

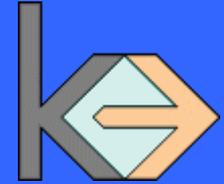
Discovering Frequent Substructures in Large Unordered Trees

Unot

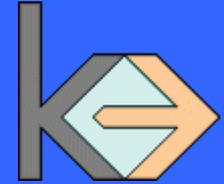
WS08/09

Seminar aus Maschinellem Lernen

Prof. J. Fürnkranz

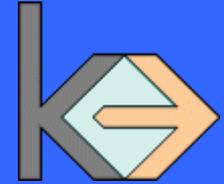


- Motivation
- Basic Definitions
- Canonical Representation
- Algorithm Unot
 - Overview
 - Enumerating Pattern
 - Compute Occurrence
- Conclusions



- Graphische Daten sind sehr üblich in vielen Bereichen:
 - Chemie, Biologie, Militär...
 - Mit Graphen lassen sich wesentlich komplexere Strukturen darstellen
 - Sind aber auch komplizierter zu handhaben als einfache Mengen
 - Graph mining: discovering patterns in large collections of graph or tree structures
 - Anwendung:
 - In der Pharmaforschung werden gemeinsame Fragmente einer Menge von Molekülen gesucht
- => ein Algorithmus für gelabelt ungeordnete Struktur

Motivation



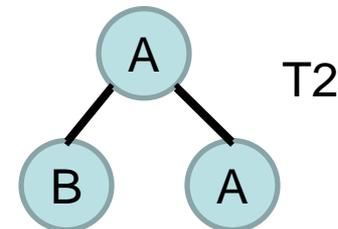
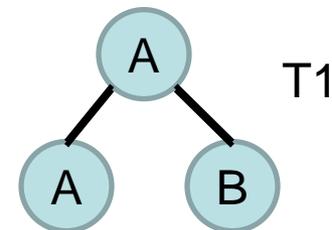
- Wie sieht es aus, wenn Algorithmus z.B. wie Freqt für ungeordnete Bäume angewendet?

⇒ Problem: Es kommt Isomorphie vor

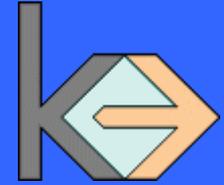
⇒ Lösung: Unot, uFreq, ...

⇒ Wie?

⇒ durch Canonical Form



Definition(1): Basics



- **Semi-strukturierte Daten:**

- Hier gelabelt ungeordnet Bäume als ein Modell von semi-strukturierte Daten und Pattern.

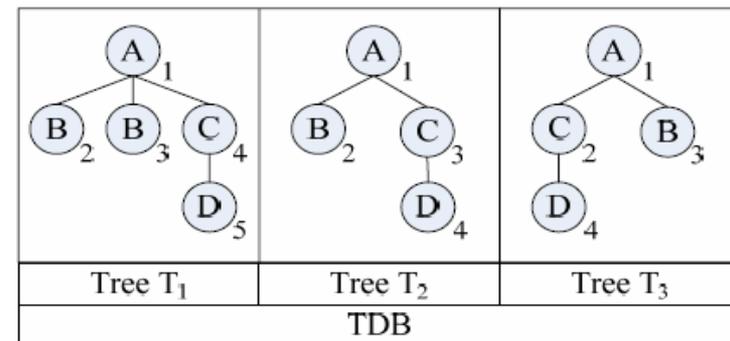
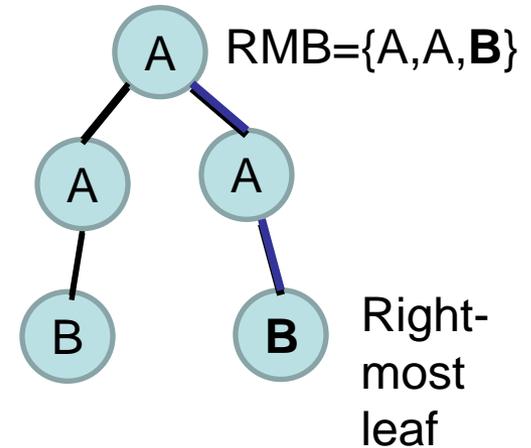
- **Labeled unordered Tree:** $T=(V,E, r,label)$

- **Labeled ordered Tree:** $T=(V,E,B, r,label)$

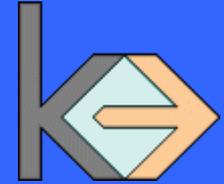
- B: Menge of Geschwester Relation von links zu rechts
- Z.B: $(v1, v2)$, $v1$ is the linke Geschwester von $v2$

- **RMB**(Right most Branch): Pfad von Wurzel zu dem meisten rechten Baumblatt (right most leaf)

