# **Frequent Connected Subgraph Mining**



Tamás Horváth

University of Bonn & Fraunhofer IAIS, Sankt Augustin, Germany tamas.horvath@iais.fraunhofer.de



Fraunhofer

# **Frequent Connected Subgraph Mining**

# Outline

- motivation, problem definition, and a negative complexity result
- mining trees with the levelwise search algorithm
- mining bounded tree-width graphs



RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT



# Notions: Isomorphism and Subgraph Isomorphism

 $G_1 = (V_1, E_1), G_2 = (V_2, E_2)$ : labeled graphs

· each vertex and edge is associated with some symbol of an alphabet

 $G_1$  is **isomorphic** to  $G_2$  if there is a bijection  $\varphi: V_1 \to V_2$  preserving

· the edges in both directions

- i.e.,  $\{u, v\} \in E_1$  if and only if  $\{\varphi(u), \varphi(v)\} \in E_2$ 

- the vertex labels
  - i.e.,  $\varphi(u)$  has the same label as u for all  $u \in V_1$
- the edge labels
  - i.e.,  $\{\varphi(u), \varphi(v)\}$  has the same label as  $\{u, v\}$  for all  $\{u, v\} \in E_1$
- $G_1$  is **subgraph isomorphic** to  $G_2$  if  $G_2$  has a subgraph isomorphic to  $G_1$





# Mining Frequent Connected Subgraphs

### frequent connected subgraph mining problem:

Given a set *D* of labeled graphs and an integer t > 0, list the set of *t*-frequent connected subgraphs w.r.t. *D* 

- t > 0 integer: *frequency threshold*
- *t*-frequent subgraph: subgraph isomorphic to at least t graphs in D

instance of the theory extraction problem:

- D : set of labeled graphs
- *L*: set of all labeled *connected* graphs
- *interestingness predicate*  $q_D$ : for a pattern  $H \in \mathcal{L}$ ,  $q_D(H)$  is true iff H is subgraph isomorphic to at least t graphs in D
  - $q_D$  is *anti-monotone*: any connected subgraph of a *t*-frequent connected graph is also *t*-frequent





5

### **Frequent Subgraph Mining: Motivation**

#### Virtual screening in drug discovery:

select a limited number of candidate compounds from millions of database molecules that are most likely to possess a desired biological activity



### **Virtual Screening in Drug Discovery**





molecules give rise to labeled undirected graphs

PhD Course, Szeged, 2012 - © T.Horváth



RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT



## **Virtual Screening in Drug Discovery**

approach [Deshpande, Kuramochi, Wale, & Karypis, 2005]:

- 1. compute the set  $\{p_1, \ldots, p_k\}$  of **frequent connected subgraphs** for the molecular graphs of the training molecules
- 2. assign a **binary colored (feature) vector** v of length k to each training molecule m with

 $v[i] := \begin{cases} 1 & \text{if } p_i \text{ is a subgraph of the molecular graph of } m \\ 0 & \text{o/w} \end{cases}$ 

 $(i=1,\ldots,k)$ 

- $\Rightarrow$  each molecule is represented by a colored vertex of the k-dimensional unit hypercube
  - green: active
  - red: inactive





### **Virtual Screening in Drug Discovery**

approach [Deshpande, Kuramochi, Wale, & Karypis, 2005] (cont'd):

- 3. compute a hyperplane h in the k-dimensional space that separates the green points from the red ones, as good as possible
  - Support Vector Machines
     [Boser, Guyon, & Vapnik, 1992; Cortes & Vapnik, 1995]:
    - choose the hyperplane with maximum distance from nearest data points
- 4. for each test molecule, compute the vector as defined in step 2 and predict its activity according to the halfspace defined by h it belongs to
  - empirical experiments: good predictive performance
  - rest of the lecture: how to perform step 1?



# **Enumeration Complexity**

the **size** of the output (theory) can be **exponential** in the size of the input *D* 

 $\Rightarrow$  the output cannot be computed in time polynomial in the size of D

#### enumeration complexities:

a set of S with N elements, say  $s_1, \ldots, s_N$ , are listed with

- **polynomial delay** if the time before printing  $s_1$ , the time between printing  $s_i$  and  $s_{i+1}$  for every i=1,...,N-1, and the termination time after printing  $s_N$  is bounded by a polynomial of the size of the input,
- **incremental polynomial time** if  $s_1$  is printed with polynomial delay, the time between printing  $s_i$  and  $s_{i+1}$  for every i=1,...,N-1 (resp. the termination time after printing  $s_N$ ) is bounded by a polynomial of the combined size of the input and the set  $s_1,...,s_i$  (resp. S),
- output polynomial time if S is printed in the combined size of the input and the entire set S

PhD Course, Szeged, 2012 - © T.Horváth





## A Negative Complexity Result

**Thm**: The frequent connected subgraph mining problem cannot be solved in output-polynomial time (unless P = NP).

### Proof:

reduction: Hamiltonian path problem

- Hamiltonian path problem:
  - *Given* a graph *G* with *n* vertices, decide whether or not *G* has a Hamiltonian path, i.e., a path containing each vertex of *G* exactly once
  - NP-complete problem





# **Proof of the Negative Complexity Result**

- $D = \{G_1, G_2\}$ , where
  - $G_1$  is an arbitrary unlabeled graph with n vertices
  - $G_2$  is an unlabeled path of length n-1 (i.e., has n vertices)
- $\mathcal{L}$  : set of all unlabeled connected graphs
- $q_D$ : 2-frequency w.r.t. D
- $\Rightarrow Th(\mathcal{L}, D, q_D)$  is the set of 2-frequent paths
  - the size of  $Th(\mathcal{L}, D, q_D)$  is polynomial in the size of  $D(|Th(\mathcal{L}, D, q_D)| \le n)$
  - $|Th(\mathcal{L}, D, q_D)| = n$  if and only if  $G_1$  has a Hamiltonian path
  - cannot be computed in output polynomial time, as otherwise the NPcomplete Hamiltonian path problem could be decided in polynomial time





# **Frequent Connected Subgraph Mining**

# Outline

- motivation, problem definition, and a negative complexity result
- mining trees with the levelwise search algorithm
- mining bounded tree-width graphs





## A Generic Levelwise Search Graph Mining Algorithm



PhD Course, Szeged, 2012 - © T.Horváth



Fraunhofer

### **Efficiency Conditions for the Generic Graph Mining Algorithm**

**Thm:** Let  $\mathcal{G}$  be a graph class satisfying

- (i) G is closed under taking subgraphs (i.e.,  $\forall G \in G$ , all subgraphs of G belong to G),
- (ii) the membership problem in  $\mathcal{G}$  (i.e., does  $G \in \mathcal{G}$  hold for any graph G) can be decided in polynomial time, and
- (iii) for every  $H, G \in \mathcal{G}$  such that H is connected, it can be decided in polynomial time whether H is subgraph isomorphic to G.

