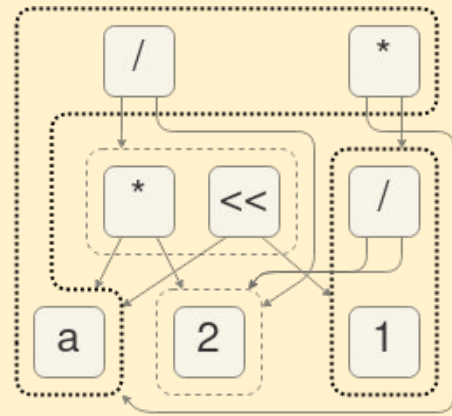
$$x * 2 \rightarrow x \ll 1$$



$$(x * y) / z \rightarrow x * (y / z)$$



$a, a * 1,$
 $a * 1 * 1, \dots$

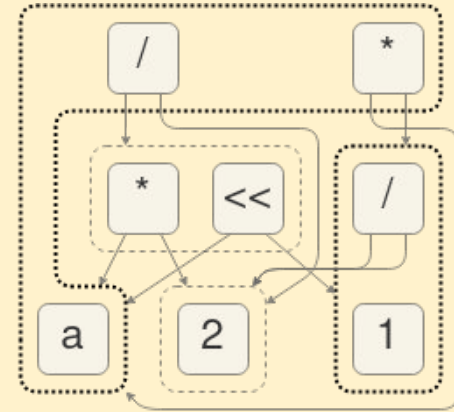


$$x / x \rightarrow 1$$

$$x * 1 \rightarrow x$$

saturation

- ✓ $x * 2 \rightarrow x \ll 1$
- ✓ $(x * y) / z \rightarrow x * (y / z)$
- ✓ $x / x \rightarrow 1$
- ✓ $x * 1 \rightarrow x$

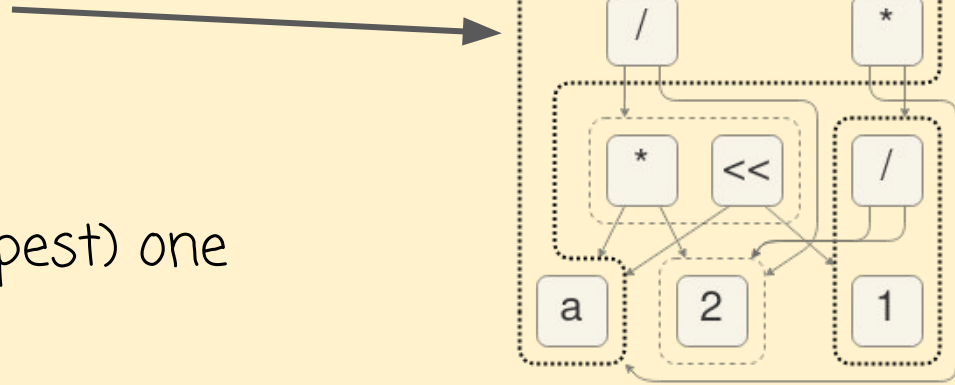


extraction

this e-class represents
 $(a * 2) / 2$, a , $a * 1$, ...

pick the smallest (cheapest) one

Knuth 76,
Generalization of Dijkstra's Algo.



