



Implementace os XPath ve vícerozměrném přístupu pro indexování XML dat

Michal Krátký¹, michal.kratky@vsb.cz
Jaroslav Pokorný², pokorny@ksi.ms.mff.cuni.cz
Tomáš Skopal¹, tomas.skopal@vsb.cz
Václav Snášel¹, vaclav.snasel@vsb.cz



¹Katedra informatiky, VŠB – Technická univerzita Ostrava

²Katedra softwarového inženýrství, Univerzita Karlova v Praze

Obsah

- Úvod – XML, dotazovací jazyky.
- Aktuální stav v oblasti indexování XML dat.
- Vícerozměrný přístup pro indexování XML dat.
- Dotazování, implementace os XPath, implementace podmnožin dalších jazyků.
- Datové struktury.
- Výsledky experimentů.
- Závěr.

Úvod

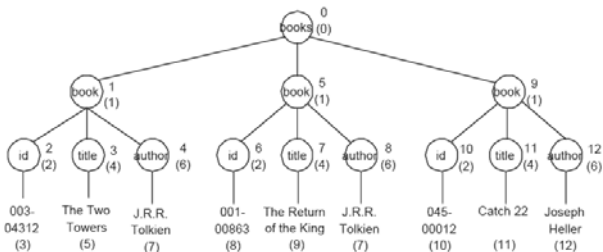
- Extensible Markup Language (XML) je značkovací jazyk vyvíjený W3C.
- “databázový pohled”: XML je jazyk pro modelování dat.
- Dokument(y) tvoří databázi, DTD (XML Schema) jsou jejími schémata.
- XML dotazovací jazyky (XPath, XQL, XQuery, ...).
- Současné přístupy (relační, objektově-relační) nejsou vhodné pro indexování XML dokumentů.
- Problémem je nutnost průchodu stromem při provádění XML dotazu.

XML dotazovací jazyky – osy XPath

- XML strom.
- 13 os XPath - relace mezi uzly dokumentu.

Osa v/a	Výsledek		
child	dítě uzlu v	following	následníci
descendant	potomek v	following-sibling	následující sourozenci
descendant-or-self	descendant + v	preceding	předchůdci
parent	rodič v	preceding-sibling	předcházející sourozenci
ancestor	předek v	attribute	atributy v
ancestor-or-self	předek + v	self	uzel v
		namespace	

Osy XPath



Aktuální stav v indexování XML dat

- API: SAX, DOM - není indexování.
- Přístupy založené na relační dekompozici.
- Trie reprezentace dokumentu.
- Vícerozměrné přístupy.
- Další: signaturové přístupy apod.

Přístupy založené na relační dekompozici

- Často velmi triviální mapování XML dokumentů do relačních tabulek (Oracle apod.).
- Neposkytuje „opravdové“ indexování XML – nelze efektivně implementovat XML dotazovací jazyky.
- Sofistikovanější přístupy: **STORED**, **XISS**, **Hybrid Storage Model**.

7/36

Vícerozměrné přístupy

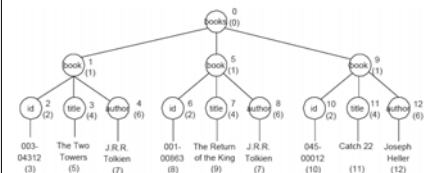
- prof. Bayer** – Indexování XML jako vícerozměrný problém. Velmi problematické indexování cest v jedné souřadnici.
- XPath Accelerator**, pro každý element je vytvořen 5-rozměrný vektor a osy XPath jsou implementovány rozsahovými dotazy.

8/36

Vícerozměrný přístup pro indexování XML dat - model

```
<!DOCTYPE books [
  <ELEMENT book*(book)>
  <ELEMENT book(title,author)>
  <!ATTLIST book id CDATA #REQUIRED>
  <ELEMENT title(#PCDATA)>
  <ELEMENT author(#PCDATA)>
]>
```

```
<?xml version="1.0" ?>
<books>
  <book id="003-04312">
    <title>The Two Towers</title>
    <author>J.R.R. Tolkien</author>
  </book>
  <book id="001-00863">
    <title>The Two Towers</title>
    <author>J.R.R. Tolkien</author>
  </book>
  <book id="045-00012">
    <title>Catch 22</title>
    <author>Joseph Heller</author>
  </book>
</books>
```



Graf je množina cest.