If the transaction graphs in D belong to G then the previous algorithm lists the frequent connected subgraphs with polynomial delay.







# **Efficiency Conditions for the Generic Graph Mining Algorithm**

### Proof (sketch):

- the cardinalities (and hence, the sizes) of the sets  $\rho(H) \cap \mathcal{G}$  and  $\rho^{-1}(H)$  in line 6 are bounded by a polynomial of the size of D,
- both sets can be computed in polynomial time,
- conditions (i) and (iii) together imply that one can define a canonical string representation for the graphs in *G* that can be computed in time polynomial in the size of *D*.
  - canonical string representation: unique modulo isomorphism (i.e., two graphs have the same canonical strings if and only if they are isomorphic)
- $\Rightarrow$  using some advanced (e.g., trie-based) data structure for the storage of  $S_l$  and the elements of  $C_{l+1}$  generated before H, conditions (i) and (ii) in line 6 can be decided in time polynomial in the size of D





# **Application: Frequent Subtree Mining in Forests**

#### problem:

Given a set D of labeled *forests* and an integer t > 0, list all trees that are subtrees of at least t forests in D.

- forest: set of vertex disjoint labeled free trees
- free tree: unordered, unrooted tree



# **Frequent Subtree Mining in Forests**

**Thm:** The frequent connected subgraph mining problem can be solved with *polynomial delay* for **forest** transaction graphs.

### proof:

each condition of the previous theorem holds, i.e.,

- (i) forests are *closed downward*,
- (ii) it can be decided in polynomial time, whether a graph G is a forest,
- (iii) *subtree isomorphism* can be decided in polynomial time
  - in time O(n<sup>2.5</sup>) [Matula,1978]
  - can further be improved by a *log* factor [Shamir &Tsur,1999]





## Subtree Isomorphism

- T tree, F forest; decide, whether T is subgraph isomorphic to F: decide for all trees T' in F, whether T is subgraph isomorphic to T'
- T and T' are trees; decide, whether T is subgraph isomorphic to T': decide for some fixed vertex u in T and for all vertices v in T', whether there is a subgraph isomorphism from the tree T rooted at u to the tree T' rooted at v that maps u to v
  - problem is reduced to subtree isomorphism between labeled, rooted, unordered trees
  - can be solved with a **bottom-up** (recursive) algorithm (next slide)





### **Bottom-up Subtree Isomorphism algorithm for Rooted Trees**

- 1.  $T_u$ : tree T rooted at u;  $u_1, \ldots, u_m$ : children of u in  $T_u$  $T'_{v}$ : tree T' rooted at v;  $v_1, \ldots, v_n$ : children of v in  $T_{v}$
- 2. decide (recursively) whether there is a subgraph isomorphism from the subtree of  $T_u$  rooted at  $u_i$  to the subtree of  $T'_v$  rooted at  $v_i$  that maps  $u_i$ to  $v_i$  for every  $i = 1, \ldots, m$  and  $j = 1, \ldots, n$ ; if so, **mark** the pair (i, j) if u and v have the same vertex label, and the edges between u,  $u_i$  and between v,  $v_j$  have the same edge labels
- 3. return YES (i.e.,  $T_u$  is subgraph isomorphic to  $T'_v$ ) if the bipartite graph

$$(\{u_1, \ldots, u_m\}, \{v_1, \ldots, v_n\}, \{\{u_i, v_j\} : (i, j) \text{ is marked }\})$$

has a (maximum) matching of size m

maximum bipartite matching can be solved in polynomial time



universitä

## Summary

- mining frequent connected subgraphs in arbitrary transaction graphs is computationally hard
  - cannot be solved in output-polynomial time (unless P = NP)
- for forest transaction graphs, the problem can be solved with polynomial delay
  - obtained by using a generic levelwise search algorithm
  - frequent patterns are not printed immediately after their generation
    - · polynomial delay vs. incremental polynomial time



# **Frequent Connected Subgraph Mining**

# Outline

- motivation, problem definition, and a negative complexity result
- mining trees with the levelwise search algorithm
- mining bounded tree-width graphs





## Positive and Negative Results So Far

frequent connected subgraph mining:

- © computationally **intractable** for **arbitrary** transaction graphs
  - cannot be solved in output-polynomial time (unless P = NP)
- © can be solved **efficiently** for **forest** transaction graphs
  - with polynomial delay

**Goal:** Generalize the positive result on forests to a broader graph class!

- What about graphs of bounded tree-width?
  - parameterized graph class (naturally) generalizing forests





## A Generic Levelwise Search Graph Mining Algorithm (recap)



PhD Course, Szeged, 2012 - © T.Horváth



RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT

Fraunhofer

# **Efficiency Conditions (recap)**

**Thm:** Let  $\mathcal{G}$  be a graph class satisfying

- (i) G is closed under taking subgraphs (i.e.,  $\forall G \in G$ , all subgraphs of G belong to G),
- (ii) the membership problem in G (i.e., does  $G \in G$  hold for any graph G) can be decided in polynomial time, and
- (iii) for every  $H, G \in \mathcal{G}$  such that H is connected, it can be decided in polynomial time whether H is subgraph isomorphic to G.

If the transaction graphs in D belong to G then the previous algorithm lists the frequent connected subgraphs with polynomial delay.

### Problem:

- <sup>(2)</sup> subgraph isomorphism is NP-complete even for graphs of tree-width 2
- © Condition (iii) can be relaxed!

PhD Course, Szeged, 2012 - © T.Horváth





# **Problem Setting**

### Given

database D: finite set of labeled graphs of bounded tree-width

• will be defined

language  $\mathcal{L}$ : set of all connected labeled graphs of bounded tree-width,

interestingness predicate  $q_D$ :  $\varphi \in \mathcal{L}$  is "interesting" w.r.t. D if it is subgraph isomorphic to at least t graphs in D

• t > 0 integer: *frequency threshold* 

**compute** the theory of D w.r.t.  $\mathcal{L}$  and  $q_D$ , i.e.,

$$Th(\mathcal{L}, D, q_D) = \{\varphi \in \mathcal{L} : q_D(\varphi) = \mathsf{true}\}$$

• special case of the theory extraction problem



### **Main Result**

Thm [H. & Ramon, 2010]:

the frequent connected subgraph mining problem can be solved in *incremental polynomial time* for graphs of *bounded tree-width* 

### significance of this result:

- efficient pattern mining is possible even for computationally hard pattern matching operators
  - subgraph isomorphism is NP-complete for bounded tree-width graphs
- first positive non-trivial result beyond trees
- positive result for a practically relevant graph class
  - e.g., molecular graphs of most pharmacological compounds have tree-width  $\leq 3$