extraction

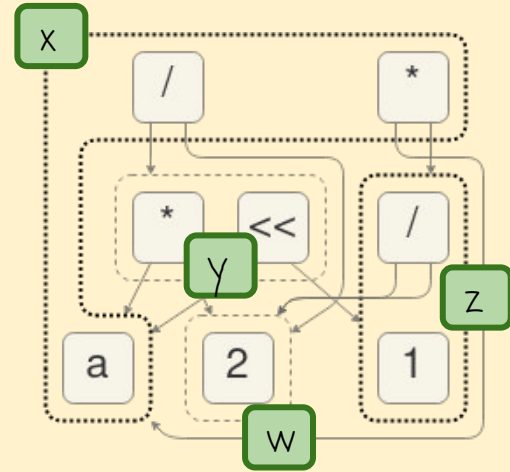
$$x = a + / (y, w) + * (x, z)$$

$$y = * (x, w) + \ll (x, z)$$

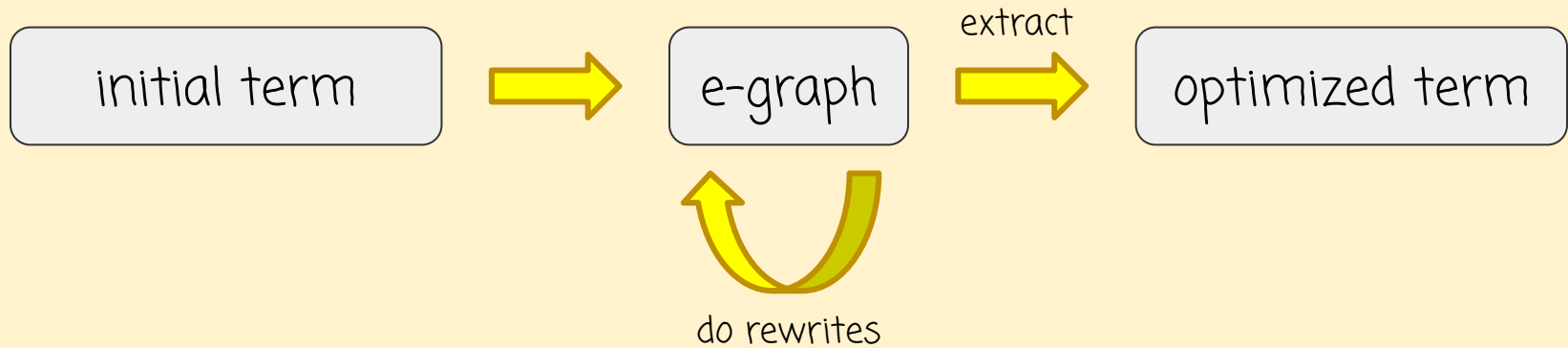
$$z = / (w, w) + 1$$

$$w = 2$$

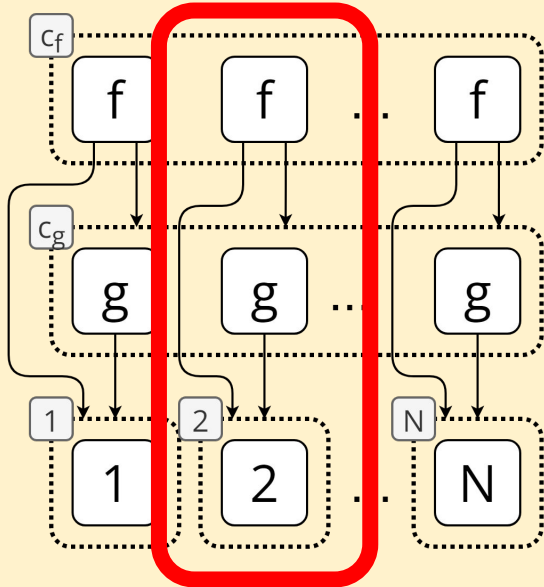
where $f(a, b) =$
min/+ semi-ring element
in terms of a, b



equality saturation



e-matching



f g e-nodes

e-classes

pattern

$f(a, g(a))$

subst

$\{a \mapsto 1\}$

$\{a \mapsto 2\}$

...

