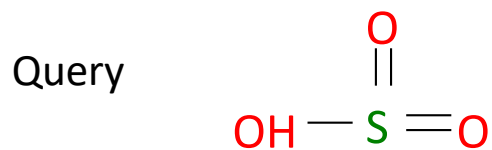# Graph Query Reformulation with Diversity

*Davide Mottin, University of Trento*

Francesco Bonchi, Yahoo Labs - Francesco Gullo, Yahoo Labs

# Pattern search

Query

O
‖
OH — S = O

PubChem Compound [ PubChem Compoun ⬍ ] OS(=O)=O

Save search   Limits   Advanced

Display Settings: ☑ Summary, 20 per page, Sorted by Default order          Send to: ☑

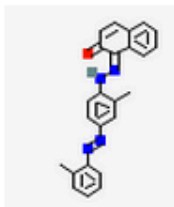Results: 1 to 20 of  510          << First  < Prev  Page [1]  of 49  Next >  Last >>

☐ 1.

- ### Too many matches
- ### Results are not grouped

ubMed (MeSH Keyword)

☐ 2.

IUPAC name: (1Z)-1-[[2-methyl-4-[(2-methylphenyl)diazenyl]phenyl]hydrazi...
Create Date: 2005-09-09
CID: 5876571
Summary   Similar Compounds   Same Parent, Connectivity   PubMed (MeSH Keyword)   Active in 7 of 209 BioAssays

☐ 3.

Basic Yellow 2; Auramine O; Auramin ...
MW: 303.829680 g/mol   MF: $C_{17}H_{22}ClN_3$
IUPAC name: 4-[4-(dimethylamino)benzenecarboximidoyl]-N,N-dimethylanilin...
Create Date: 2005-08-08

db Trento

# Finding specializations



510 matches

*Specializations*

382 matches

46 matches

448 matches

46 matches

114 matches

Graph Query Reformulation with Diversity – Davide Mottin, Francesco Bonchi, Francesco Gullo
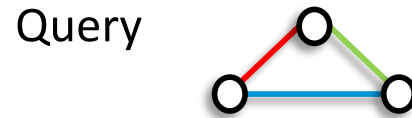
db Trento

# Applications

- Finding groups of molecules having a particular reagent
- Analyze a set of proteins to find diseases
- Workflow optimization
- Anomalies detection in a network
- Finding similar 3D shape search

Graph Query Reformulation with Diversity – Davide Mottin, Francesco Bonchi, Francesco Gullo

db Trento

# Dealing with specializations in web and relational data

- Faceted Search
  - present aspects of the results [Roy08]

- Query reformulation
  - Modify some of the query conditions
    - In structured databases [Mishra09]
    - In web search [Dang10]

## Frist Study of Problem on GRAPHS

db Trento

# Graph Query Reformulation

Query

Results

$R_1$ $R_2$ $R_3$ $R_4$ $R_5$

Reformulations:
query
supergraphs

Exponential number
of reformulations

...

Graph Query Reformulation with Diversity – Davide Mottin, Francesco Bonchi, Francesco Gullo

dbTrento

# Challenges

- The number of reformulation is exponential

- Quantify the interestingness of a reformulation

- Finding query reformulations is **NP**-complete

db Trento

# Our Approach

## Graph Query Reformulation with Diversity

- Finds *k meaningful* specializations efficiently
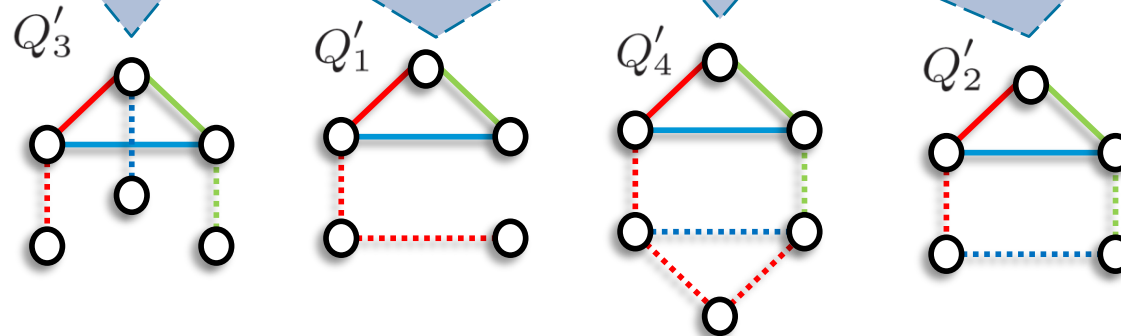
db Trento

