

Session 9

Graph Data 2

Jinsoo Lee, Wook-Shin Han, Romans Kasperovics,
Jeong-Hoon Lee: An In-depth Comparison of
Subgraph Isomorphism Algorithms in Graph
Databases

An In-depth Comparison of Subgraph Isomorphism Algorithms in Graph Databases

- データベース中のグラフ d からクエリグラフ g と同型のサブグラフを検索するアルゴリズムの性能を比較

- 比較対象

- Ullmannの手法 (1976)

小さな
グラフ用

VF2 [Cordella, PAMI04]

QuickSI [Shang, PVLDB08]

GraphQL[He, SIGMOD08]

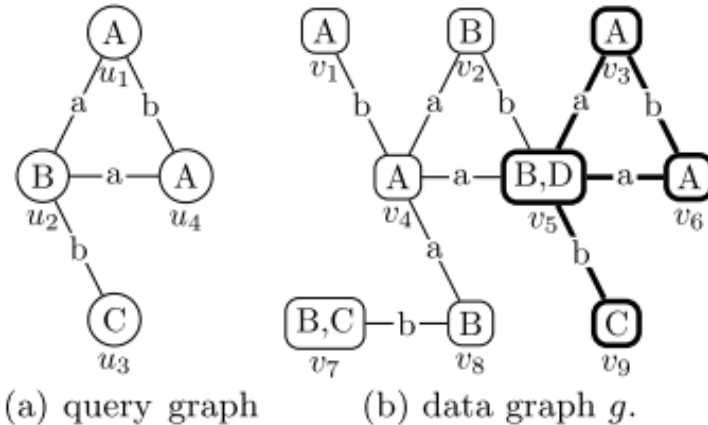
GADDI[Zhang, EDBT09]

SPath[Zhao, PVLDB10]

大規模
グラフ用

- 得られた結果

- QuickSIがGraphQL, GADDI, Spathより性能がよい
- QuickSI, VF2, GADDIは木データに対するサブグラフを妥当な時間内に発見できなかった
- GraphQLは唯一すべてのクエリセットを処理できるが、多くのデータセットでQuickSIより遅い
- SPathは一貫してGraphQLより遅い(使ってる言語違うじゃんというツッコミ)
- 全てのアルゴリズムはjoin order selectionと同じ問題を持つ



アルゴリズムに共通するフレームワーク

Algorithm 1 GENERICQUERYPROC

Input: query graph q

Input: data graph g

Output: all subgraph isomorphisms of q in g

```
1:  $M := \emptyset$ ;  
2: for each  $u \in V(q)$  do  
3:    $C(u) := \text{FILTERCANDIDATES}(q, g, u, \dots)$ ;  
    $[[ \forall v \in C(u)((v \in V(g)) \wedge (L(u) \subseteq L(v))) ]]$   
4:   if  $C(u) = \emptyset$  then  
5:     return;  
6:   end if  
7: end for  
8:  $\text{SUBGRAPHSEARCH}(q, g, M, \dots)$ ;
```

Subroutine $\text{SUBGRAPHSEARCH}(q, g, M, \dots)$

```
1: if  $|M| = |V(q)|$  then  
2:   report  $M$ ;  
3: else  
4:    $u := \text{NEXTQUERYVERTEX}(\dots)$ ;  
    $[[ u \in V(q) \wedge \forall (u', v) \in M(u' \neq u) ]]$   
5:    $C_R := \text{REFINECANDIDATES}(M, u, C(u), \dots)$ ;  
    $[[ C_R \subseteq C(u) ]]$   
6:   for each  $v \in C_R$  such that  $v$  is not yet matched do  
7:     if  $\text{ISJOINABLE}(q, g, M, u, v, \dots)$  then  
8:        $[[ \forall (u', v') \in M((u, u') \in E(q) \implies$   
          $(v, v') \in E(g) \wedge L(u, u') = L(v, v')) ]]$   
9:        $\text{UPDATESTATE}(M, u, v, \dots)$ ;  
        $[[ (u, v) \in M ]]$   
10:       $\text{SUBGRAPHSEARCH}(q, g, M, \dots)$   
11:       $\text{RESTORESTATE}(M, u, v, \dots)$ ;  
        $[[ (u, v) \notin M ]]$   
12:     end if  
13:   end for  
14: end if
```

再帰関数

•再帰の回数をいかに減らすかがキモだと思っていたのだが...

