# Relational E–Matching
## Simpler, Faster, and Optimal

**Yihong Zhang**, University of Washington
PLDI 2021 Student Research Competition

# E-Graphs are everywhere!

Spores [VLDB '20]

Szalinski [PLDI '20]

Herbie [PLDI '15]

TenSat [MLSys '21]

Diospyoros [ASPLOS '21]
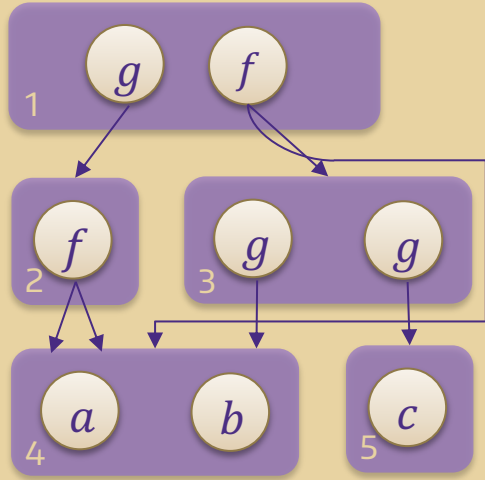
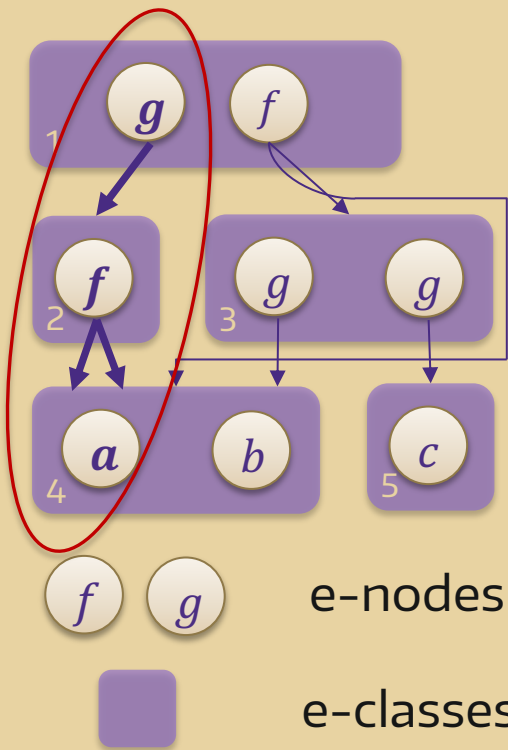Glenside [MAPS '21]

egg [POPL '20]

Z3

CVC4

Metatheory.jl

...

# E-Graphs
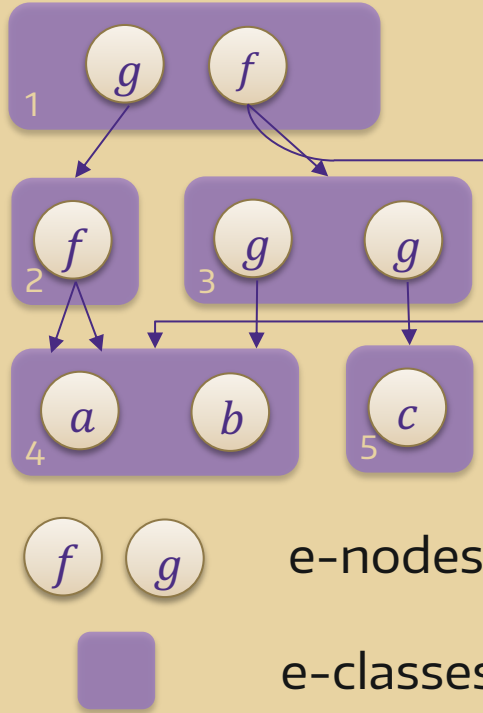


$$g\big(f(a,a)\big)$$

e-class 1
represents

f  g   e-nodes

e-classes

# E-Graphs



$$g\big(f(a,a)\big)$$

e-class 1
represents

e-nodes

e-classes

# E-Graphs



e-nodes

e-classes

e-class 1 represents

$g\big(f(a,a)\big)$    $f\big(a,g(a)\big)$
$g\big(f(a,b)\big)$    $f\big(a,g(b)\big)$
$g\big(f(b,a)\big)$    $f\big(b,g(a)\big)$
$g\big(f(b,b)\big)$    $f\big(b,g(b)\big)$
$f\big(a,g(c)\big)$
$f\big(b,g(c)\big)$

exponentially many terms!

