

# Color Refinement as an Index for Evaluating Acyclic Conjunctive Queries

---

Cristian Riveros   [Benjamin Scheidt](#)   Nicole Schweikardt

AlMoTh'24

Humboldt-Universität zu Berlin



# What the next 20 minutes will be about

The Problem

Query Evaluation

Color Refinement

Main Result: Using the Computed Coloring as an Index

Outlook



Lots of characters

# Who could expect to become a godparent?

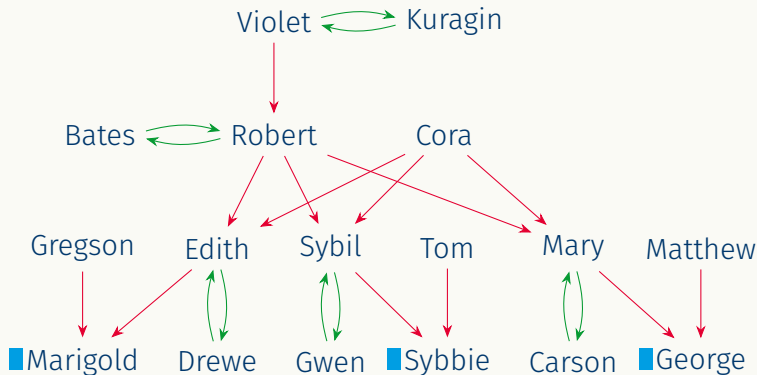
Parent	of	Parent	of	Fond	of
Violet	Robert	:	:	Mary	Carson
Robert	Mary	Mary	George	Carson	Mary
Robert	Edith	Matthew	George	Kuragin	Violet
Robert	Sybil	Sybil	Sybbie	Violet	Kuragin
Cora	Mary	Tom	Sybbie	Robert	Bates
Cora	Edith	Edith	Marigold	Bates	Robert
Cora	Sybil	Gregson	Marigold	Edith	Drewe
				Drewe	Edith
				Sybil	Gwen
				Gwen	Sybil

---

<b>Baby</b>	George	Marigold	Sybbie
-------------	--------	----------	--------

---

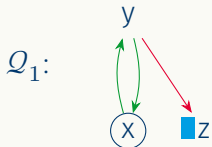
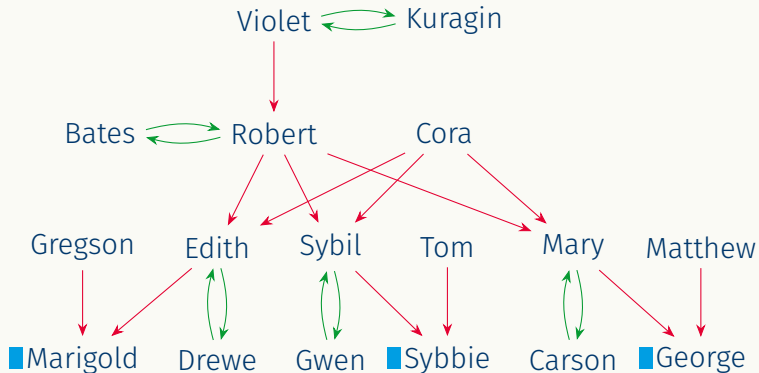
# The Database as a Graph



Who could expect to become a godparent?

$$Q_1 := \text{Ans}(x) \leftarrow F(x, y), F(y, x), P(y, z), B(z)$$

# Who could expect to become a godparent?



$$[[\mathcal{Q}_1]](D) = \{\text{Drewe, Gwen, Carson}\}$$

A quick (and simplified) recap on

★ Acyclic Conjunctive Queries ★

# Evaluating Acyclic Conjunctive Queries

## Conjunctive Queries (CQ, ACQ)

A CQ  $Q$  with  $\text{free}(Q) = \bar{x}$  looks like this:

$$\text{Ans}(\bar{x}) \leftarrow \alpha_1(\bar{x}_1), \alpha_2(\bar{x}_2), \dots, \alpha_n(\bar{x}_n)$$

$Q$  is acyclic, if it has a join tree.

In particular, a **binary**  $Q$  is acyclic, if it “looks like a tree”.