$\{a \mapsto N\}$

terms

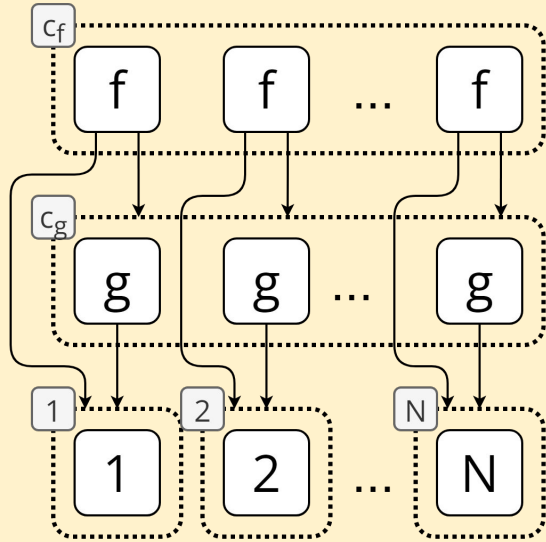
$f(1, g(1))$

$f(2, g(2))$

...

$f(N, g(N))$

e-matching



pattern
 $f(\alpha, g(\alpha))$

for e-class c in e-graph E :

```
for  $f$ -node  $n_1$  in  $c$ :
```

```
  subst = {root  $\mapsto$   $c$ ,  $\alpha \mapsto n_1.child_1$ }
```

```
  for  $g$ -node  $n_2$  in  $n_1.child_2$ :
```

```
    if subst[ $\alpha$ ] =  $n_2.child_1$ :
```

```
      yield subst
```

N^2 time, but only N matches

e-matching

- existing impls are backtracking based & complex
- doesn't help with equality constraints
- no data complexity results
 - NP-hard in pattern size... e-graph size??

more than rewriting

- there's more than syntactic rewriting
- sometimes, it's useful to consider semantics
 - $17 + 32 \rightarrow 49, \dots$
- constant folding, nullability, tensor shape, non-zero, interval arithmetic, etc, ...

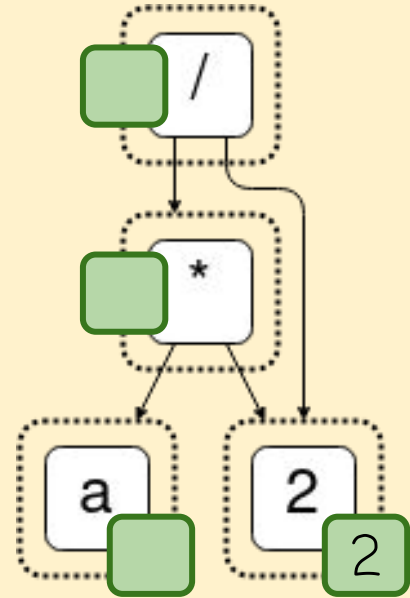
more than rewriting

analyses modulo equality

- uniform interface that works in many cases
- an understanding of analyses mean

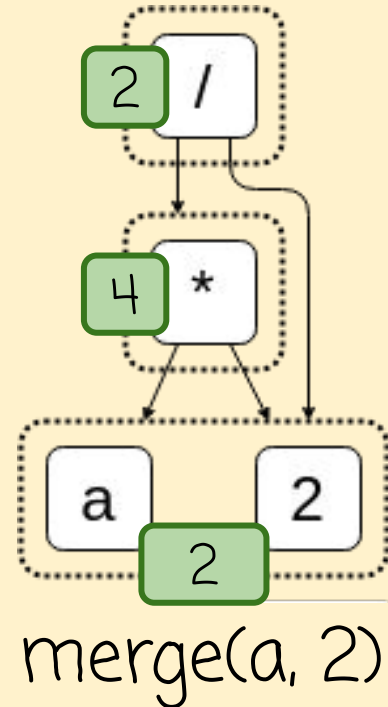
constant folding

- Option<Number> per eclass
- try to eval new e-nodes
- Option "or" on merge



constant folding

- Option<Number> per eclass
- try to eval new e-nodes
- Option "or" on merge
- it propagates up!