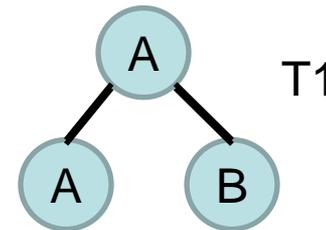
- **Unordered database** : is endliche Menge $\mathcal{D}=(D_1...D_m)$ von (geordnete) Bäume
 - Wobei D_i is data tree



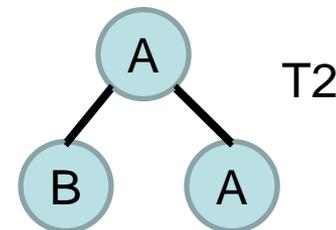
Canonical Form (1)



- gelabelten geordnete Bäume als **Representation** von ungeordneten Bäume
- zwei geordnete Bäume T_1, T_2 äquivalent, $T_1 \equiv T_2$, wenn sie representieren gleiche ungeordnete Bäume \Rightarrow **Isomorphie**
- **Depth-label sequence $C(T)$**
 $C(T) = ((\text{dep}(v_1), \text{label}(v_1)), \dots, (\text{dep}(v_k), \text{label}(v_k)))$, wobei jeder Komponent (d, l) als „depth-label Paar“

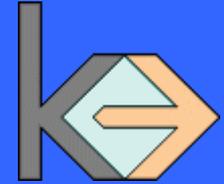


$$C(T_1) = ((0,A), (1,A), (1,B))$$

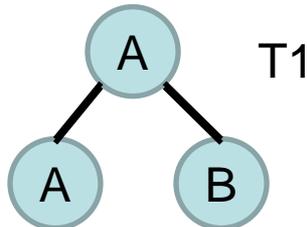


$$C(T_2) = ((0,A), (1,B), (1,A))$$

Canonical Form(2)



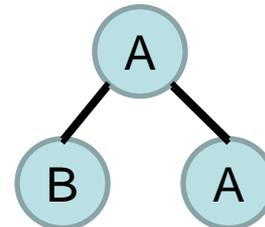
- **Vergleichen** von zwei depth-label Paaren,
 - **$(d1, l1) > (d2, l2)$**
 - wenn $d1 > d2$ oder
 - $d1 = d2$ and $l1 > l2$
- Beispiel:
 $(1, A) > (0, A)$, wegen $d1 = 1 > 0 = d2$
 $(1, A) > (1, B)$, wegen $l1 = A > B = l2$
und $d1 = d2$
- „**Heavy than**“: **$C(T1) \geq_{lex} C(T2)$** , wenn **i) oder ii)**
 - **i)** Vergleich mit depth-label Paar



T1

\geq_{lex}
x

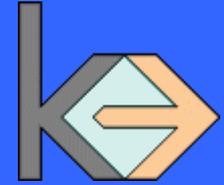
T2



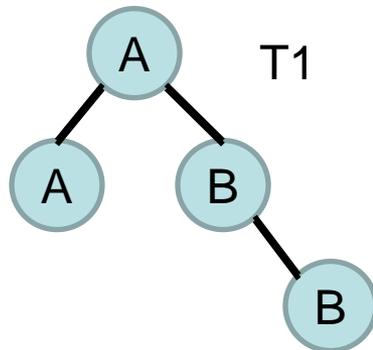
$C(T1) = ((0, A), (1, A), (1, B))$

$C(T2) = ((0, A), (1, B), (1, A))$

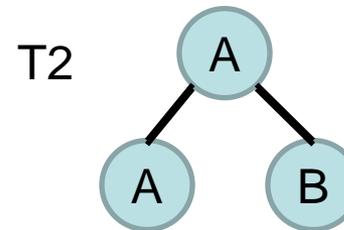
Canonical Form(3)



- ii) durch Längen der $C(T)$: T_2 ist Prefix von T_1



\geq_{le}
 x

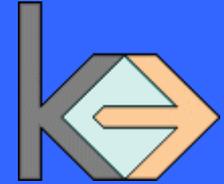


$C(T_1) = ((0,A), (1,A), (1,B), (2,B))$

$C(T_2) = ((0,A), (1,B), (1,A))$

- Ein Baum T ist in **Canonical Form**, wenn $C(T)$ ist **heaviest** von allen zu T äquivalenten Bäumen

Canonical Form(4)