# Finding Meaningful Specializations

Query

Results

Find **k meaningful** specializations:
1. Span **all** the results
2. Present **different** aspects of the results

**?**

$cov(Q)$

$Q'$

$Q'_3$  $Q'_1$  $Q'_4$  $Q'_2$

Graph Query Reformulation with Diversity – Davide Mottin, Francesco Bonchi, Francesco Gullo

db Trento

# Graph Query Reformulation with Diversity

**Problem**

Find a set $\mathcal{Q}^*$ of $k$ specializations of $Q$ that maximize a combination of **coverage** and **diversity**

$$f(\mathcal{Q}) = cov(\mathcal{Q}) + \lambda \sum_{Q',Q'' \in \mathcal{Q}} div(Q', Q'')$$

$$\mathcal{Q}^* = \underset{\mathcal{Q} \subseteq \mathbb{S}_Q}{\arg\max} \quad f(\mathcal{Q})$$

$$\text{subject to} \qquad |\mathcal{Q}| = k.$$

**Theorem (NP-hardness)**

The problem reduces to **MAX-SUM Diversification** Problem, so it is NP-hard

db Trento

# Solution: Greedy Algorithm

**Greedy**

While k-specializations are not found

1. Find the specialization leading to the maximum increment of the objective function (marginal gain)

2. Add the specialization to the results

**Theorem**
The algorithm is a ½-approximation

Finding the maximum gain is #P-complete [Valiant79]

**Solution**

**Fast_MMPG**: Branch and bound algorithm with quality guarantees

db Trento

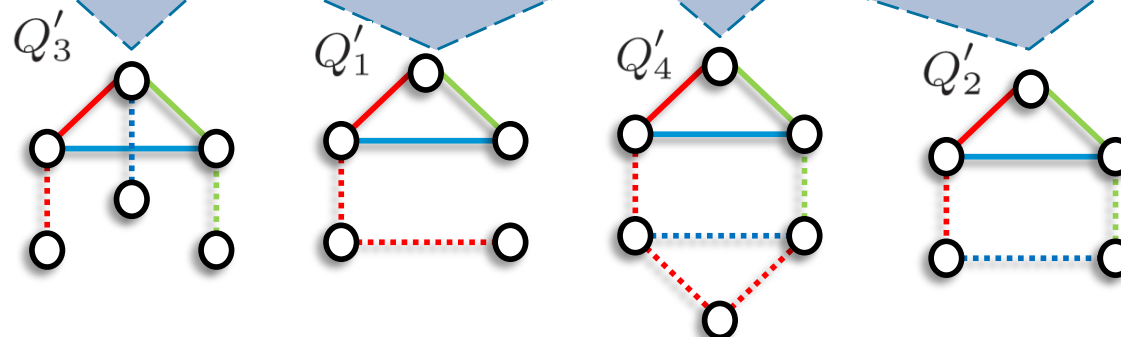# The multiplicity vector



| 2 | 3 | 3 | 3 | 1 |

Results

$R_1$ $R_2$ $R_3$ $R_4$ $R_5$

$Q'_3$ $Q'_1$ $Q'_4$ $Q'_2$

Output set of specializations

Graph Query Reformulation with Diversity – Davide Mottin, Francesco Bonchi, Francesco Gullo

dbTrento

# Upper bound on the Marginal gain

**Lemma**

The marginal gain increases if the multiplicity of the considered item is where $\leq \frac{|\mathcal{Q}|}{2}$ $|\mathcal{Q}|$ is the number of reformulations in the reformulated set constructed so far.