9/36

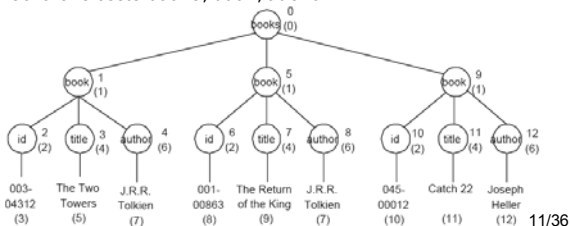
Cesty, značkové cesty

- $id_u(u_i)$ – jedinečné číslo uzlu u_i
- cesta**: $p = id_u(u_0), id_u(u_1), \dots, id_u(u_{l_p})$, $s \in X_p$,
- značková cesta**: $p = s_0, s_1, \dots, s_{l_p} \in X_{p_n}$,
- délka značkové cesty: l_p
- délka cesty: l_p nebo l_{p+1}

10/36

Cesty, značkové cesty

- Cesty $0, 1, 2, '003-04212'; 0, 5, 6, '001-00863'$ a $0, 9, 10, '045-00012'$ náležejí ke značkové cestě $books, book, id$,
- ...
- Cesty $0, 1, 4, 'J.R.R. Tolkien'; 0, 5, 8, 'J.R.R. Tolkien'$ a $0, 9, 12, 'Joseph Heller'$ náležejí ke značkové cestě $books, book, author$.



11/36

Body reprezentující cesty a značkové cesty

Definice 1 (bod n -rozměrného prostoru reprezentující značkovou cestu).

Mějme n -rozměrný diskrétní prostor značkových cest $\Omega_{p_n} = D^n$, $|D| = 2^D$. Mějme značkovou cestu $p_n = s_0, s_1, \dots, s_{l_p} \in X_{p_n}$, kde l_p je délka značkové cesty. **Bod n -rozměrného prostoru reprezentující značkovou cestu** je definován $t_{p_n} = (id_t(s_0), id_t(s_1), \dots, id_t(s_{l_p})) \in \Omega_{p_n}$, kde $id_t(s_i)$ je jedinečné číslo řetězce s_i , $id_t(s_i) \in D$. $n = \max(l_{p_n} + 1, p_n)$, $1 \leq i \leq |X_{p_n}|$. Každé značkové cestě p_n přiřadíme jedinečné číslo $id_{p_n}(p_n)$. ■

Definice 2 (bod n -rozměrného prostoru reprezentující cestu).

Mějme n -rozměrný diskrétní prostor cest $\Omega_p = D^n$, $|D| = 2^D$. Mějme cestu $p = id(u_0), id(u_1), \dots, id(u_{l_p})$, $s \in X_p$ a příslušnou značkovou cestu p_n s jedinečným číslem $id_{p_n}(p_n)$. **Bod n -rozměrného prostoru reprezentující cestu** je definován $t_p = (id_{p_n}(p_n), id_u(u_0), \dots, id_u(u_{l_p}), id_t(s)) \in \Omega_p$. $n = \max(l_{p_n} + 2, p_n)$, $1 \leq i \leq |X_p|$. ■

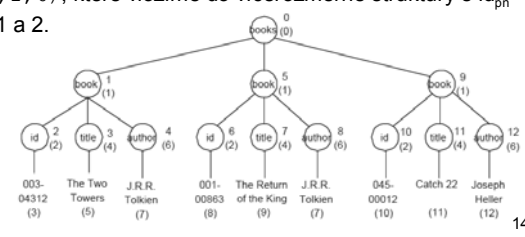
12/36

Indexy

- **Index termů** – uložení veškerých řetězců XML dokumentu s_i a jejich id_i (s_i).
- **Index značkových cest** – uložení bodů reprezentující značkové cesty.
- **Index cest** – uložení bodů reprezentující cesty.

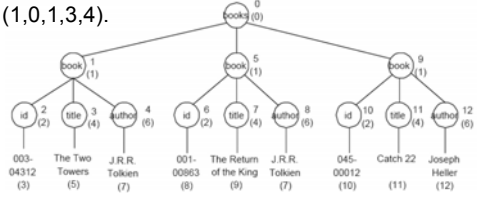
Příklad – index značkových cest

- $books, book, id; books, book, title$ a $books, book, author$. Pomocí id_i názvů elementů a atributů vytvoříme vektory $(0, 1, 2); (0, 1, 4)$ a $(0, 1, 6)$, které vložíme do vícerozměrné struktury s id_{pn} 0, 1 a 2.