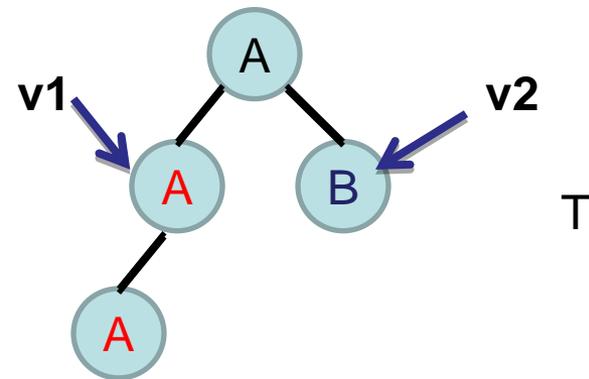


- Jeder kanonische Baum T muss „left-heavy“ halten.
- **left heavy**: für beliebige Geschwister Knoten v_1, v_2 impliziert

$$C(T(v_1)) \geq_{\text{lex}} C(T(v_2))$$

=> Kriterien für Rightmost Expansion

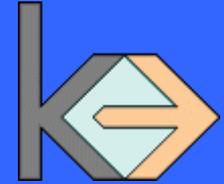
$$C(T) = ((0,A), (1,A), (1,B))$$



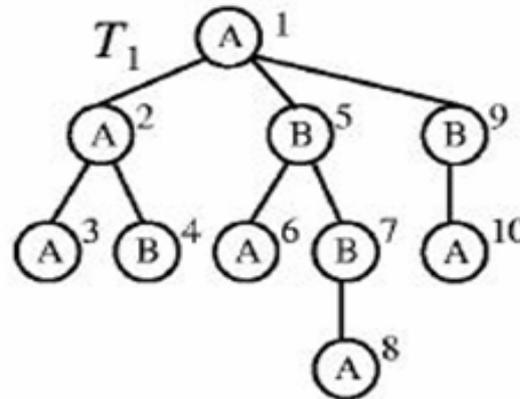
$$C(T(v_1)) = ((0,A), (1A))$$
$$\geq_{\text{lex}}$$

$$C(T(v_2)) = ((0,B))$$

Beispiel: Left-heavy condition

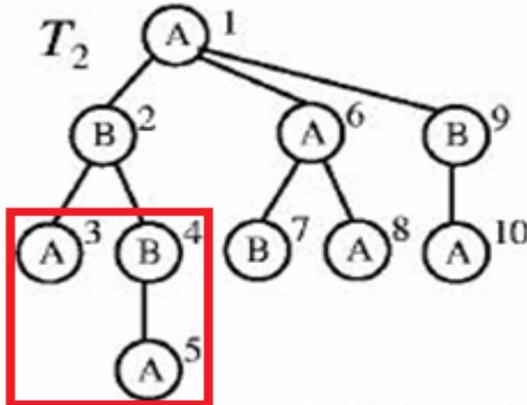


T1, T2, T3 sind äquivalent
Annahme: A>B>C

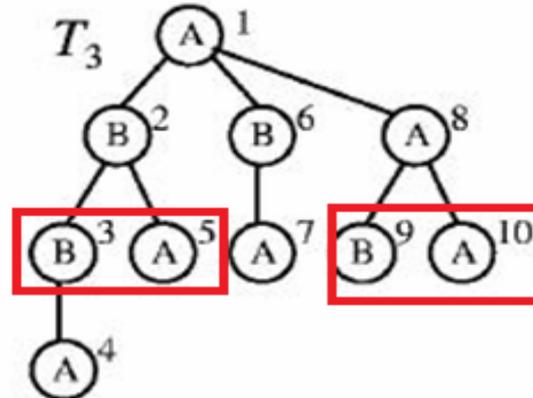


T1 ist left heavy
und kanonisch

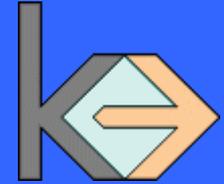
(0,A,1,A,2,A,2,B,1,B,2,A,2,B,3,A,1,B,2,A)



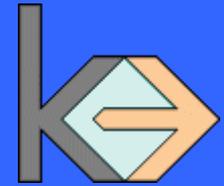
(0,A,1,B,2,A,2,B,3,A,1,A,2,B,2,A,1,B,2,A)



(0,A,1,B,2,B,3,A,2,A,1,B,2,A,1,A,2,B,2,A)



- Übersicht:
- Unot: findet alle kanonische Representation für frequent ungeordneten Trees in a vorgegebenem Datenbank \mathcal{D} von Data Trees
- 2 steps:
- i) Enumeration von allen Subtree als Canonical Representation
- => FindAllChildren
- li) Computation von ihre Occurrences
- => UpdateOcc



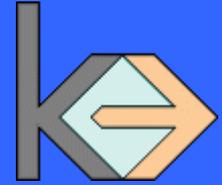
Algorithm Unot($\mathcal{D}, \mathcal{L}, \sigma$)

Input: the database $\mathcal{D} = \{D_1, \dots, D_m\}$ ($m \geq 0$) of labeled unordered trees, a set \mathcal{L} of labels, and the minimum frequency threshold $0 \leq \sigma \leq 1$.

Output: the set $\mathcal{F} \subseteq \mathcal{C}$ of all frequent unordered trees of size at most N_{\max} .