**Upper bound** : is the value of the objective function considering only results with multiplicity $\leq \frac{|\mathcal{Q}|}{2}$

**Theorem**

For a reformulation $Q' \in \mathbb{S} \setminus \mathcal{Q}$ it holds that

$$\max_{Q'' \in \mathcal{T}_{Q'}} \Delta_f(\mathcal{Q}, Q'') \leq \overline{\Delta_f}(\mathcal{Q}, Q') =$$

$$= \frac{1}{2}\vec{u}_{\mathcal{Q}} \cdot \vec{x}_{Q^*} + \lambda \left( \|\vec{m}_{\mathcal{Q}}\| + |\mathcal{Q}| \times \|\vec{x}_{Q^*}\| - 2\vec{m}_{\mathcal{Q}} \cdot \vec{x}_{Q^*} \right).$$

db Trento

# Upper bound

$$\leq \frac{|\mathcal{Q}|}{2}$$

| 1 | 2 | 1 | 1 | 1 |
|---|---|---|---|---|

Results

$R_1$   $R_2$   $R_3$   $R_4$   $R_5$

$Q'_3$   $Q'_1$   $Q'_2$

Output set of
reformulations

dbTrento

# The Fast_MMPG Algorithm

$\overline{\Delta}_f(\mathcal{Q}, Q_1') =$ upper bound

$\Delta_f(\mathcal{Q}, Q_1') =$ marginal gain

$Q$

$\overline{\Delta}_f(\mathcal{Q}, Q_1') = 30$
$\Delta_f(\mathcal{Q}, Q_1') = 18$

$\overline{\Delta}_f(\mathcal{Q}, Q_2') = 21$

$\overline{\Delta}_f(\mathcal{Q}, Q_3') = 26$
$\Delta_f(\mathcal{Q}, Q_3') = 20$

Until the reformulation with the max upper bound and marginal gain is not found

$Q_1'$   $Q_2'$   $Q_3'$

1. Expand the reformulation with the max upper bound
2. Prune Reformulations with marginal gain smaller than the upper bound so far

$Q_{11}'$   $Q_{12}'$   $Q_{31}'$   $Q_{32}'$

$\overline{\Delta}_f(\mathcal{Q}, Q_{11}') = 22$

$\Delta_f(\mathcal{Q}, Q_{11}') = 22$

$\overline{\Delta}_f(\mathcal{Q}, Q_{12}') = 18$

$\overline{\Delta}_f(\mathcal{Q}, Q_{31}') = 18$

$\overline{\Delta}_f(\mathcal{Q}, Q_{32}') = 16$

dbTrento

# Experimental Setup

- ## Datasets:

  - AIDS: 10k chemical compounds

  - Financial: 17k transaction workflows

  - Web: 13k interactions with a recommender system
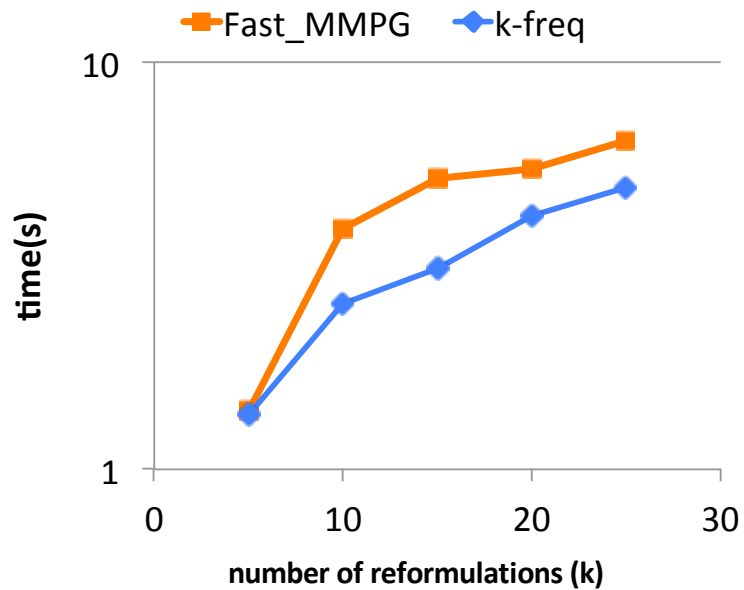
- ## Baseline algorithms:

  - k-freq: returns top-k frequent supergraphs of a query

  - LIndex: informative patterns index

- ## Experiments:

  - Time and objective function value varying k, query size, $\lambda$
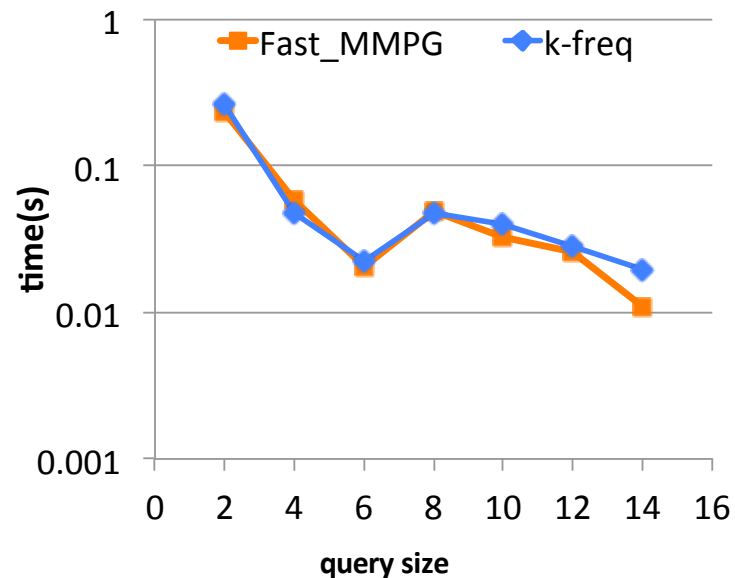
  - Anecdotal

  - Scalability

db Trento

# Time Comparison



**Number of reformulations**
1. k-freq runs only slighly faster
2. Time increases linearly in k
3. Fast_MMPG has real-time performance

**Query size**
1. Fast_MMPG comparable to k-freq
2. Time decreases with query size (less reformulations)

dbTrento

# Objective function gain

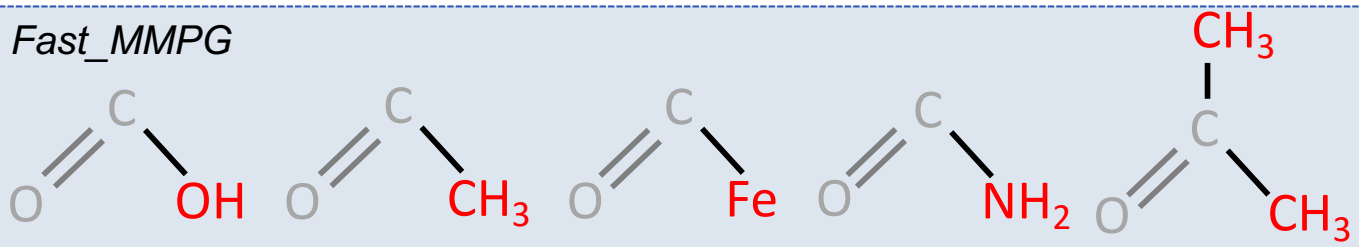| | $\lambda$ | | | | |
|---|---|---|---|---|---|
| | 0 | 0.01 | 0.05 | 0.1 | 0.5 |
| Fast_MMPG | 433 | 613 | 1 345 | 2 260 | 9 566 |
| k-freq | 409 | 540 | 1 063 | 1718 | 6 954 |
| $gain$ (%) | 6 | 12 | 21 | 24 | 27 |

**Analysis**

1. Lambda correctly moves the objective function towards diversity
2. k-freq only captures coverage

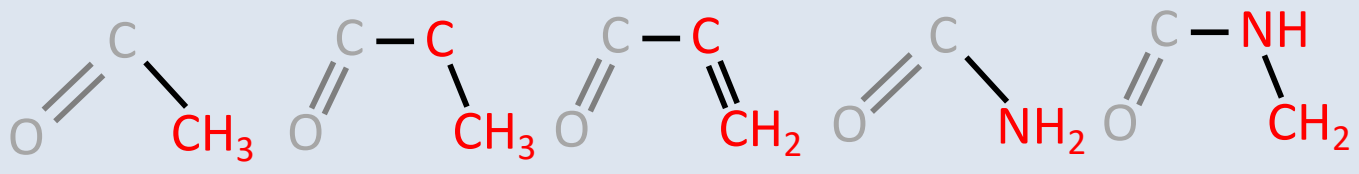$$f(\mathcal{Q}) = cov(\mathcal{Q}) + \lambda \sum_{Q',Q'' \in \mathcal{Q}} div(Q', Q'')$$