## Theorem (Yannakakis 1981)

*There exists an algorithm that computes  $\llbracket Q \rrbracket(D)$  in time  $\mathcal{O}(|D| \cdot |\llbracket Q \rrbracket(D)|)$  for every database  $D$  and every ACQ  $Q$ .*

**Data Complexity** — Assume  $Q$  to be constant.



# Constant Delay Enumeration: Refined runtime

## Yannakakis

$\text{enum}(Q_1, D)$   
... ⚙ ...  $\mathcal{O}(??)$   
out 1: Drewe  
... ⚙ ...  $\mathcal{O}(??)$   
out 2: Gwen  
... ⚙ ...  $\mathcal{O}(??)$   
out 3: Carson  
... ⚙ ...  $\mathcal{O}(??)$   
out: EOE

$\left. \begin{array}{l} \mathcal{O}(??) \\ \mathcal{O}(??) \\ \mathcal{O}(??) \\ \mathcal{O}(??) \end{array} \right\} \mathcal{O}(|D| \cdot \|Q_1\|(D))$

## Bagan, Durand, Grandjean

$\text{enum}(Q_1, D)$   
... ⚙ ...  $\mathcal{O}(|D|)$   
eop, num: 3  
... ⚙ ...  $\mathcal{O}(1)$   
out 1: Drewe  
... ⚙ ...  $\mathcal{O}(1)$   
out 2: Gwen  
... ⚙ ...  $\mathcal{O}(1)$   
out 3: Carson  
... ⚙ ...  $\mathcal{O}(1)$   
out: EOE

$\left. \begin{array}{l} \mathcal{O}(|D|) \\ \mathcal{O}(1) \\ \mathcal{O}(1) \\ \mathcal{O}(1) \\ \mathcal{O}(1) \end{array} \right\} \mathcal{O}(|D| + \|Q_1\|(D))$

### Theorem (Bagan, Durand, Grandjean 2007)

*There is an enumeration algorithm that enumerates  $\llbracket Q \rrbracket(D)$  with constant delay after preprocessing in time  $\mathcal{O}(|D|)$ , for every  $D$  and every free-connex ACQ  $Q$ .*

### Definition (Free-connex ACQs)

A CQ  $Q$  is free-connex acyclic, if  $Q$  and  $Q + R(\bar{x})$  are acyclic, where  $\bar{x} = \text{free}(Q)$ .

In particular: A **binary**  $Q$  is free-connex acyclic, if it is acyclic and  $\text{free}(Q)$  induces a connected subgraph.

Can we do better?

# Our Idea: Precompute an index

## Indexing Phase

Build a data structure  $DS_D$  in time  $t$ .

## Evaluation Phase

Enumerate  $\llbracket Q \rrbracket(D)$  with constant delay after preprocessing in  $\mathcal{O}(|DS_D|)$ .

## Theorem

For every **binary** free-connex ACQ  $Q$ , we can enumerate  $\llbracket Q \rrbracket(D)$  with constant delay after  $\mathcal{O}(|D_{\text{col}}|)$  preprocessing.

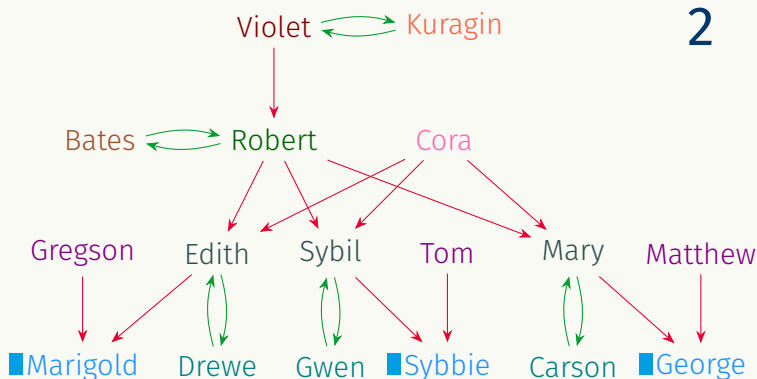
$D_{\text{col}}$  is constructed using *Color Refinement*.