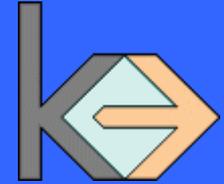
Method:

1. $\mathcal{F} := \emptyset$; $\alpha := \lceil |\mathcal{D}| \sigma \rceil$; //Initialization
 2. For any label $\ell \in \mathcal{L}$, do:
 $T_\ell := (0, \ell)$; /* 1-pattern with copy depth 0 */
 Expand($T_\ell, \mathcal{O}, 0, \alpha, \mathcal{F}$);
 3. Return \mathcal{F} ; //The set of frequent patterns
-



- Input
 - Datenbank von Data Trees
 - Menge von Labels
 - Minimum support threshold σ
- Output
 - Frequent Subtrees die support threshold erreichen
- Initialisierung

Rightmost Expansion: Expand



Procedure $\text{Expand}(S, \mathcal{O}, c, \alpha, \mathcal{F})$

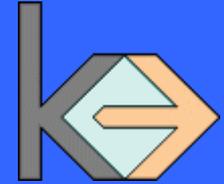
Input: A canonical representation $S \in \mathcal{U}$, the embedding occurrences $\mathcal{O} = EO^{\mathcal{D}}(S)$, and the copy-depth c , nonnegative integer α , and the set \mathcal{F} of the frequent patterns.

Method:

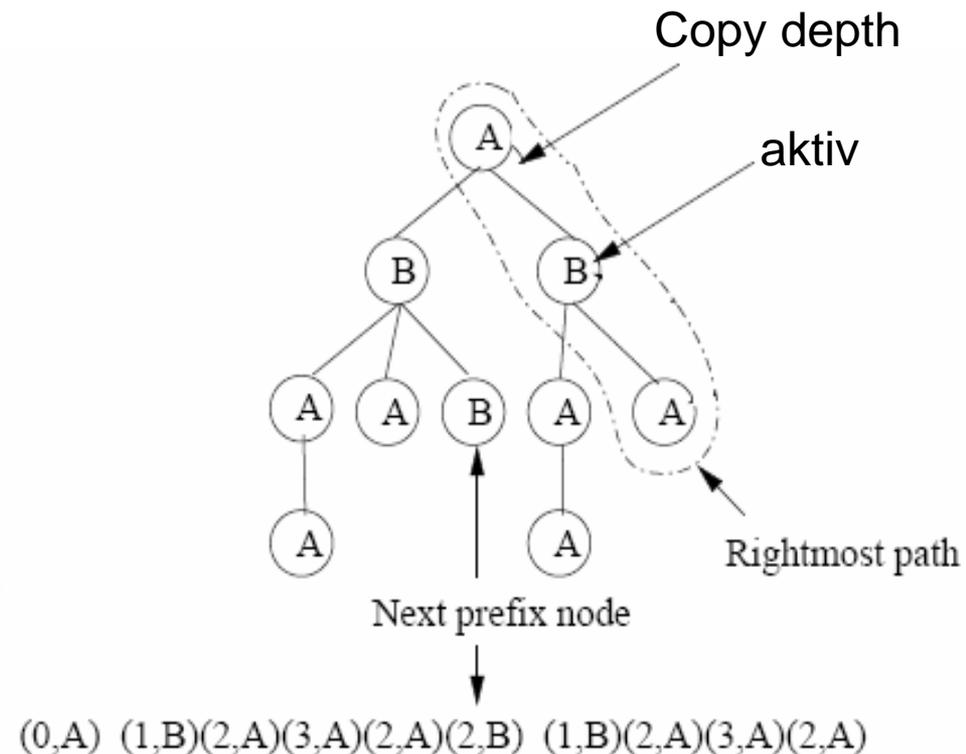
- If $(|\mathcal{O}| < \alpha)$ then return;
- Else $\mathcal{F} := \mathcal{F} \cup \{S\}$;
- For each $\langle S \cdot (i, \ell), c_{\text{new}} \rangle \in \text{FindAllChildren}(S, c)$, do;
 - $T := S \cdot (i, \ell)$;
 - $\mathcal{P} := \text{UpdateOcc}(T, \mathcal{O}, (i, \ell))$;
 - $\text{Expand}(T, \mathcal{P}, c_{\text{new}}, \alpha, \mathcal{F})$;

-
- 1) FindAllChildren abrufen um alle Subtrees zu finden
 - 2) UpdateOcc abrufen für jeden Subtree, um jeweilige Occurence zu rechnen
 - 3) test, ob der gefundene Subtree frequent ist
 - 4) neue frequente Subtree in Output Mengen hinzufügen

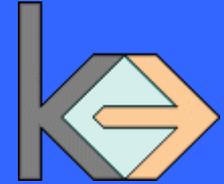
Copy Depth



- Interne Knoten am RMB r_i ist **aktiv**, wenn $C(R_i) = C(L_i)$
- Copy Depth: höchste aktive Knoten
- R_{ml} als spezielle Copy Depth wenn es keine andere gibt.