## Example



### **NCI Chemical Dataset:**

250251 compounds

tree-width #molecules

0	13	isolated vertices
1	21950	trees
2	221675	mostly outerplanar
3	6548	
≥4	65	





# **Outline for the Rest (Technical Part) of this Topic**

- tree-width
- subgraph isomorphism for bounded tree-width graphs
- remarks and open problems







## Tree-width (Robertson & Seymour, 1986)

- tree decomposition of a graph G = (V, E): a pair  $TD = (T, \mathcal{X})$ , where
  - T = (I, F) is an unordered tree,
  - $\mathcal{X} = \{ bag(i) : i \in I \}$  is a family of subsets of V, called **bags**, s.t.
    - \*  $\bigcup_{i\in I} \mathsf{bag}(i) = V$ ,
    - \* for every  $\{u, v\} \in E$  there is an  $i \in I$  with  $\{u, v\} \subseteq bag(i)$ ,
    - \* for all  $i, j, k \in I$ , if j is on the path from i to k in T then  $bag(i) \cap bag(k) \subseteq bag(j)$
- tree-width of TD:  $\max_i |\mathsf{bag}(i)| 1$
- tree-width of G: minimum tree-width over any TD
- graphs of bounded tree-width: graphs with tree-width bounded by some constant  $\boldsymbol{k}$

PhD Course, Szeged, 2012 - © T.Horváth





## Tree-width (Robertson & Seymour, 1986)

- **measure** of the tree-likeness of graphs
  - e.g., the tree-width of trees is 1 and the tree-width of cycles is 2
- **useful** tool in the design of algorithms because
  - many computationally hard problems on graphs become polynomial for graphs of bounded tree-width
  - many practically relevant graph classes have small tree-width
    - e.g., *k*-outerplanar graphs have tree-width at most 3*k*-1













### Some Properties of Bounded Tree-width Graphs

- ③ the class of bounded tree-width graphs is *closed downward* 
  - any subgraph of a graph of tree-width at most *k* has tree-width at most *k*
- © membership problem can be decided in linear time
  - for **constant** *k*, one can decide in linear time, whether a graph has tree-width at most *k*, and if so, compute a tree-decomposition of tree-width at most *k* 
    - [Bodlaender, 1996]
- Subgraph isomorphism remains NP-complete for graphs of bounded tree-width
  - NP-complete if the pattern is not k-connected or has more than O(k) vertices of unbounded degree; o/w it can be decided in polynomial time
    - [Gupta & Nishimura, 1996]







# Mining Bounded Tree-Width Graphs

- ⇒ generic mining algoritm: candidate generation/test is **not** directly applicable
  - because subgraph isomorphism is NP-complete
- ⇒ polynomial delay: open question

### What about incremental polynomial time?

- modify the generic levelwise search graph mining algorithm
  - slide 34
  - **changes:** print frequent patterns directly after their generation (slide 35)







## A Generic Levelwise Search Graph Mining Algorithm (recap)



PhD Course, Szeged, 2012 - © T.Horváth



RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT

Fraunhofer

# **Modified Levelwise Search Graph Mining Algorithm**

**Input** : database D of transaction graphs and integer t > 0**Output**: all frequent connected subgraphs

MAIN:

```
1: let S_0 \subseteq \mathcal{G} be the set of frequent graphs consisting of a single labeled vertex
```

```
2: print S<sub>0</sub>
```

3: for 
$$(l := 0; S_l \neq \emptyset; l := l + 1)$$
 do

4: 
$$C_{l+1} := S_{l+1} := \emptyset$$

5: forall  $P \in S_l$  do

6: forall 
$$H \in \rho(P) \cap \mathcal{G}$$
 satisfying (i)  $H \notin C_{l+1}$  and (ii)  $\rho^{-1}(H) \subseteq S_l$  do

- 7: add H to  $C_{l+1}$
- 8: **if** SUPPORTCOUNT(H)  $\geq t$  **then**
- 9: **print** *H* and add it to  $S_{l+1}$
- 10: add H to  $S_{l+1}$

#### SUPPORT COUNT(H, D, t):

- 1: counter := 0
- **2:** forall G in D do
- 3: if  $H \preccurlyeq G$  then
- 4: counter := counter + 1
- 5: return counter

PhD Course, Szeged, 2012 - © T.Horváth



## **Efficiency Conditions for the Modified Graph Mining Algorithm**

**Thm:** Let  $\mathcal{G}$  be a graph class satisfying

- (i)  $\mathcal{G}$  is closed under taking subgraphs,
- (ii) the membership problem in  $\mathcal{G}$  can be decided in polynomial time,
- (iii\*) for any strong candidate pattern H and transaction graph  $G, H \preccurlyeq G$  can be decided in time polynomial in the combined size of G and the set of frequent patterns generated before H, and
- (iv) isomorphism for  $\mathcal{G}$  can be decided in polynomial time.

If the transaction graphs in D belong to G then the previous algorithm lists the frequent connected subgraphs in incremental polynomial time.

proof: exercise


# **Application to Mining Bounded Tree-Width Graphs**

Let  $\mathcal{G}$  be the class of graphs of tree-width at most k for some constant k.

- (i)  $\mathcal{G}$  is closed under taking subgraphs.
- (ii) The membership problem in  $\mathcal{G}$  can be decided in polynomial time.
  - one can decide in linear time, whether a graph has tree-width at most k [Bodlaender, 1996]
- (iv) Isomorphism between graphs of tree-width at most k can be decided in polynomial time [Bodlaender, 1990].
- ⇒ to show that frequent connected subgraph mining in bounded tree-width graphs can be done in incremental polynomial time, we need to prove condition (iii\*)







# Mining Bounded Tree-width Graphs

H: strong candidate pattern generated by levelwise search

- G: transaction graph
  - both of tree-width at most k
- $\mathcal{F}_H$ : H + set of all frequent connected subgraphs generated by the modified levelwise search algorithm before H
- **Thm:** Given *H* and *G* above, it can be decided in time  $poly(size(G), size(\mathcal{F}_H))$ , whether *H* is subgraph isomorphic to *G* 
  - poly(size(G), size(F<sub>H</sub>)): polynom of the combined size of G and the set of frequent patterns computed before H
  - incremental polynomial time

rest of this and next talk(s): proof sketch

PhD Course, Szeged, 2012 - © T.Horváth





### Main Idea of the Proof

to decide subgraph isomorphism from a strong candidate pattern H into a transaction graph G, we can utilize the information computed earlier for the subgraphs of H

- strong candidate pattern: each of its subgraphs is frequent
- only non-redundant information, certain set of tuples, is actually required
- number of non-redundant tuples is bounded by  $poly(size(G), size(\mathcal{F}_H))$
- we can identify a superset  ${\cal S}$  of the set of non-redundant tuples such that
  - the size of S is bounded by  $poly(size(G), size(\mathcal{F}_H))$
  - each tuple in S can be computed in time  $poly(size(G), size(\mathcal{F}_H))$