Příklad – index cest

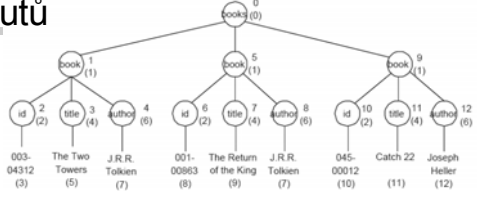
- Např. cesta k hodnotě The Two Towers. Jedná se o značkovanou cestu $book, book, title$ s id_{pn} 1. Po vložení jedinečného čísla značkové cesty id_{pn} , jedinečných čísel elementů id_i a id_i termu The Two Towers získáme vektor $(1, 0, 1, 3, 4)$.



Dotazy na hodnoty elementů a atributů

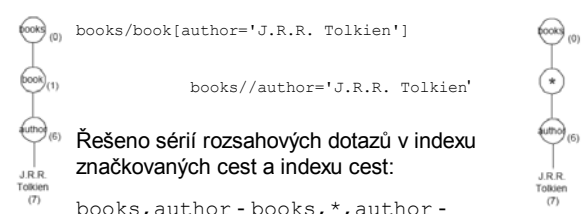
- Fáze dotazování:
 - Nalezení id_i termů dotazu v indexu termů.
 - Nalezení id_{pn} značkových cest dotazu v indexu značkových cest.
 - Nalezení vektorů v indexu cest.
- Problém vykonávání dotazů je převeden na problém definování a realizace bodových a rozsahových dotazů.

Dotazy na hodnoty elementů a atributů



- XPath dotaz: `books/book[author="Joseph Heller"]`
 - id termů z **indexu termů**,
 - id 2 značkové cesty `books/book/author` z **indexu značkových cest**: bodový dotaz $(0, 1, 6)$,
 - vektory z **indexu cest**: rozsahový dotaz $(2, 0, 0, 12) \times (2, \max, \max, 12)$.

Složitější dotazy



Řešeno sérií rozsahových dotazů v indexu značkových cest a indexu cest:
`books, author - books, *, author - books, *, ..., *, author` – počet rozsahových dotazů je n .

Složitější (disjunktní) rozsahové dotazy je možné realizovat v datové struktuře najednou.

[Implementace os XPath 1]

- Z indexu cest získáme body reprezentující cesty – získáme předky uzlu u : $id_u(u_0), \dots, id_u(u_{l-1})$. Osy `parent`, `ancestor` a `ancestor-or-self` uzlu u jsou tedy realizovány přímo.
- `descendent`: $[0, id_u(u_0), \dots, id_u(u_{l-1}), id_u(u), 0, \dots, 0] \times [max_D, id_u(u_0), \dots, id_u(u_{l-1}), id_u(u), max_D, \dots, max_D]$
- `child` – naivní přístup - $(l(u)+3)$. souřadnice

19/36

[Implementace os XPath 2]

- `preceding-siblings`: $[0, id_u(u_0), \dots, id_u(u_{l-1}), 0, 0, \dots, 0] \times [max_D, id_u(u_0), \dots, id_u(u_{l-1}), id_u(u)-1, max_D, \dots, max_D]$
- `following-siblings`: $[0, id_u(u_0), \dots, id_u(u_{l-1}), id(u)+1, 0, \dots, 0] \times [max_D, id_u(u_0), \dots, id_u(u_{l-1}), max_D, max_D, \dots, max_D]$
- `preceding`, `following`
- Osy `child`, `preceding-siblings` a `following-siblings` je nutné řešit složitějším způsobem.

20/36

[Komplikovanější dotazy]

- Dotazy na hodnoty a realizace osy XPath.
- např. `books/book[author='Joseph Heller']/title`
- Kombinace dříve popsaných technik:
 - dotaz na hodnotu,
 - zjištění sourozence.