Enumeration: FindAllchildren



Procedure FindAllChildren(T, c) :

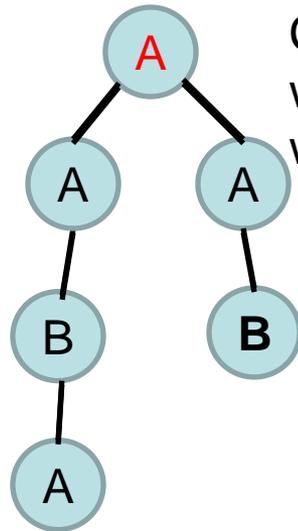
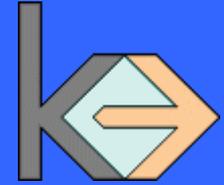
Method : Return the set *Succ* of all pairs $\langle S, c \rangle$, where S is the canonical child tree of T and c is its copy depth generated by the following cases:

Case I : If $C(L_k) = C(R_k)$ for the copy depth k :

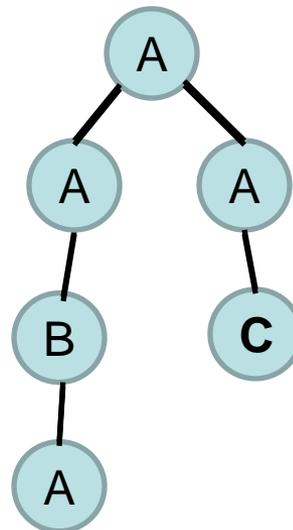
- The canonical child trees of T are $T \cdot (1, \ell_1), \dots, T \cdot (k+1, \ell_{k+1})$, where $\text{label}(r_i) \geq \ell_i$ for every $i = 1, \dots, k+1$. The trees $T \cdot (k+2, \ell_{k+2}), \dots, T \cdot (g+1, \ell_{g+1})$ are not canonical.
- The copy depth of $T \cdot (i, \ell_i)$ is $i-1$ if $\text{label}(r_i) = \ell_i$ and i otherwise for every $i = 1, \dots, k+1$.

Case II : If $C(L_k) \neq C(R_k)$ for the copy depth k :

- Let $m = |C(R_k)| + 1$ and $w = (d, \ell)$ be the m -th component of $C(L_k)$ (the next position to be copied). The canonical child trees of T are $T \cdot (1, \ell_1), \dots, T \cdot (d, \ell_d)$, where $\text{label}(r_i) \geq \ell_i$ for every $i = 1, \dots, d-1$ and $\ell \geq \ell_d$ holds.
- The copy depth of $T \cdot (i, \ell_i)$ is $i-1$ if $\text{label}(r_i) = \ell_i$ and i otherwise for every $i = 1, \dots, d-1$. The copy depth of $T \cdot (d, \ell_d)$ is k if $w = v$ and d otherwise.

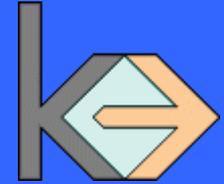


Copy depth bleibt bei (0, A),
wenn rechte Subtree kopiert
weiter



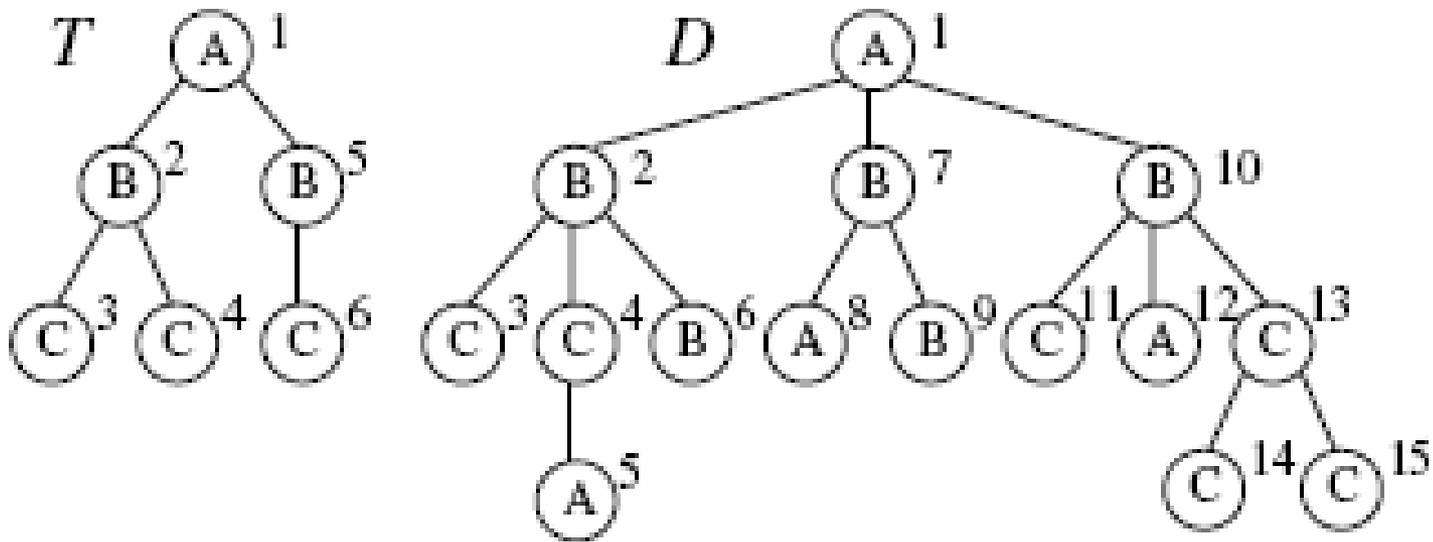
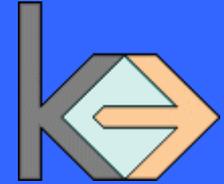
Kopiert nicht
mehr weiter,
copy depth
ändert