# **Preprocessing Step**

compute a nice tree-decomposition TD(G) for every transaction graph G

- nice tree-decomposition: rooted tree with 3 different types of nodes
  - leaf node: it has no children
  - separator node z: has a single child z' such that

 $\mathsf{bag}(z) \subseteq \mathsf{bag}(z')$ 

– join node z: has two children  $z_1$  and  $z_2$  such that

 $\mathsf{bag}(z) = \mathsf{bag}(z_1) \cup \mathsf{bag}(z_2)$ 

- size of TD(G) is linear in the size of G
- · can be computed in linear time for graphs of bounded tree-width

use TD(G) for the entire mining process

PhD Course, Szeged, 2012 - © T.Horváth





#### **Nice Tree-Decomposition**





41



RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT

universität**bon** 

# Iso-Quadruples

given pattern graph H, transaction graph G, and tree-decomposition TD(G)

- iso-quadruple of H relative to a node z of TD(G):  $(S, \mathcal{D}, K, \psi)$  such that
  - $S \subseteq V(H)$  satisfying  $|S| \leq k+1$ ,
  - $\mathcal{D}$  is a subset of the set of connected components of the induced subgraph  $H[V(H)\setminus S]$ ,
  - K is the induced subgraph  $H[S \cup V(\mathcal{D})]$ ,
  - $\psi$  is an injective function mapping S to  $\mathrm{bag}(z)$



# Iso-Quadruples

*z*-characteristic of H for a node z in TD(G):

iso-quadruple  $(S, \mathcal{D}, K, \psi)$  of H relative to z such that there is a subgraph isomorphism  $\varphi$  from K into  $G_{[z]}$  satisfying

- $\varphi(u) = \psi(u)$  for every  $u \in S$ ,
- $\varphi(v) \not\in \mathsf{bag}(z)$  for every  $v \in V(\mathcal{D})$

**induced subgraph** of *G* defined by the union of the bags of *z*'s descendants (*z* is also a descendant of itself)

**Lemma:** *H* is subgraph isomorphic to *G* if and only if there exists an *r*-characteristic  $(S, \mathcal{D}, \boldsymbol{H}, \psi)$  for the root *r* of TD(G)

proof: exercise

- $\Rightarrow$  check whether an iso-quadruple relative to the root is a characteristic
  - · the number of iso-quadruples to be checked for the root is bounded by

 $O(|V(H)|^{k+1})$ 





# **Computing Characteristics**

[Matoušek & Thomas, 1992; also Hajiaghayi & Nishimura, 2007]: compute the set of *z*-characteristics for every node *z* in *TD(G)* with **dynamic programming**:

- postorder traversal of *TD(G)* 
  - straightforward for *leaf* nodes (next slide),
  - use only the characteristics of the child(ren) for separator and join nodes (next slides)

**notations:** for pattern graph *H*, transaction graph *G*, both of bounded tree-width, nice tree-decomposition TD(G), and node *z* in TD(G):

- $\Gamma(H,z)$  denotes the set of iso-quadruples of H relative to z
- $\Gamma_{ch}(H,z)$  denotes the set of *z*-characteristics of *H*



# **Computing Characteristics: Leaf Nodes**

**Leaf Lemma:** Let G and H be graphs of bounded treewidth and z be a leaf node in TD(G). Then, for every  $(S, \mathcal{D}, K, \psi) \in \Gamma(H, z)$ 

 $(S, \mathcal{D}, K, \psi) \in \Gamma_{\mathrm{ch}}(H, z) \iff$ 

 $\mathcal{D} = \emptyset$  and  $\psi$  is a subgraph isomorphism from H[S] to G[bag(z)]

#### proof: exercise

•  $\mathcal{D} = \emptyset$  implies K = H[S]



### **Computing Characteristics: Leaf Nodes**



### **Computing Characteristics: Separator Nodes**

**Separator Lemma:** Let H, G be as in the previous case and z be a separator node in TD(G) with (a single) child z'. Then for all  $(S, \mathcal{D}, K, \psi) \in \Gamma(H, z)$ 

 $(S, \mathcal{D}, K, \psi) \in \Gamma_{\mathrm{ch}}(H, z) \iff \exists (S', \mathcal{D}', K', \psi') \in \Gamma_{\mathrm{ch}}(H, z') \text{ such that }$ 

- $\bullet \ S=\{v\in S':\psi'(v)\in \mathsf{bag}(z)\},$
- $\mathcal{D}'$ : set of all connected components C of  $H[V(H) \setminus S']$  such that C is a subgraph of (a connected component in)  $\mathcal{D}$ ,
- $\psi(v) = \psi'(v)$  for every  $v \in S$ .

Proof: see [Hajiaghayi & Nishimura, 2007]



universitatbo



# **Computing Characteristics: Separator Nodes**







## **Computing Characteristics: Join Nodes**

**Join Lemma:** Let H, G be as in the previous cases and z be a join node in TD(G) with children  $z_1$  and  $z_2$ . Then for all  $(S, \mathcal{D}, K, \psi) \in \Gamma(H, z)$ 

 $\begin{aligned} (S,\mathcal{D},K,\psi) \in \Gamma_{\mathrm{ch}}(H,z) \iff \\ \exists (S_1,\mathcal{D}_1,K_1,\psi_1) \in \Gamma_{\mathrm{ch}}(H,z_1) \land \exists (S_2,\mathcal{D}_2,K_2,\psi_2) \in \Gamma_{\mathrm{ch}}(H,z_2) \text{ s.t.} \end{aligned}$ 

- $S_i = \{v \in S : \psi(v) \in bag(z_i)\}$  for i = 1, 2,
- the connected components of  $\mathcal D$  is partitioned into  $\mathcal D_1$  and  $\mathcal D_2$ 
  - i.e.,  $\mathcal{D}_1$  and  $\mathcal{D}_2$  are disjoint and  $\mathcal{D}=\mathcal{D}_1\cup\mathcal{D}_2$

RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT

- $\psi_i(v) = \psi(v)$  for every  $v \in S_i$  (i = 1, 2),
- $\psi$  preserves the labels and the edges

Proof: see [Hajiaghayi & Nishimura, 2007]

#### **Computing Characteristics: Join Nodes**







# Example



- all vertices in H and G have the same label (not denoted)
- edge labels are denoted by colors (i.e., there are 3 edge labels)





# Example (cont'd) – Nice Tree-Decomposition





RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT

universität**bon** 



# **Computing Characteristics**

problem: the number of iso-quadruples for separator and join nodes can be exponentially large

- Thm: For graphs of bounded tree-width and bounded degree, the set of z-characteristics can be computed in polynomial time for every node z
  - [Matoušek & Thomas, 1992]
- we cannot use this result
  - no additional assumption besides bounded tree-width







# **Equivalent Iso-Quadruples**

- $H_1, H_2, G$ : connected graphs of bounded treewidth
- TD(G): nice tree-decomposition of G computed in the preprocessing step
- z: node in TD(G)
- $\xi_1 = (S_1, \mathcal{D}_1, K_1, \psi_1) \in \Gamma(H_1, z); \xi_2 = (S_2, D_2, K_2, \psi_2) \in \Gamma(H_2, z)$
- **Def.:**  $\xi_1$  is equivalent to  $\xi_2$  if there is an isomorphism  $\pi$  between  $K_1$  and  $K_2$  such that
  - $\pi$  is a bijection between  $S_1$  and  $S_2$
  - $\psi_1(v) = \psi_2(\pi(v))$  for every  $v \in S_1$

Notation:  $\xi_1 \equiv \xi_2$ 



# **Equivalent Iso-Quadruples**

Prop. 1: Equivalence of iso-quadruples can be decided in polynomial time.