21/36

[Další XML dotazovací jazyky]

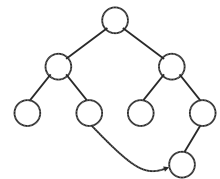
▪ XML-QL

```

where
<list><person>
  <name>Radim</name>
  <surname>Horáček</surname>
  <email>še<email>
</person></list> in ...
construct $e

```

ID, IDREF



▪ XQuery

```

document("data.xml")//article/[title='XML Book']/
  first_author/@author_id->author//surname

```

22/36

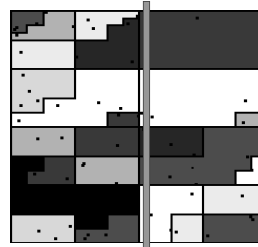
[Indexující datové struktury]

- Perzistentní vícerozměrné datové struktury – (B)UB-stromy, varianty R-stromů.
- Poskytují bodové a rozsahové dotazy.
- Problémy:
 - indexování vektorů různé dimenze,
 - úzké rozsahové dotazy.

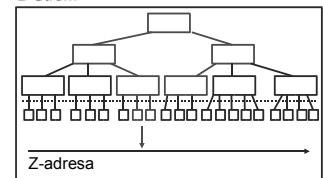
23/36

[(B)UB-strom]

UB-Strom

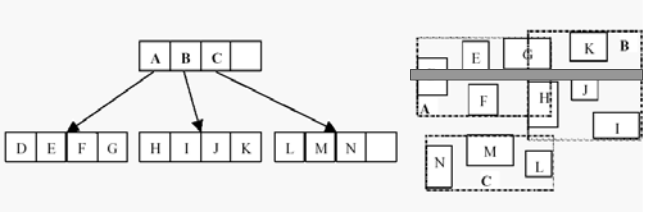


B-Strom



24/36

R-strom



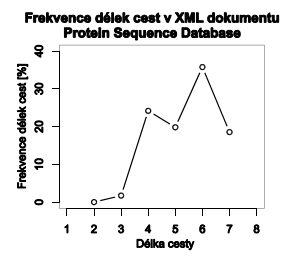
Indexování bodů různé dimenze

- Body reprezentující cesty a značkové cesty mají různou dimenzi.

	Maximální hloubka	Průměrná hloubka
DocBook	16	8.56
Shakespeare	6	4.77
XHTML	15	5.47

Vícerozměrné lesy

- BUB-les, R-les
- Každý strom lesa indexuje prostor různé dimenze.
- Např. XML dokument Protein Sequence Database: $n=7$ a $n=9$.

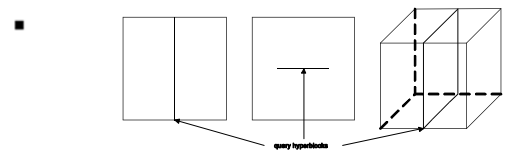


Úzké rozsahové dotazy

Definice 2 (Úzký rozsahový dotaz).

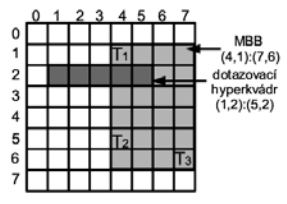
Mějme n -rozměrný diskrétní prostor $\Omega = D^n$. Dotazovací hyperkvádř je definován dvěma body $QL = (ql_1, ql_2, \dots, ql_n)$ a $QH = (qh_1, qh_2, \dots, qh_n)$, kde $\forall i : ql_i \leq qh_i$. Nechť konstanty ψ a ϕ jsou mohutnosti dosti malého ($\psi \rightarrow \min(D)$) resp. velkého ($\phi \rightarrow |D|$) intervalu $\subseteq D$. Rozsahový dotaz nazveme úzký pokud:

- $\forall i : qh_i - ql_i \leq \psi \vee qh_i - ql_i \geq \phi$.
- Označme počet dimenzí, pro které platí $qh_i - ql_i \leq \psi$ a n_ϕ pro které platí $qh_i - ql_i \geq \phi$. Pro úzký rozsahový dotaz musí platit $1 < n_\psi < n \wedge 1 < n_\phi < n$.

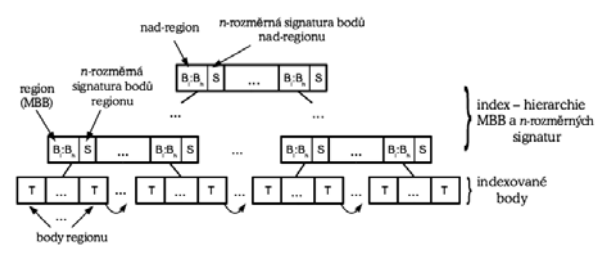


Signaturové vícerozměrné stromy

- Regiony aproximují přibližnou distribuci dat.
- Regiony protnuté dotazovacím hyperkvádrem jsou prohledány.
- S rostoucí dimenzí poměr relevantních (N_R) a protnutých (N_I) regionů $c_R \ll 1$.
- Vložení signatury dochází k přesnější aproximaci distribuce dat.