# Identify similar nodes using Color Refinement

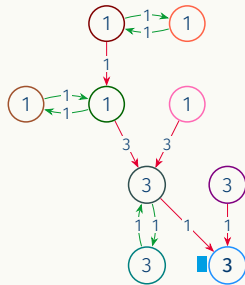
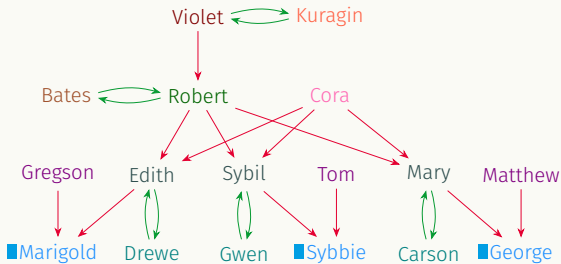
## Color Refinement

*Initial:*  $u, v$  get different colors, iff different unary relations.

*Iterate:*  $u, v$  get different colors, if colors of neighbors differ.

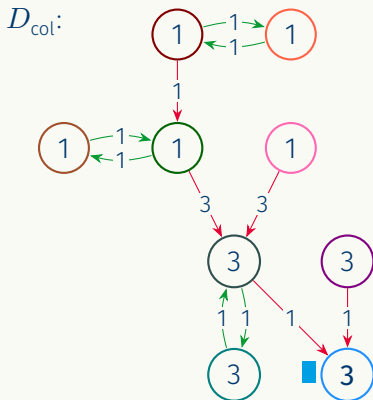
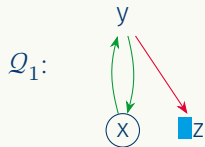


# Create the color database $D_{\text{col}}$



Evaluating queries using  $D_{\text{col}}$ .

# Evaluate $\mathcal{Q}_1$ over $D_{\text{col}}$



`enum( $\mathcal{Q}_1, D_{\text{col}}$ )`

... ⚙ ...

$\mathcal{O}(|D_{\text{col}}|)$

`eop, num: 1`

... ⚙ ...

$\mathcal{O}(1)$

`out 1: ●`

`list( $D, \bullet$ )`

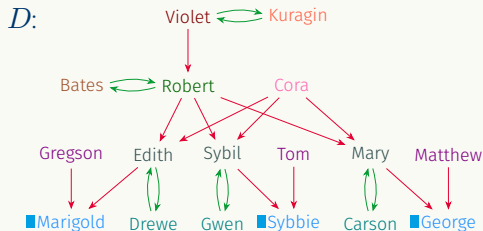
... ⚙ ...

$\mathcal{O}(1)$

`out: EOE`



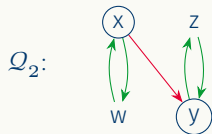
# Evaluate $Q_1$ over database



$\text{list}(D, \bullet) \cdots \text{gear} \cdots \text{Drewe} \cdots \text{gear} \cdots \text{Gwen} \cdots \text{gear} \cdots \text{Carson}^1$   
 $\mathcal{O}(1)$                        $\mathcal{O}(1)$                        $\mathcal{O}(1)$

Enumerated  $\llbracket Q_1 \rrbracket(D)$  with const. delay after  $\mathcal{O}(|D_{\text{col}}|)$  preprocessing!  
[1]: Need additional data structures to do this with constant delay.

# Evaluate $\mathcal{Q}_2$ over $D_{\text{col}}$



`enum( $\mathcal{Q}_2, D_{\text{col}}$ )`

... ⚙️ ...

$\mathcal{O}(|D_{\text{col}}|)$

`eop, num: 2`

... ⚙️ ...

$\mathcal{O}(1)$

`out 1: ●, ●`

`list( $D, P(\bullet, \bullet)$ )`

... ⚙️ ...

$\mathcal{O}(1)$

`out 2: ●, ●`

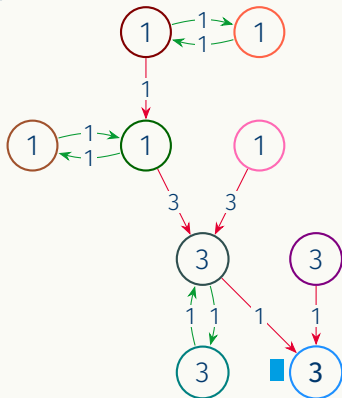
`list( $D, P(\bullet, \bullet)$ )`

... ⚙️ ...