Proof: exercise

**Prop. 2:** If  $\xi_1, \xi_2$  in the above definition are equivalent then  $\xi_1 \in \Gamma_{ch}(H_1, z) \iff \xi_2 \in \Gamma_{ch}(H_2, z)$ 

Proof: exercise

⇒ it suffices to store only one representative z-characteristic for each equivalence class of the set of z-characteristics



RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT



## Non-Redundant Iso-Quadruples

- $\Rightarrow$  **given** a strong candidate pattern H generated by levelwise search and a node z in TD(G), it is **sufficient to check** only those iso-quadruples of  $\Gamma(H, z)$  for which there exists no equivalent z-characteristic for some frequent pattern listed before H
  - non-redundant iso-quadruples of H relative to z
  - $\Gamma_{nr}(H, z)$ : set of non-redundant iso-quadruples of H relative to z
  - $\Gamma_{nr,ch}(H,z)$ : set of non-redundant *z*-characteristics of *H*



### Non-Redundant Iso-Quadruples

**Prop.:** Let *H* be a strong candidate pattern, *G* be a transaction graph, both of bounded treewidth, *z* be a node in TD(G), and  $\xi \in \Gamma(H, z)$ . Then

$$\xi \in \Gamma_{\mathsf{ch}}(H, z) \iff \exists \xi' \in \bigcup_{P \in \mathcal{F}_H} \Gamma_{\mathsf{nr}, \mathsf{ch}}(P, z) \text{ such that } \xi' \equiv \xi$$

-  $\mathcal{F}_H$ : {*H*}  $\cup$  set of frequent patterns computed before *H* 

Proof: exercise





# Non-Redundant Iso-Quadruples

#### Theorem

(i) The number of non-redundant iso-quadruples of H relative to z is bounded by

## $O(|V(H)|^{k+1}).$

- instead of computing  $\Gamma_{nr}(H, z)$ , we will efficiently compute a superset  $\Gamma_{f}(H, z) \supseteq \Gamma_{nr}(H, z)$  of cardinality bounded by  $O(|V(H)|^{k+1})$
- ⇒ we must test only for polynomially many iso-quadruples whether they are characteristics
- (ii) It can be decided in time polynomial in the combined size of G and the set of frequent patterns listed before H whether an iso-quadruple in  $\Gamma_{f}(H, z)$  is a characteristic.





**Lemma:** Let *H* be a strong candidate pattern generated by levelwise search, *z* be a node in TD(G), and  $\xi = (S, \mathcal{D}, K, \psi) \in \Gamma_{nr}(H, z)$ . Then, for every vertex  $v \in V(H) \setminus V(K)$  it holds that

(i) the degree of v in H is at least 2 and

(ii) H has no cycle containing v



Proof: (blackboard)





**Lemma:** Let *H* be a strong candidate pattern generated by levelwise search, *S* be a subset of V(H) having at most k + 1 elements, and  $C_A$  be the maximal subset of the set of connected components of  $H[V(H) \setminus S]$ ) such that for every  $C \in C_A$ , all vertices of *C* satisfy (i) and (ii) in the lemma on the previous slide. Then

 $|\mathcal{C}_A| \leq k$ .







- feasible iso-quadruples: for a strong candidate pattern H generated by levelwise search and for a node z in TD(G), an iso-quadruple  $\xi = (S, \mathcal{D}, K, \psi) \in \Gamma(H, z)$  is called feasible if it satisfies the previous two conditions, i.e.,
  - (i) the degree of v in H is at least 2 and
  - (ii) H has no cycle containing v

for every  $v \in V(H) \setminus V(K)$ 

- set of feasible iso-quadruples in  $\Gamma(H, z)$  is denoted by  $\Gamma_{f}(H, z)$
- $\Rightarrow \Gamma_{\rm nr}(H,z) \subseteq \Gamma_{\rm f}(H,z)$





**Thm.:** Let *H* be a strong candidate pattern generated by levelwise search and let *z* be a node of TD(G) for some transaction graph *G*. Then

 $|\Gamma_{\mathsf{f}}(H,z)| \le O(|V(H)|^{k+1})$ 

**proof:** How many different  $\xi = (S, \mathcal{D}, K, \psi) \in \Gamma_{f}(H, z)$  can we have at most?

- S can be chosen in at most  $|V(H)|^{k+1}$  different ways,
- $\mathcal{D}$  can be chosen in at most  $2^{|\mathcal{C}_A|} \leq 2^k$  different ways (lemma on slide 61),
- there are at most (k+1)! injective functions from H[S] to G[bag(z)],
- $\Rightarrow$  we have

```
|\Gamma_{f}(H,z)| \le 2^{k} \cdot (k+1)! \cdot |V(H)^{k+1}|,
```

from which the claim directly follows, as k is assumed to be a constant.





# Claim (ii): Algorithm Computing Feasible Characteristics

#### Input:

- connected strong pattern graph H generated by levelwise search and transaction graph G, both of tree-width at most k for some constant k,
- a nice tree-decomposition TD(G) of G with nodes associated with

1

$$\bigcup_{P \in \mathcal{F}_H \setminus \{H\}} \Gamma_{\mathsf{f},\mathsf{ch}}(P,w)$$

for every node w in TD(G), and

• a node z in TD(G)

**Output:** set  $\Gamma_{f,ch}(H, z)$  of feasible *z*-characteristics

#### Algorithm: next slide





IAIS

# Claim (ii): Algorithm FeasibleCharacteristics



# I. Computing Feasible Iso-Quadruples (Line 2)

2:  $\Gamma_{f}(H, z) := FEASIBLEISOQUADRUPLES(H, bag(z))$ 

**Lemma:**  $\Gamma_{f}(H, z)$  can be computed in time  $O(|V(H)|^{k+1}|)$ 

**Proof:** By slides 62 and 63, the algorithm below is correct and computes  $\Gamma_{\rm f}(H,z)$  in time  $O(|V(H)|^{k+1}|)$ .