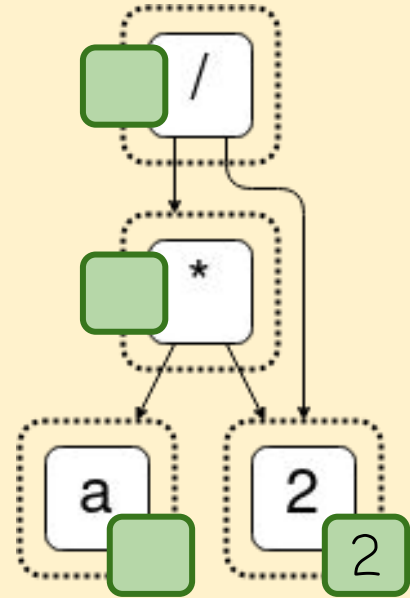


e-class analysis

- 1 fact per e-class from a join-semilattice D
- $\text{make}(n) \rightarrow d_c$
 - make a new analysis value for a new e-node
- $\text{join}(d_{c_1}, d_{c_2}) \rightarrow d_c$
 - combine two analysis values
- $\text{modify}(c) \rightarrow c'$
 - change the e-class (optionally)

constant folding

- $D = \text{Option}\langle\text{Number}\rangle$
- `make = eval`
- `join = option "or"`
- `modify = add the constant`



detour: intervals

x in $[0, 1]$
 y in $[1, 2]$

$x + y$ in $[1, 3]$

detour: intervals

x in $[0, 1]$
 y in $[1, 2]$

$$1 - 2y / (x + y) \quad \text{in} \quad [-3, 1/3]$$

$$= (x - y) / (x + y) \quad \text{in} \quad [-2, 0]$$

$$= 2x / (x + y) - 1 \quad \text{in} \quad [-1, 1]$$

intervals modulo equality

$$\begin{array}{l} x \text{ in } [0, 1] \\ y \text{ in } [1, 2] \end{array}$$

$$\begin{array}{l} 1 - 2y / (x + y) \text{ in } [-3, 1/3] \\ = (x - y) / (x + y) \text{ in } [-2, 0] \\ = 2x / (x + y) - 1 \text{ in } [-1, 1] \end{array} \left. \vphantom{\begin{array}{l} 1 - 2y / (x + y) \\ = (x - y) / (x + y) \\ = 2x / (x + y) - 1 \end{array}} \right\} [-1, 0]$$

e-class analysis uses

- lift program analyses to e-graphs
- conditional & dynamic rewrites
 - $x / x = 1$ iff $x \neq 0$
- can express other e-graph "hacks"
 - on-the-fly **extraction**

e-class analysis invariant

for each e-class

fixed point

$$\forall c \in G. \quad d_c = \bigvee_{n \in c} \text{make}(n) \quad \text{and} \quad \text{modify}(c) = c$$

Analysis data is LUB
(lattice properties)

egg: fast & easy e-graphs

- Rust library for generic e-graphs and eqsat
- packaged and documented: <https://docs.rs/egg>
- tutorials, industrial and academic users

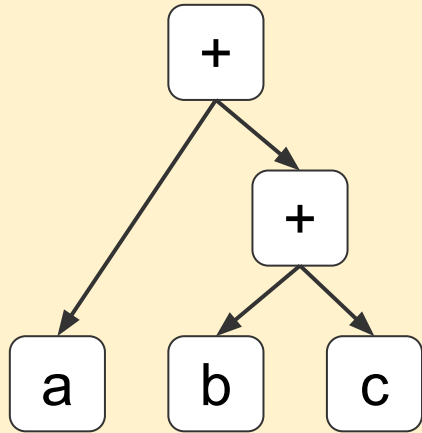
e-graphs, four ways

- e-graphs, as data structures
- e-graphs, as datalog
- ... but fast
- ... as datalog again

why?

- equality saturation is monotonic (-ish)
 - more equalities, more terms, no "destructive" rewrites
- e-matching is the bottleneck
 - it's a backtracking search for substitutions that satisfy a formula...

schema



"+" table

b	c	bc
a	bc	abc

rewrites as rules

- example: $x + (y + z) \rightarrow (x + y) + z$
- $+(x, y, \underline{xy}), +(xy, z, \underline{\text{root2}}) \leftarrow +(x, yz, \text{root}), +(y, z, yz)$
 - tempting to put root there
- not full existentials, just ADTs
 - existentials always "resolved" by FD, need a hashmap

what about the "e"?