• FILTERCANDIDATES

- クエリ内のノード u に対する候補ノード $C(u)$ を検索する

• NEXTQUERYVERTEX

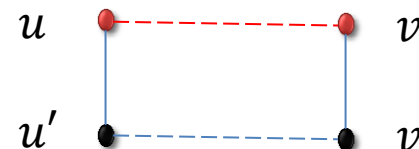
- u' の次にマッチングするノード u を選択する

• REFINECANDIDATES

- 候補ノード $C(u)$ を $C_R \subset C(u)$ に絞り込む

• ISJOINABLE

- クエリ内の (u, u') とDB内の (v, v') がマッチするか



サブグラフ同型アルゴリズム その1 (Ullmannの手法をベースにしたもの)

	Ullmann	VF2	QuickSI
FILTER CANDIDATES	ラベルが一致するノードを候補にする	(Ullmannと同様)	(Ullmannと同様)
NEXT QUERYVERTEX (u' の次ノード u)	u' の次に入力されたノード	既にマッチしたクエリサブグラフ M_q に隣接するノード	ノードラベルと辺ラベルが低頻度な隣接ノード
REFINE CANDIDATE	u より次数の低い $v \in C(u)$ を消す	u が M_q に隣接することを利用した枝刈(※1)	最小全点木を利用した枝刈(※2)
ISJOINABLE	u の隣接ノードをすべてチェックする	(Ullmannと同様)	必要ない (実は結構重要)

※1: 後で書く

※2: クエリグラフで最小全点木上の u の親ノード u' に対応するDBグラフ上のノード v' と v がつながっていなければ消す

近接シングネチャ (GraphQL, SPath)

- GraphQL neighborhood signature

- 隣接ノードのラベルの集合

$$sig_{GraphQL}(u_2) = \{A, A, C\}$$

$$sig_{GraphQL}(v_2) = \{A, B, D\}$$

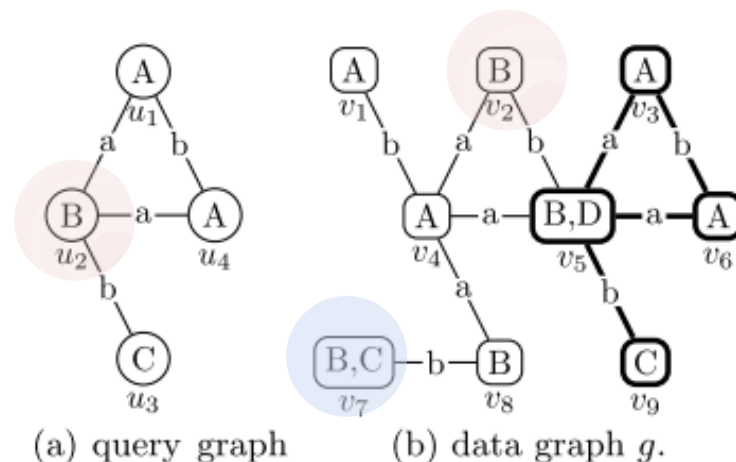
- SPath neighborhood signature

- 近接ノードに関するtripleの集合

- $(d, l, c) : u$ からの最短距離が d でラベルが l である点が c 個ある

$$sig_{GraphQL}(v_7) = \{(1, B, 1)\} \quad (k_0 = 1 \text{ の場合})$$

$$sig_{GraphQL}(v_7) = \{(1, B, 1), (2, A, 1)\} \quad (k_0 = 2 \text{ の場合})$$



サブグラフ同型アルゴリズム その2 (シグネチャを使う方法)

- シグネチャを使って早めに候補を少なくし、再帰の回数をできるだけ削減する

	GADDI	GraphQL	SPath
FILTER CANDIDATES	(省略)	Ullmannの手法＋枝刈 i) $sig_{GraphQL}(u) \not\subseteq sig_{GraphQL}(v)$ な v を排除 ii) 擬似同型を使った枝刈 (※1)	Ullmannの手法＋枝刈 ・ 距離 k までのノード数の合計が u より少ないラベルがあったら排除
NEXT QUERYVERTEX (u'の次ノードu)	(省略)	中間結果のサイズを最小にする隣接ノードを選択	<hr/> Subroutine SUBGRAPHSEARCH (q, g, M, k_0) 1: if $ M = V(q) $ then 2: report M ; 3: else 4: $p_q :=$ NEXTQUERYPATH (q, g, k_0); 5: $P :=$ GETCANDIDATEPATHS (p_q, M, C); 6: for each $p_g \in P$ do 7: if ISJOINABLE (p_q, p_g, M) then 8: UPDATESTATE (p_q, p_g, M); 9: SUBGRAPHSEARCH (q, g, M, k_0); 10: RESTORESTATE (p_q, p_g, M); 11: end if 12: end for 13: end if <hr/>
REFINE CANDIDATE	(省略)		
ISJOINABLE	(省略)		

ノード単位でなく
Path単位でマッチングする

Table 1: Summary of the real-world datasets.