**Input:** connected graph H and bag(z) of a node z in TD(G)**Output:**  $\Gamma_{\rm f}(H,z)$ 

- 1:  $Y := \emptyset$
- 2: forall  $S \subseteq V(H)$  satisfying  $|S| \leq |\mathsf{bag}(z)|$  do
- 3: let C be the set of connected components of  $H[V(H) \setminus S]$
- 4: compute the subset  $C_A \subseteq C$  defined on slide 61
- 5: forall  $\mathcal{D}' \subseteq \mathcal{C}_A$  do
- $6: \qquad \mathcal{D} := \mathcal{C} \setminus \mathcal{D}'$
- 7: forall subgraph isomorphisms  $\psi: S \to bag(z)$  from H[S] to G[bag(z)] do
- 8: add  $(S, \mathcal{D}, H[S \cup V(\mathcal{D})], \psi)$  to Y
- 9: return Y

PhD Course, Szeged, 2012 - © T.Horváth





# II. Leaf Nodes (Line 5)

5: if  $\xi$  satisfies the condition of the Leaf Lemma then add  $\xi$  to  $\Gamma_{f,ch}(H, z)$ 

**condition** of the Leaf Lemma for  $\xi = (S, \mathcal{D}, K, \psi)$ :  $\mathcal{D} = \emptyset$  and  $\psi$  is a subgraph isomorphism from H[S] to G[bag(z)]

Lemma: The condition of the Leaf Lemma can be decided in constant time.

**Proof:**  $\psi$  is already a required subgraph isomorphism

• see lines 7-8 on the previous slide



#### III. Separator Nodes (Lines 7-9)

- // z is a separator node in TD(G) with child  $z^\prime$
- 7:  $\Gamma_{f,ch}(H, z') := \mathsf{FEASIBLECHARACTERISTICS}(H, G, TD(G), z')$
- 8: if  $\exists \xi' \in \mathfrak{S}(\xi)$  and  $\exists \xi'' \in \Gamma_{f,ch}(H, z')$  such that  $\xi' \equiv \xi''$  then
- 9: add  $\xi$  to  $\Gamma_{\rm f,ch}(H,z)$

let  $\xi \in \Gamma_{f}(H, z)$  be of the form  $\xi = (S, \mathcal{D}, K, \psi)$ 

 $\underbrace{\mathfrak{S}(\xi)}_{\text{of the Separator Lemma}, \text{ i.e.,}} \mathfrak{S}(\xi) \in \Gamma(H, z') \text{ satisfying the conditions}$ 

(S.a) 
$$S = \{v \in S' : \psi'(v) \in bag(z)\},\$$

(S.b)  $\mathcal{D}' = \{C : C \text{ is a connected component of } H[V(H) \setminus S'] \text{ and } C \text{ is a subgraph of } \mathcal{D}\},\$ 

(S.c)  $\psi(v) = \psi'(v)$  for every  $v \in S$ .





### III. Separator Nodes (Lines 7-9)

**Lemma:** Suppose that  $\Gamma_{f,ch}(P, z')$  has been computed correctly by the alg. on Slide 65 for every  $P \in \mathcal{F}_H$ . Then the alg. on Slide 65 computes  $\Gamma_{f,ch}(H, z)$ correctly in time polynomial in |V(H)|.

Proof (sketch): (correctness)

**Claim:**  $\mathfrak{S}(\xi) \subseteq \Gamma_{\mathfrak{f}}(H, z')$  for all  $\xi \in \Gamma_{\mathfrak{f}}(H, z)$ 

**Proof (sketch):** let  $\xi \in \Gamma_{f}(H, z)$  with  $\xi = (S, \mathcal{D}, K, \psi)$ 

- ⇒ *K* is a subgraph of *K'* for all  $\xi' = (S', \mathcal{D}', K', \psi') \in \mathfrak{S}(\xi)$  by (S.b) (see previous slide)
- $\Rightarrow$  all  $v \in V(H) \setminus V(K')$  satisfy the two conditions on slide 62

 $\Rightarrow$   $\xi'$  is feasible

⇒ correctness follows together with Prop. 2 (Slide 56) and the Separator lemma

PhD Course, Szeged, 2012 - © T.Horváth



#### III. Separator Nodes (Lines 7-9)

Proof (sketch) of the lemma on the previous slide: (efficiency)

- **1.** compute  $\mathfrak{S}(\xi)$ 
  - $\mathfrak{S}(\xi)$  has at most  $(k+1)! \cdot |V(H)|^{k+1}$  elements and it can be computed in time polynomial in |V(H)|
- **2.** for all  $\xi' \in \mathfrak{S}(\xi)$ , check whether  $\xi' \equiv \xi''$  for some  $\xi'' \in \Gamma_{f,ch}(H, z')$ 
  - $|\Gamma_{f,ch}(H,z')| \le |\Gamma_{f}(H,z')| = O(|V(H)|^{k+1})$  by the theorem on Slide 63
  - equivalence between iso-quadruples can be decided in polynomial time (Prop. 1 on Slide 56)
  - $\Rightarrow\,$  the condition above for  $\xi'$  can be decided in polynomial time



70

universita

### IV. Join Nodes (Lines 11-14)

- $/\!/ z$  is a join node in TD(G) with children  $z_1$  and  $z_2$
- 11:  $\Gamma_{f,ch}(H, z_1) := \mathsf{FEASIBLECHARACTERISTICS}(H, G, TD(G), z_1)$
- 12:  $\Gamma_{f,ch}(H, z_2) := \mathsf{FEASIBLECHARACTERISTICS}(H, G, TD(G), z_2)$
- 13: if  $\exists \xi_1 \in \bigcup_{P \in \mathcal{F}_H} \Gamma_{\mathsf{f},\mathsf{ch}}(P,z_1) \land \exists \xi_2 \in \bigcup_{P \in \mathcal{F}_H} \Gamma_{\mathsf{f},\mathsf{ch}}(P,z_2)$  such that  $\xi \equiv \bigoplus_{\xi}(\xi_1,\xi_2)$  then
- 14: add  $\xi$  to  $\Gamma_{\rm f,ch}(H,z)$
- let  $\xi = (S, \mathcal{D}, K, \psi) \in \Gamma_{f}(H, z)$  and  $\xi_{i} = (S_{i}, \mathcal{D}_{i}, K_{i}, \psi_{i}) \in \Gamma_{f}(H, z_{i})$  (i = 1, 2)//  $\psi$  is guaranteed a subgraph isom. from H[S] to G[bag(z)] by slide 66
- $\xi_1$  and  $\xi_2$  are **join consistent** with  $\xi$  if they satisfy the following conditions in the **Join Lemma**:
  - (J.a)  $S_i = \{v \in S : \psi(v) \in bag(z_i)\}$  for i = 1, 2,
  - (J.b) the connected components of  $\mathcal{D}$  are partitioned into  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , and
  - (J.c) for i = 1, 2,  $\psi_i(v) = \psi(v)$  for every  $v \in S_i$ .