Signaturový R-strom



Výsledky experimentů

- Databáze bílkovin z XML UW projektu:
 - velikost souboru: 683MB,
 - počet elementů: 21 305 818,
 - počet atributů: 1 290 647.
 - maximální délka cesty 7.
- BUB-les, R*-les, Signaturový BUB-les a R*-les. Index struktury: stromy indexující prostory dimenze 7 a 9.

31/36

Výsledky experimentů

Tabulka 1. Charakteristika indexujících datových struktur.

Dimenze prostoru n	Počet bodů	Velikost indexu [MB]					
		BUB-strom	Signaturový BUB-strom	R*-strom	Signaturový R*-strom	BUB-strom	Signaturový R*-strom
7	8 268 357	440.9	471 [+7%]	478.6	512.2 [+7%]		
9	8 739 522	562.1	635.2 [+13%]	603.1	680.7 [+13%]		

Dimenze prostoru n	Počet vnitřních uzlů				Počet listových uzlů			
	BUB-strom	Signaturový BUB-strom	R*-strom	Signaturový R*-strom	BUB-strom	Signaturový BUB-strom	R*-strom	Signaturový R*-strom
7	10 917	22 432	15 731	33 186	214 842	234 412	256 520	258 187
9	17 751	36 412	24 750	55 750	270 065	298 142	318 370	331 474

32/36

Výsledky experimentů

ProteinDatabase/ProteinEntry/[reference/refinfo/authors/author='Smith, E.L.']

Tabulka 2. Charakteristika testovacích úzkých rozsahových dotazů.

Množina dotazů	Prostor dimenze n	Velikost výsledku	N_I		N_R		c_R	
			BUB-strom	R*-strom	BUB-strom	R*-strom	BUB-strom	R*-strom
1	7	5	116	828	1	1	0.009	0.0012
2	7	3 397	139 927	1 542	4 251	717	0.030	0.47
3	9	8	243	641	15	7	0.060	0.01
4	9	2 794	112 927	136	3 745	136	0.033	1

33/36

Výsledky experimentů

Tabulka 3. Statistika realizace rozsahových dotazů v indexu cest.

Množina dotazů	Procento prohledávaných listových uzlů				DAC				Čas vykonání dotazu [s]			
	BUB-strom	Sign. BUB-strom	R*-strom	Sign. R*-strom	BUB-strom	Sign. BUB-strom	R*-strom	Sign. R*-strom	BUB-strom	Sign. BUB-strom	R*-strom	Sign. R*-strom
1	0.05	0.04	0.32	0.061	324	319	960	74	0.13	0.1	0.08	0.015
2	2.3	0.4	0.60	0.45	15 275	11 24	1 671	1 043	5.8	0.9	0.19	0.094
3	0.04	0.03	0.20	0.013	513	489	817	477	0.2	0.2	0.20	0.140
4	2.2	0.4	0.05	0.03	14 451	1 245	220	184	5.6	0.8	0.017	0.015
Průměr	1.15	0.22	0.29	0.14	7 566	794	917	445	2.9	0.5	0.12	0.07

34/36

Závěr

<http://www.cs.vsb.cz/arg>



- Dotazování na částečnou shodu hodnot elementů a atributů.
- Komprese datové struktury.
- Implementace složitějších rozsahových dotazů.
- Kombinace s přístupy indexování nestrukturovaných dokumentů v IR.
- Implementace podmnožiny dotazovacích jazyků jako např. XQuery.

35/36

Reference

- M. Krátký, J. Pokorný, T. Skopal, V. Snášel: *The Geometric Framework for Exact and Similarity Querying XML data. In Proceedings of EurAsia-ICT 2002*. Shiraz, Iran, Springer Verlag, LNCS 2510.
- M. Krátký, T. Skopal, and V. Snášel: *Multidimensional Term Indexing for Efficient Processing of Complex Queries*. *Kybernetika, Journal of the Academy of Sciences of the Czech Republic*, 2003, accepted.

36/36