Definition(2): Occurrences

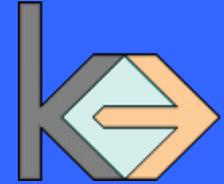


- Occurrence: T occurs in D
 - Es existiert mapping $\phi : V_D \dashrightarrow V_T$
 - die Eltern Relation beibehalten
 - Und Labels beibehalten
- Total occurrence: is the k Tuple
 $Toc(\phi) = \langle \phi(1), \dots, \phi(k) \rangle \in (V_D)^k$
- Embedding occurrence: is the set
 $Eoc(\phi) = \{ \phi(1), \dots, \phi(k) \}$
gehört zu V_D
- Variant: total, embedding, root, document.
- Root occurrence: $Roc(\phi) = \phi(1) \in V_D$

Beispiel: Occurrences



Updateocc (1)



Algorithm UpdateOcc(T, \mathcal{O}, d, ℓ)

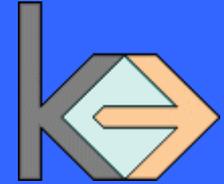
Input: the rightmost expansion of a pattern S , the embedding occurrence list $\mathcal{O} = EO^{\mathcal{D}}(S)$, the depth $d \geq 1$ and a label $\ell \in \mathcal{L}$ of the rightmost leaf of T .

Output: the new list $\mathcal{P} = EO^{\mathcal{D}}(T)$.

Method:

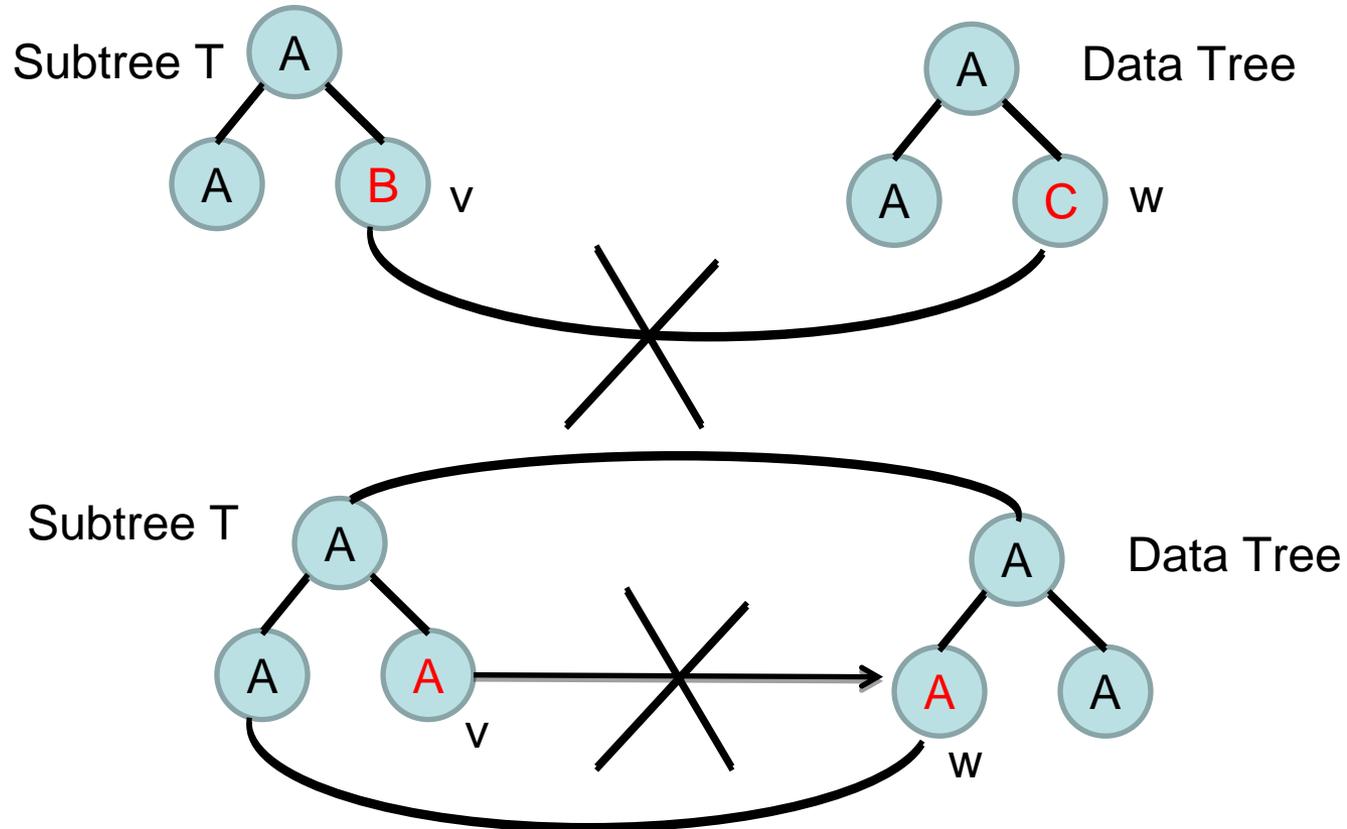
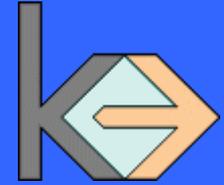
- $\mathcal{P} := \emptyset$;
- For each $\varphi \in \mathcal{O}$, do:
 - + $x := \varphi(r_{d-1})$; /* the image of the parent of the new node $r_d = (d, \ell)$ */
 - + For each child y of x do: **Bedingung 3)**
 - If $label_D(y) = \ell$ and $y \notin E(\varphi)$ then **Bedingung 1), 2)**
 - $\xi := \varphi \cdot y$ and $flag := true$;
 - Else, skip the rest and continue the for-loop;
 - For each $i = 1, \dots, d - 1$, do:
 - If $C(L_i) = C(R_i)$ but $\xi(left_i) \neq \xi(right_i)$ then **Bedingung 4)**
 - $flag := false$, and then break the inner for-loop;
 - If $flag = true$ then $\mathcal{P} = \mathcal{P} \cup \{\xi\}$;
- Return \mathcal{P} ;

Updateocc (2)

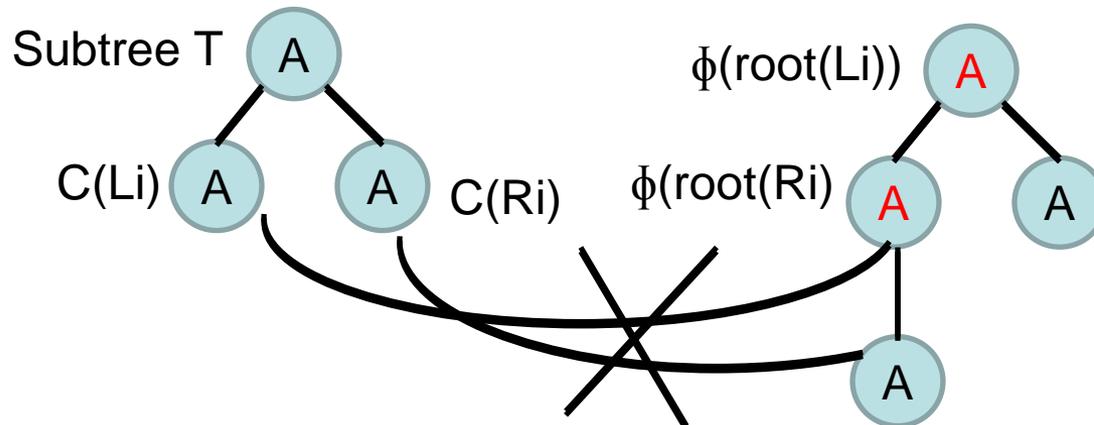
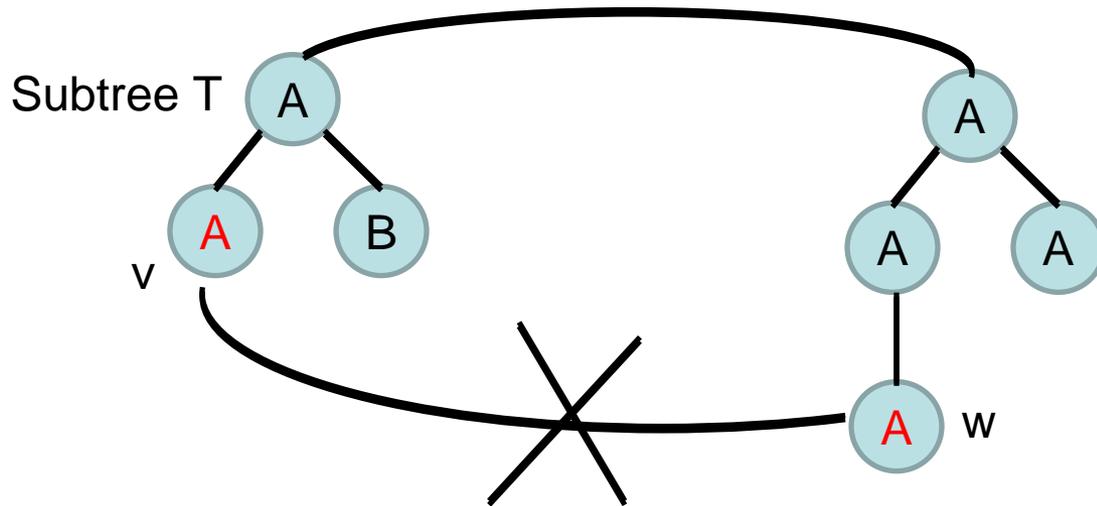
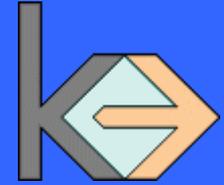