#### IV. Join Nodes (Lines 11-14): The Join Operator

let z be a join node with children  $z_1$  and  $z_2$ , and let

• 
$$\xi_0 = (S_0, \mathcal{D}_0, K_0, \psi_0) \in \Gamma(H_0, z_0),$$

•  $\xi_1 = (S_1, \mathcal{D}_1, K_1, \psi_1) \in \Gamma(H_1, z_1)$  and  $\xi_2 = (S_2, \mathcal{D}_2, K_2, \psi_2) \in \Gamma(H_2, z_2)$ 

for some graphs  $H_0$ ,  $H_1$ , and  $H_1$ 

• assumption (w.l.o.g.):  $K_0, K_1$ , and  $K_2$  are pairwise vertex disjoint

the join of  $\xi_1$  and  $\xi_2$  w.r.t.  $\xi_0$ , denoted  $\bigoplus_{\xi_0} (\xi_1, \xi_2)$ , is a **quadruple**  $\xi = (S, \mathcal{D}, K, \psi)$  with

- $S = \{w_u : u \in \psi_1(S_1) \cup \psi_2(S_2)\}$  is a set of **new** vertices,
- $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2$ ,
- K is the graph ... (continued on the next slide)




### IV. Join Nodes (Lines 11-14): The Join Operator

- K is the graph obtained from the union of S and D by labeling all  $w_u \in S$  with the label of the corresponding vertex u and connecting
  - (i)  $w_{\psi_i(u)}$  and  $w_{\psi_i(v)}$  by an edge associated with the label of  $\{u, v\}$  for all  $\{u, v\} \in E(K_i)$  such that  $u, v \in S_i$  (i = 1, 2),
  - (ii)  $w_{\psi_i(u)}$  and v by an edge associated with the label of  $\{u, v\}$  for all  $\{u, v\} \in E(K_i)$  such that  $u \in S_i$  and  $v \in V(\mathcal{D}_i)$  (i = 1, 2), and
  - (iii)  $w_{\psi_1(u)}$  and  $w_{\psi_2(v)}$  by an edge associated with the label of  $\{u', v'\}$ for all  $\{u', v'\} \in E(K_0)$  such that  $\psi_0(u') = \psi_1(u)$ ,  $\psi_0(v') = \psi_2(v)$  for some  $u \in S_1$ ,  $v \in S_2$ , and  $\psi_1(u), \psi_2(v) \notin \psi_1(S_1) \cap \psi_2(S_2)$ ,
- $\psi: S \to X_z$  is defined by  $\psi(w_u) = u$  for every  $w_u \in S$ .





### IV. Join Nodes (Lines 11-14): The Join Operator

**Prop.:** Let *z* be a join node with children  $z_1$  and  $z_2$ , and let  $\xi_i = (S_i, \mathcal{D}_i, K_i, \psi_i) \in \Gamma_{ch}(H_i, z_i)$  for some graph  $H_i$  (i = 0, 1, 2) such that  $z_0 = z$ , and  $K_0, K_1$ , and  $K_2$  are pairwise vertex disjoint. Then  $\bigoplus_{\xi_0} (\xi_1, \xi_2) = (S, \mathcal{D}, K, \psi)$  is uniquely defined and it is an element of  $\Gamma(P, z)$  for some pattern *P*.

#### Proof (sketch):

- uniqueness is straightforward
- any graph P satisfying K = P[V(K)] is appropriate
- $|S| \le k+1$
- since the  $\psi_i$  's are all injective,  $\psi$  is also injective







### IV. Join Nodes (Lines 11-14)

- **Lemma:** Let z be a join node in TD(G) with children  $z_1$  and  $z_2$ , and suppose that  $\Gamma_{f,ch}(P, z_i)$  has been computed correctly by Alg. FEASIBLECHARAC-TERISTICS for every  $P \in \mathcal{F}_H$  and i = 1, 2. Then Alg. FEASIBLECHARACTERISTICS computes  $\Gamma_{f,ch}(H, z)$  correctly in time polynomial in the combined size of H and the set of frequent patterns computed before H.
- **Proof (sketch):** For the correctness, let  $\xi = (S, \mathcal{D}, K, \psi) \in \Gamma_{f}(H, z)$ . Then, as  $\psi$  is a subgraph isomorphism from H[S] to G[bag(z)] by the alg. on slide 66, we have

$$\begin{split} \xi \in \Gamma_{\mathsf{f},\mathsf{ch}}\left(H,z\right) \\ \iff & \exists \xi_i \in \Gamma_{\mathsf{ch}}(H,z_i) \text{ for } i=1,2 \text{ such that } \xi_1,\xi_2, \text{ and } \xi \\ & \text{ satisfy the conditions of the Join Lemma} \\ \iff & \exists \xi_i \in \bigcup_{P \in \mathcal{F}_H} \Gamma_{\mathsf{f},\mathsf{ch}}(P,z_i) \text{ for } i=1,2 \text{ with } \xi \equiv \oplus_{\xi}(\xi_1,\xi_2) \end{split}$$



### IV. Join Nodes (Lines 11-14)

#### Proof (sketch) of the lemma on the previous slide: (efficiency)

- 1. the join operator can be computed in polynomial time
- **2.** the number of pairs  $\xi_1$ ,  $\xi_2$  to be considered is bounded by

 $\mathsf{poly}(\mathsf{size}(H),\mathsf{size}(\mathcal{F}_H))$ 

- the last equivalence on the previous slides implies quadratic time algorithm in  $\text{size}(\mathcal{F}_H)$
- in fact, it can be done in time **linear** in  $size(\mathcal{F}_H)$
- equivalence of iso-quadruples can be decided in polynomial time (Prop. 1 on Slide 56)





# **Putting Together**

- **Thm:** The algorithm on the next two slides lists frequent connected subgraphs in incremental polynomial time.
- **Proof:** Using the previous results, it follows by induction on the depth of the treedecomposition of the transaction graph.



universitatbo



## The Mining Algorithm

**Input:** database DB of graphs of tree-width at most k and integer t > 0**Output:** all *t*-frequent connected subgraphs

- 1: // (preprocessing of the transaction graphs)
- 2: forall G in DB do
- 3: compute a nice tree-decomposition TD(G) of G
- 4: forall node z in TD(G) do  $\Sigma_{G,z} := \{(\emptyset, \emptyset, \emptyset, \emptyset)\}$
- 5: // (computing frequent subgraphs consisting of a single vertex)
- 6:  $S_0 = \emptyset$
- 7: forall graphs H consisting of a single labeled vertex do  $PROCESS(H, S_0)$
- 8: // (computing frequent subgraphs consisting of at least one edge)