- e-graph is an equivalence relation
 - congruence?
- pattern matching modulo equivalence
- equivalence is user-extensible!
 - think EGDs from chase

eq(x, y)

- just make an equivalence relation
 - symmetric, transitive, reflexive
- all joins modulo eq
 - $R(x, y), R(y, z)$ becomes
 $R(x, \underline{y_1}), eq(\underline{y_1}, \underline{y_2}), R(\underline{y_2}, z)$

rewrites with eq

- example: $x + (y + z) \rightarrow (x + y) + z$
- $+(x, y, \underline{xy}), +(xy, z, \underline{root2}), eq(xy1, xy2) \leftarrow$
 $+(x, yz1, root), +(y, z, yz2), eq(yz1, yz2)$
- non-linear patterns tend to be cyclic
 - consider $x + (y + x)$

congruence

- $eq(z1, z2) \leftarrow$
 $+ (x1, y1, z1), + (x2, y2, z2), eq(x1, y1), eq(x2, y2)$

doesn't work

- too slow
- various tricks don't fix it
 - specialized eqrel a la Souffle,
 - subsumption
 - see Yihong Zhang's thesis

lattices

- downside of e-class analyses: there's only one
- datalog has nice, cooperating "analyses"
 - mutually recursive rules
- requires recursive aggregation
 - LowerBound(expr, number)

e-graphs, four ways

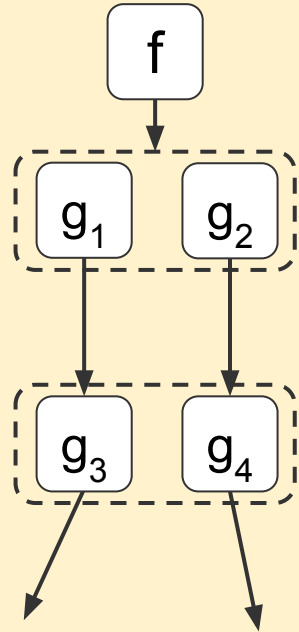
- e-graphs, as data structures
- e-graphs, as datalog
- ... but fast PLDI 23
- ... as datalog again

what was the problem?

- $eq(x, y)$ too slow
 - n^2 size, etc.
- no canonicalization!

if you don't canonicalize...

- e-matching only yields canonical entries!
- $f(g(g(x)))$
- $f(g_1(g_3(...)))$, $f(g_2(g_3(...)))$,
 $f(g_1(g_4(...)))$, $f(g_2(g_4(...)))$



canonicalize

- use a union-find to define $\text{leader}(x)$
 - $\text{leader}(x) = \text{leader}(y)$ iff $\text{eq}(x, y)$
- e-graph: explicit maintenance
 - if $x = y$, replace x, y , with $\text{leader}(x)$
 - collapse e-nodes $f(x), f(y)$ to $f(\text{leader}(x))$
 - massively shrinks the e-graph

canonicalize the db

- could use some form of subsumption
 - $f(\text{leader}(x)) :- f(x)$
 - $f(x) \leq f(y) :- x = \text{leader}(x), x \neq y$
 - way too slow/hacky to implement in, e.g. Souffle
- let's just do what e-graphs do
 - congruence closure, "rebuilding"