# E-Matching

> E-matching is the query that finds patterns in an e-graph.



... responds to

$$f(a, g(a)), f(a, g(b))$$
$$f(b, g(a)), f(b, g(b)).$$

All witnessed by subst. $\{\alpha \mapsto 4\}$.

> Responsible for 60-90% of the run time.

# Existing E-Matching Algorithms

$f(\boldsymbol{\alpha}, g(\boldsymbol{\alpha}))$ $\Longrightarrow$

**for** all e-class $c$ **in** e-graph $E$:
👉 **for** all $f$ node $n_1$ **in** $c$:
    subst $= \{\boldsymbol{\alpha} \mapsto n_1.child_1, \textbf{root} \mapsto c\}$
👉 **for** all $g$ node $n_2$ **in** $n_1.child_2$:
    **if** subst$[\boldsymbol{\alpha}] = n_2.child_1$:
        **yield** subst

## Quadratic runtime!
### Yet at most $O(N)$ terms.

# Relational E–Matching

**E–Matching**

$\subseteq$

**E–Graphs**

<span style="color:darkred">Relational queries only involving joins e.g.,</span>
$$Q(a, c) \leftarrow R(a, b), S(b, c)$$

**Conjunctive Queries**

**Relational Databases**
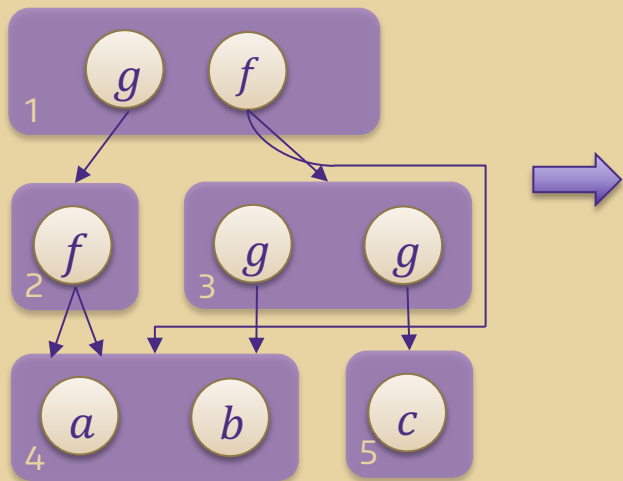
**Simpler**

# Comparison to Existing E-Matching

## Existing E-Matching

✗ Top-down backtracking search only.

✗ Exploits structural constraints only.

✗ No theoretical guarantee.

## Relational E-Matching

✓ Top-down, bottom-up, middle-out, etc. depending on the query optimizer.

✓ Exploits both structural constraints and equality constraints.

✓ Achieves optimality by adapting results from database research.

# E–Graphs → Relational Databases

Relations representing function symbols

$R_f$

| eclass-id | child$_1$ | child$_2$ |
|-----------|-----------|-----------|
| 1 | 4 | 3 |
| 2 | 4 | 4 |

# Faster E-Matching → Conjunctive Queries

$$f(\boldsymbol{\alpha}, g(\boldsymbol{\alpha}))$$

$$Q(\text{root}, \boldsymbol{\alpha}) \leftarrow$$
$$R_f(\text{root}, \boldsymbol{\alpha}, \boldsymbol{x}), R_g(\boldsymbol{x}, \boldsymbol{\alpha})$$

**for** all e-class $c$ **in** e-graph $E$:
  **for** all $f$ node $n_1$ **in** $c$:
    subst $= \{\boldsymbol{\alpha} \mapsto n_1.child_1, \textbf{root} \mapsto c\}$
    **for** all $g$ node $n_2$ **in** $n_1.child_2$:
      **if** subst$[\boldsymbol{\alpha}] = n_2.child_1$:
        **yield** subst