9: for 
$$(l := 0; S_l \neq \emptyset; l := l + 1)$$
 do

- $10: \quad C_{l+1} := S_{l+1} := \emptyset$
- 11: forall  $P \in S_l$  do
- 12: forall  $H \in \rho(P)$  satisfying tree-width $(H) \leq k \wedge H \notin C_{l+1} \wedge \rho^{-1}(H) \subseteq S_l$  do
- 13: add H to  $C_{l+1}$
- 14:  $\mathsf{PROCESS}(H, S_{l+1})$

// next slide



## Function Process (Lines 7 and 14)

#### function PROCESS(H, S)

- 1: counter := 0
- 2: forall G in DB do
- 3: unmark G
- 4: r := root of TD(G)
- $\Gamma_{f,ch}(H,r) := FEASIBLECHARACTERISTICS(H, G, TD(G), r)$ 5: // (Slide 66)
- 6: if  $\exists (S, \mathcal{D}, H, \psi) \in \Gamma(H, r)$  equivalent to some  $\xi \in \Sigma_{G, r} \cup \Gamma_{\mathsf{f.ch}}(H, r)$  then // (H is subgraph isomorphic to G)
- 7: counter := counter + 1
- 8: mark G
- 9: if counter > t then
- 10: **print** H and add it to S
- forall G in DB such that G is marked do 11:
- forall node z in TD(G) do  $\Sigma_{G,z} := \Sigma_{G,z} \cup \Gamma_{f,ch}(H,z)$ 12:



universitatbo



// (H is frequent)

### Example

#### mining problem:

list all 1-frequent connected subgraphs of the database consisting of the single transaction graph *G*:

- i.e., all subtrees
- all vertices in *H* and *G* have the same label (not denoted)
- edge labels are denoted by colors (i.e., there are 3 edge labels)
  - see also the previous example







# Example (cont'd)

Steps 1 - 4 of the alg. on slide 78:

- compute nice tree-decomposition of G
- assign the empty iso-quadruple to each node





universitätbor

## **Example: Feasible Characteristics**

Steps 6-7 of the alg. on Slide 78

pattern  $H_0: \bullet_u$ 

 $(S, \mathcal{D}, K)$ -triples for possible feasible iso-quadruples:

•  $(\emptyset, \emptyset, \emptyset), (\emptyset, H_0, H_0)$ 

a

b

C

е

•  $(\{u\}, \emptyset, H_0)$ 

transaction

graph G:



d



 $(\{x\}, y, H_1, x \mapsto b)$ 

 $(\{y\}, x, H_1, y \mapsto b)$ 

 $(\{u\}, \emptyset, H_0, u \mapsto b)$ 

 $(\emptyset, H_0, H_0, \emptyset)$ 

 $(\emptyset, H_0, H_0, \emptyset)$ 

 $(\emptyset, H_0, H_0, \emptyset)$ 

 $(\emptyset, \emptyset, \emptyset, \emptyset)$ 

 $(\emptyset, \emptyset, \emptyset, \emptyset)$ 

 $(\emptyset, \emptyset, \emptyset, \emptyset)$ 

 $(\{u\}, \emptyset, H_0, u \mapsto c)$ 

 $(\{u\}, \emptyset, H_0, u \mapsto c)$ 

 $(\{u\}, \emptyset, H_0, u \mapsto d)$ 

 $\{\{u\}, \emptyset, H_0, u \mapsto c\}$ 

IAIS

83

 $(\emptyset, H_0, H_0, \emptyset)$ 

 $(\emptyset, \emptyset, \emptyset, \emptyset)$ 

 $(\{x,y\}, \emptyset, H_1, x \mapsto b, y \mapsto c)$ 

 $(\{x, y\}, \emptyset, H_1, x \mapsto c, y \mapsto b)$ 

 $(\emptyset, H_0, H_0, \emptyset)$ 

(0, 0, 0, 0)

 $(\{u\}, \emptyset, H_0, u \mapsto b)$ 

 $(\{u\}, \emptyset, H_0, u \mapsto c)$ 

## Example (cont'd)

 $H_1 \in S_1 \iff (\{y\}, \{x\}, H_1, y \mapsto b) \\ (\{y\}, \{x\}, H_1, y \mapsto b) \\ (\{u\}, \{u\}, H_0, u \mapsto a) \\ (\{u\}, \emptyset, H_0, u \mapsto a) \\ (\{u\}, u \mapsto a)$ Step 12 of the alg. on slide 78  $(\{u\}, \emptyset, H_0, u \mapsto b)$  $\rho(\bullet) = \{\bullet - \bullet, \bullet - \bullet, \bullet - \bullet\}$  $(\emptyset, \emptyset, \emptyset, \emptyset)$ **凡** {a,b} pattern  $H_1$ : • • •  $(\{u\}, \emptyset, H_0, u \mapsto a)$  $(S, \mathcal{D}, K)$ -triples for possible {a} {b}  $(\emptyset, \emptyset, \emptyset, \emptyset)$ feasible iso-quadruples: {b,c} •  $(\emptyset, \emptyset, \emptyset)$ ,  $(\emptyset, H_1, H_1)$ •  $(\{x\}, \emptyset, x), (\{x\}, y, H_1)$  $(\{u\}, \emptyset, H_0, u \mapsto b)$ (0, 0, 0, 0){b} {c} •  $(\{y\}, \emptyset, y), (\{y\}, x, H_1)$ •  $(\{x, y\}, \emptyset, H_1)$ {c,d} а  $(\{u\}, \emptyset, H_0, u \mapsto d)$  $(\emptyset, \emptyset, \emptyset, \emptyset)$ {d} {c} transaction b graph G: {c,e} С  $(\{u\}, \emptyset, H_0, u \mapsto c)$  $\{u\}, \emptyset, H_0, u \mapsto e$  $(\emptyset, \emptyset, \emptyset, \emptyset)$ е PhD Course, Szeged, 2012 - © T.Horváth

universitätbor RHEINISCHE FRIEDRICH-WILHELMS-UNIVERSITÄT

 $(\{x\},\{y\},H_1,x\mapsto b)$ 

Frequent Connected Subgraph Mining



# Example (cont'd)

and so on...

- **notice** that the algorithm could further be improved because there are redundant characteristics
  - because feasible iso-quadruples are processed
  - the number of redundant characteristics is polynomial in the combined size of the input and the set of previously generated frequent pattern
  - using some advanced data structure, redundant characteristics can be removed in time polynomial in the input







## Summary

- efficient pattern mining is possible even for computationally hard matching operators
- the technique might be of some independent interest and useful to design efficient algorithms if straightforward dynamic programming requires exponential space
- the positive theoretical result of this lecture is not always practical
  - e.g., for k > 4 (or 5?), no practical algorithm is known for deciding whether a graph has tree-width at most k
  - for k < 4: fast algorithm [Arnborg, Corneil, Proskurowski, 1987]
    - chemical graphs of pharmacological compounds have mostly tree-width at most 3

**open problem:** Is it possible to mine frequent connected subgraphs in graphs of bounded tree-width with **polynomial delay**?