no more eq

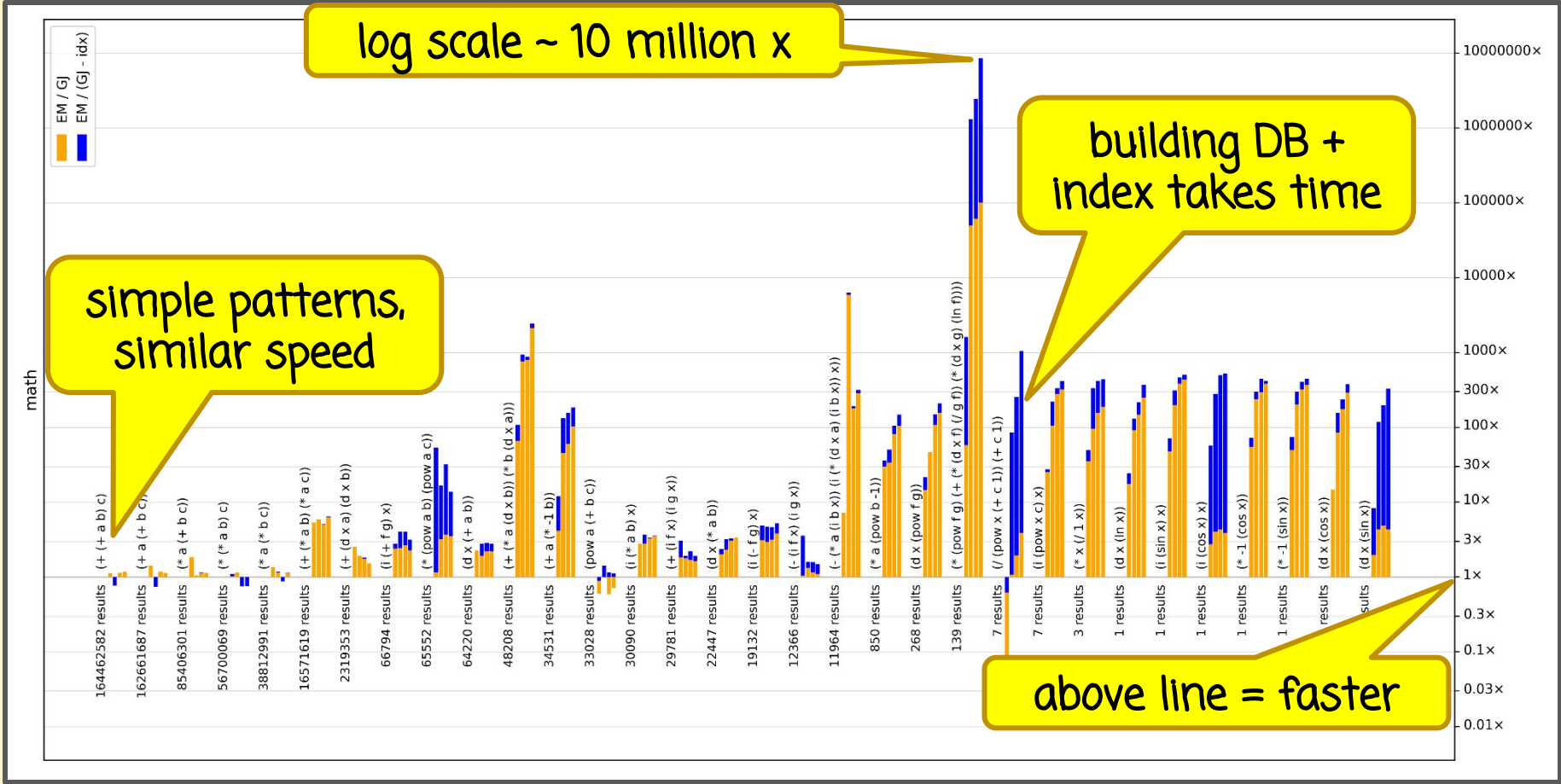
- eq relation/joins are gone!
- "semantic" equality becomes "syntactic" (again!)
- $R(x, y_1), R(y_2, z), \text{eq}(y_1, y_2)$ becomes
 $R(x, y), R(y, z)$

egglog

- datalog + functions + extensible equality
- examples
 - datalog: reachability
 - datalog + equality: reachability with node merging
 - eqsat: simple arithmetic optimization

egglog > eqsat

- simple implementation
 - separate optimization pass
- multipatterns
 - $a \times b = \text{split}(1, (a ++ c) \times b)$
 - $a \times c = \text{split}(2, (a ++ c) \times b)$
- incrementality via semi-naive