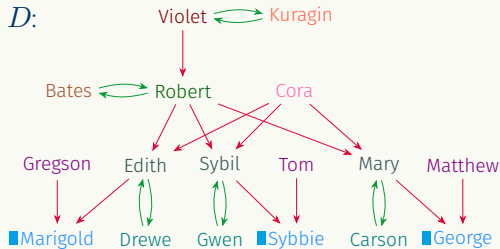
$\mathcal{O}(1)$

`out: EOE`

$D_{\text{col}}$ :



# Evaluate $Q_2$ over database



$\text{list}(D, P(\bullet, \bullet)) \cdots \text{gears} \cdots (\text{Violet}, \text{Robert})$

$\text{list}(D, P(\bullet, \bullet)) \cdots \text{gears} \cdots (\text{Robert}, \text{Edith})$

$\cdots \text{gears} \cdots (\text{Robert}, \text{Sybil})$

$\cdots \text{gears} \cdots (\text{Robert}, \text{Mary})$

$\mathcal{O}(1)$

(Need additional data structures to do this with constant delay.)

I lied to you. 👉👈

This is just the idea.

It is more complicated — but not much.

## Theorem

*For every binary free-connex ACQ  $\mathcal{Q}$ , after  $\mathcal{O}(|D_{\text{col}}|)$  preprocessing, we can return the number of result tuples  $|\llbracket \mathcal{Q} \rrbracket(D)|$  and enumerate them with constant delay.*

- $D_{\text{col}}$  has to be computed only once:  $\mathcal{O}(|D| \cdot \log |D|)$ .
- $|D_{\text{col}}|$  at most  $\mathcal{O}(|D|)$ , smaller if db has “good” structure.

## Outlook:

- Generalize to arbitrary schemas (not trivial).
- Only  $d$  iterations works for queries of height  $\leq d$ .
- DBs are rarely static. How to avoid recomputing the index?

## Further Reading



Christoph Berkholz, Fabian Gerhardt, Nicole Schweikardt.  
“Constant Delay Enumeration for Conjunctive Queries: a  
Tutorial”.

2020. DOI: [10.1145/3385634.3385636](https://doi.org/10.1145/3385634.3385636).



Sandra Kiefer.

“The Weisfeiler-Leman Algorithm: An Exploration of Its  
Power”.

2020. DOI: [10.1145/3436980.3436982](https://doi.org/10.1145/3436980.3436982).

## Some Related Work



Jeroen Bollen, Jasper Steegmans, Jan Van Den Bussche, Stijn Vansummeren.

“Learning Graph Neural Networks Using Exact Compression”.

2023. DOI: [10.1145/3594778.3594878](https://doi.org/10.1145/3594778.3594878).



Martin Grohe, Kristian Kersting, Martin Mladenov, Erkal Selman.

“Dimension Reduction via Colour Refinement”.

2014. DOI: [10.1007/978-3-662-44777-2\\_42](https://doi.org/10.1007/978-3-662-44777-2_42).



Moe Kayali, Dan Suciu.

“Quasi-Stable Coloring for Graph Compression: Approximating Max-Flow, Linear Programs, and Centrality”.

2022. DOI: [10.14778/3574245.3574264](https://doi.org/10.14778/3574245.3574264).

## Other References



Guillaume Bagan, Arnaud Durand, Etienne Grandjean.  
“On Acyclic Conjunctive Queries and Constant Delay Enumeration”.

2007. DOI: [10.1007/978-3-540-74915-8\\_18](https://doi.org/10.1007/978-3-540-74915-8_18).



Holger Dell, Martin Grohe, Gaurav Rattan.  
“Lovász Meets Weisfeiler and Leman”.

2018. DOI: [10.4230/LIPIcs.ICALP.2018.40](https://doi.org/10.4230/LIPIcs.ICALP.2018.40).



Zdeněk Dvořák.

“On recognizing graphs by numbers of homomorphisms”.

2010. DOI: [10.1002/jgt.20461](https://doi.org/10.1002/jgt.20461).



Mihalis Yannakakis.

“Algorithms for Acyclic Database Schemes”.

1981.



# Color Refinement counts homomorphisms

Why does this work?

**Theorem (Dvořák 2010; Dell, Grohe, Rattan 2018)**

*Two vertices  $u, v \in V(G)$  have different colors if, and only if, there is a rooted tree  $(T, r)$  such that*

$$\begin{aligned} & |\{h \in \text{hom}(T, G) \mid h(r) = u\}| \\ & \neq |\{h \in \text{hom}(T, G) \mid h(r) = v\}| \end{aligned}$$