Dataset	AIDS	NASA	Yeast	Human
# of graphs	10000	36790	1	1
# of vertices	2~214	2~889	3112	4675
# of edges	1~217	1~888	12519	86282
Avg. degree	1.95	1	8.05	36.82
Max. degree	11	245	168	771
# of distinct v. labels	51	117302	184	90
# of distinct e. labels	4	0	0	0
Avg. # of labels per vertex	1	1	7.55	4.63

評価

問合せ処理

AIDS: 小さく疎なグラフ

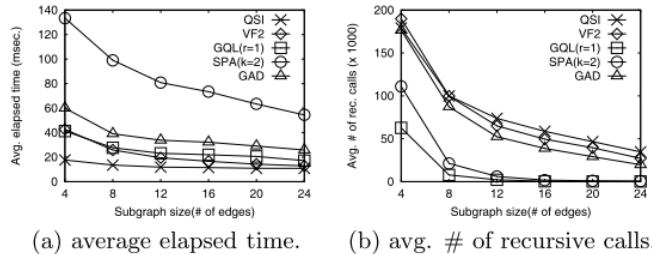


Figure 6: AIDS dataset.

NASA: ラベルの種類が多い木構造データ

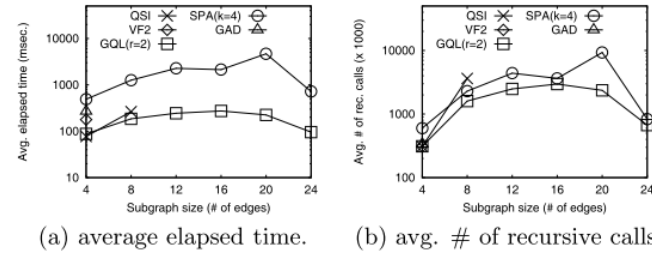


Figure 7: NASA dataset.

- SPathは再帰回数は減っているが一番遅い
- QuickSIは再帰回数多いが速い
 - ISJOINABLEをしなくていいのが効いている
- GraphQLはFILTERCANDIDATESの処理が全体の63.16%
 - SPathは61.23%
- SPathはPathベースのマッチングでオーバーヘッドがかかっている

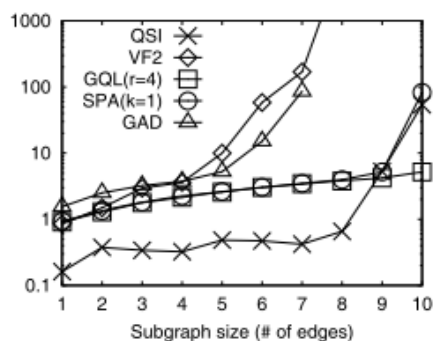
- VF2, QuickSI, GADDIは指数オーダで増加して結果が得られなくなった
- GraphQL, SPathは最後まで処理できた

シグネチャによる早めの枝刈が効いている

- GraphQLはSPathより平均10.38倍速い
 - シグネチャの生成時間(5.5倍)
 - FILTERCANDIDATESの時間(2.45倍)

評価

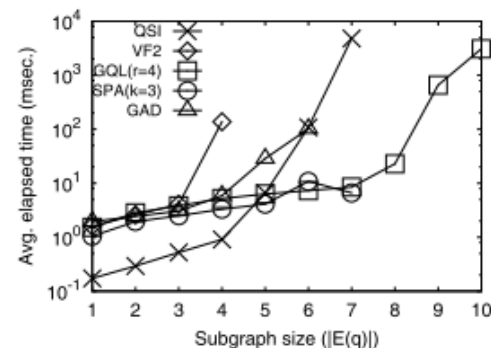
Yeast:1個の大きく密なグラフ



(a) subgraph queries.

- QuickSI, GraphQL, SPath以外は小さなサブグラフで計算不能に
- QuickSIが最も性能がよい
- シグネチャによる枝刈があまり効果がない
 - 半数くらいにしか減らない

Human:Yeastより大きく密な1つのグラフ



(a) subgraph queries.

- GraphQL以外は小さなサブグラフで計算不能に
- クエリサイズが大きいと探索空間サイズが大きくなるため
- join order selection problem