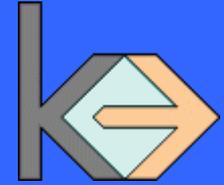
- Ein Mapping $\phi = \phi \bullet w$ ist ein kanonische total Occurrence von Subtree T in \mathcal{D} , wobei ϕ ist Occurrence Liste der Vorfahren von Subtree T , wenn es folgende Bedingungen erfüllt:
 - 1) $\text{Label}_T(v) = \text{Label}_{\mathcal{D}}(w)$
 - 2) w trat noch nicht in der Occurrence Liste der Vorfahren von Subtree T
 - 3) wenn v ist ein Kinder von r_i , muss w auch ein Kinder von $\phi(r_i)$
 - 4) wenn $C(Li) = R(i)$, impliziert $\phi(\text{root}(Li)) = \phi(\text{root}(Ri))$
 - Nachbildung von Li sowie Ri in \mathcal{D} haben auch gleichen Elternknoten.

Beispiel

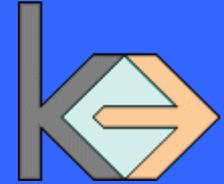


Beispiel

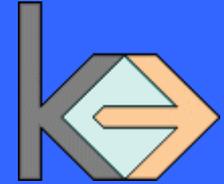




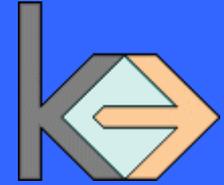
- Konstante Zeitkomplexität $O(1)$ für die Enumeration jedes Subtree
- Enumeration für jedes frequent pattern T in $O(kb^2m)$, wobei k und b jeweils Größe und Branching Fator von T , sowie m ist total occurrences of T in Data Trees



- Depth-first Traverse
- Enumeration based on Canonical Representation => Isomorphie Problem vermeiden (effizienter, schneller)
- Nutzen von Rightmost Expansion
- Nutzen von Occurrence Liste zum Rechnen von Menge von frequent canonical Subtrees

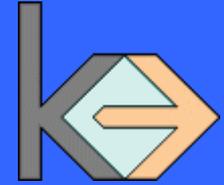


- [1] Asai, T., Abe, K., Kawasoe, S., Arimura, H.: Efficient Substructure Discovery from Large Semi-Structured Data, 2002.
- [2] K. Abe, S. Kawasoe, T. Asai, H. Arimura, and S. Arikawa. Optimized Substructure Discovery for Semi-structured Data, 2002.
- [3] Asai, T., Abe, K., Kawasoe, S., Arimura, H.: Efficient Substructure Discovery from Large Semi-structured Data
- [4] S. Nakano, T. Uno, Efficient Generation of Rooted Trees, *NII Technical Report, 2002*



Thank you ! 😊

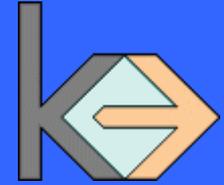
Comparison uFreq and Unot



- same:
- Rightmost expansion
- Use occurrence list

- Difference:
- Freq: rekord occurrence list of every pattern (more memory cost, more complex occurrence list)
- Unot: rekord occurrence list of every canonical pattern

Basic definition



- Pre-order-string: same to depth label sequece