db Trento

# Qualitative evaluation

*Query*  $C = O$

*Fast_MMPG*

$$O = C - OH \quad O = C - CH_3 \quad O = C - Fe \quad O = C - NH_2 \quad O = C \begin{array}{c} CH_3 \\ CH_3 \end{array}$$

*k-freq*

$$O = C - CH_3 \quad O = C - C - CH_3 \quad O = C - C = CH_2 \quad O = C - NH_2 \quad O = C \begin{array}{c} NH \\ CH_2 \end{array}$$
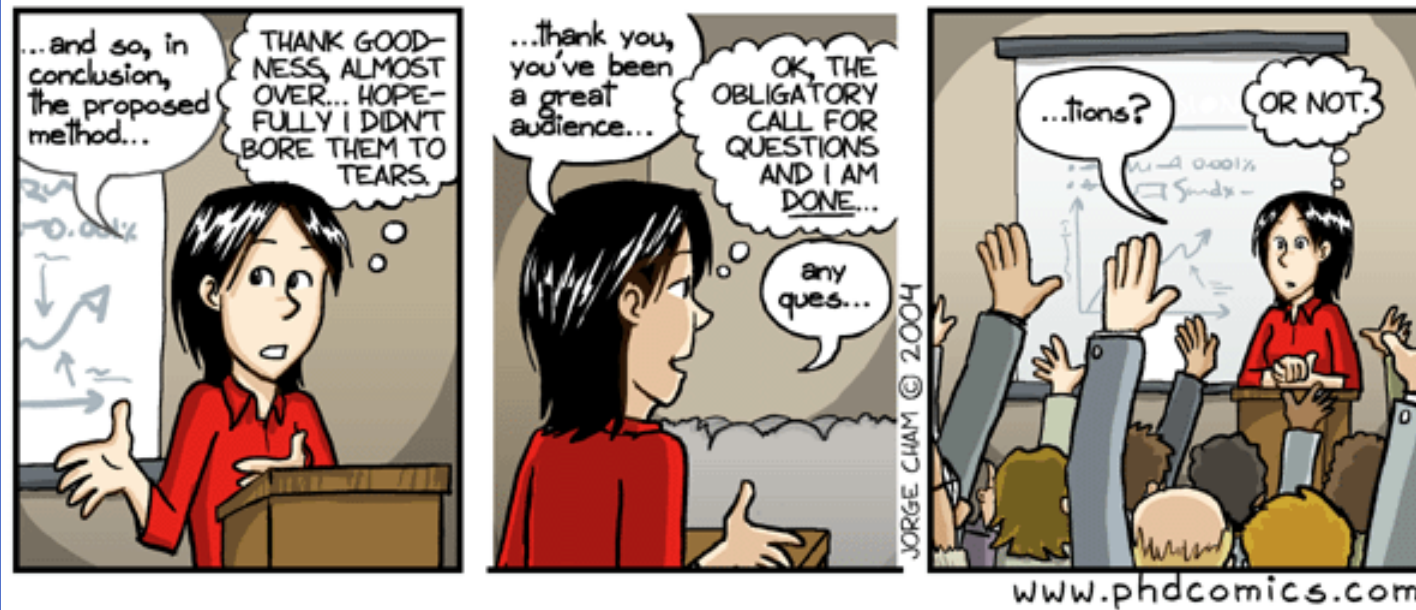
**Analysis**

- k-freq finds reformulation of the same superquery
- Fast_MMPG returns reformulations with more diversified structures

db Trento

# Conclusions

- First study of the problem in **graph databases**
- **Principled** objective function optimizing **coverage** and **diversity**
- Algorithmic solutions **with quality guarantees** and **real time responses**

Thank you!

# Questions?

db Trento