**Quadratic**

**for** all $r_f$ **in** $R_f$:
  **if** $\left(r_f.child_2, r_f.child_1\right)$ **in** $R_g$:
    **yield** $\{\textbf{root} \mapsto R_f.eclass\text{-}id,$
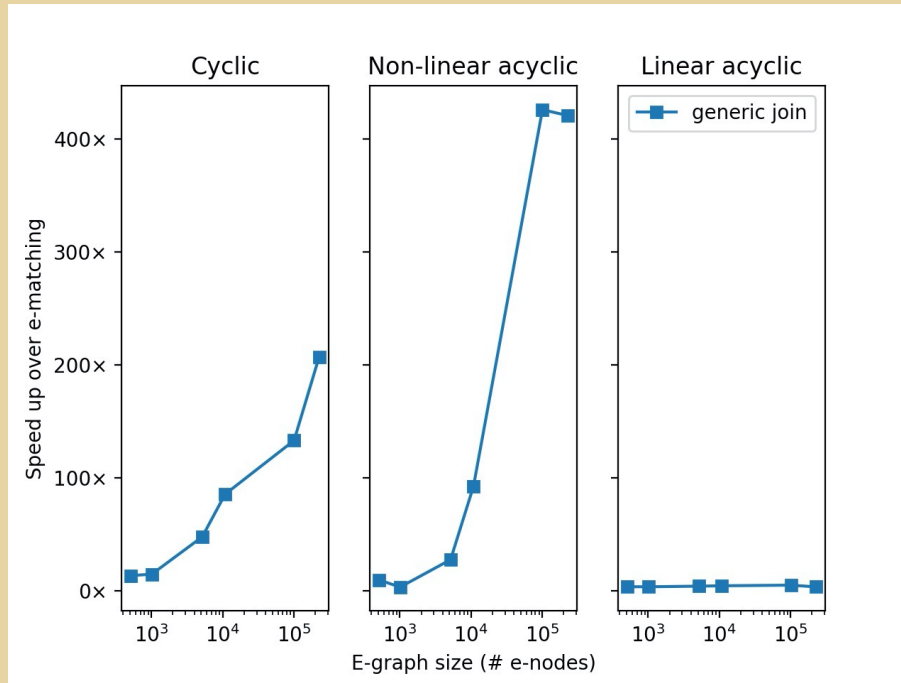        $\boldsymbol{\alpha} \mapsto R_f.child_1\}$

**Linear**

# Which CQ Algorithm to Use?

> Traditional two-way join algorithms like hash join and merge-sort join are provably suboptimal.

> We propose to use the generic join algorithm to solve the generated conjunctive query.

> Fix a pattern $p$, let $M(p, E)$ be the set of substitutions yielded by e-matching on an e-graph $E$ with size $n$, relational e-matching runs in time

$$\tilde{O}\left(\max_E |M(p, E)|\right).$$

**Optimal**

# (Preliminary) Evaluations



**Speed–ups over existing e-matching algorithms**
**(de Moura and Bjørner)**

- Asymptotic speed-ups on patterns with equality constraints (the cyclic and non-linear acyclic patterns).
- Similar performance on patterns with no equality constraints (the linear pattern).

# Thank you!

# Why Faster?

$f(\boldsymbol{\alpha}, g(\boldsymbol{\alpha}))$

enumerates all terms of the form $f(\boldsymbol{\alpha}, g(\boldsymbol{\beta}))$, and check if $\boldsymbol{\alpha} = \boldsymbol{\beta}$ only before yielding.

$Q(\mathrm{root}, \boldsymbol{\alpha}) \leftarrow$
$R_f(\mathrm{root}, \boldsymbol{\alpha}, \boldsymbol{x}), R_g(\boldsymbol{x}, \boldsymbol{\alpha})$

builds indices on both $\boldsymbol{\alpha}$ and $\boldsymbol{x}$ ; only enumerates terms where constraints on $\boldsymbol{\alpha}$ and $\boldsymbol{x}$ are both satisfied.

Equality constraints are exploited to prune the search space!

# Which CQ Algorithm to Use?

> Traditional two-way join algorithms like hash join and merge-sort join are not optimal.

> For example, for conjunctive query whose corresponding hypergraph is cyclic, two-way join algorithms will enumerate asymptotically more atoms than needed.

  – $(\alpha + \beta) \times (\alpha + \gamma)$.



Hypergraph of $(\alpha + \beta) \times (\alpha + \gamma)$

# Generic Join Algorithm

> Multi-way join algorithm that avoids enumerating unnecessarily large intermediate relations.

> Worst-case optimal & efficient for both cyclic and acyclic queries.