functions + equality

- $f(a) = b, f(c) = d$
- equality is extensible! user asserts $a = c$
 - what happens?
 - what happens in datalog when $f(x, y)$ and $f(x, z)$
 - lattices / semirings

merge expressions

- (function foo (i64) i64 :merge (min old new))
- example
 - (set (foo 7) 5)
 - (set (foo 7) 4)
 - (foo 7) = 4

merge expressions

- also works for conflicts coming from equality
- ex: interval arithmetic
 - (function hi (Expr) Rational :merge (min old new))
 - (function lo (Math) Rational :merge (max old new))

terms?

- (function mul (Expr Expr) Expr ...
 - :merge (union old new)
 - :default (mkset))
- congruence!
- labeled null? dynamic lattice?

the chase

- very similar to the Skolem chase
- EGDs capture FDs and "extensible" equality
- egglog has "stable" equality in a way
 - the UF prevents oscillation that can result in the chase

what's wrong?

- Built-in eq relation is really special
 - Requires infrastructure to do canonicalization
- semantic questions
 - why does this work? what's special about union-find?
- limits us to equality, and only one version of it

review

good

- fast queries
- fast congruence
- functions modulo eq
 - for e.g. e-class analysis

bad

- no terms!
 - only families of terms
- semantic questions
- only supports eq
 - rebuilding, extraction are special

e-graphs, four ways

- e-graphs, as data structures
- e-graphs, as datalog
- ... but fast
- ... as datalog again



WARNING
braindump ahead

criteria

- no "built-in" equality
 - what are the features needed to support this?
- can't lose terms
- fast e-matching, congruence, etc.
 - via canonicalization?
- other relations
 - partial orders, indexed equality, etc.

new encoding

- can't canonicalize the term tables
- capture the good part of the union-find
- "term tables" Add, Mul are immutable
 - a "mutable" union-find captures the eqrel
 - canonicalization rules create new terms

new encoding

- $(x + y) + z \rightarrow \text{Add}(x, y, xy), \text{Add}(xy, z, \text{root})$
- $\text{Add}(x, y, xy_1), \text{Add}(xy_2, z, \text{root}),$
- $\text{leader}(xy_1, xy_2),$
- $\text{leader}(x, x), \text{leader}(y, y), \text{leader}(z, z)$

egglog views

- views modulo "joining relations" like eq
- $\text{AddEq} = \{ \text{leader}(x), \text{leader}(y), \text{leader}(x) \mid (x, y, z) \text{ in Add} \}$
- some kind of notion of monotonicity
 - why does this work? what's the algebraic structure?

monotonicity?

- some kind of notion of monotonicity
 - why does this work? what's the algebraic structure?
- leader is also aggregation over eqrel

payoff

- other relations than eq
 - non-symmetric: reduction relations
 - other equalities
 - Indexed equality: $\text{eq}(\text{expr}, \text{expr}) \rightarrow \text{semiring}$
 - context as a set of assumptions

example

- $a + 3b = a + b \pmod{3}$
- $\text{eq}(x, y, \text{mod})$
- $\text{plus}(a, b, \underline{ab}), \text{eq}(ab, \text{root}, 3) \leftarrow$
 $\text{plus}(a, m1, \text{root}), \text{eq}(m1, m2, T), \text{mul}(3, b, m2)$
- $\text{eq}(x, y, f), \leftarrow \text{eq}(x, y, k), \text{factor}(k, f)$

payoff

- provenance
- semi-naive
 - can you re-derive congruence closure?
- top-down via demand transformation
- closer to "regular" datalog
 - other applications?

questions

- can any of this be implemented efficiently?
 - what algebraic structures are compatible with UF?
 - difficult IVM problem
 - queryable semirings?
- how to control the application of rules?
 - Demand transformation, explicit schedules, etc...
- a more flexible notion of monotonicity
- termination?