



Physische XML-Speicherstrukturen und Indizes

Holger Meyer

Database Research Group

University of Rostock

E-Mail: hm@GUUG.de

1. Überblick XML-Dateiorganisation und Zugriffspfade
2. Werteindizes
3. Volltextindizes
4. XML-Dateiorganisation
5. Struktur- und Pfadindizes
6. Zusammenfassung

- ⑥ *Werteindex*, atomar Elementinhalte, Attributwerte, strukturierte Anteile eines XML-Dokumentes
- ⑥ *Volltextindex* einzelnen Worte des *Volltextes*, ganzen Sätzen oder Phrasen einer natürlichen Sprache, Zugriff mit *Information Retrieval*-Techniken.
- ⑥ *Struktur- und Pfadindex* effiziente Auswertung von Anfragen unter Einbeziehung der Dokumentstruktur.
- ⑥ physische *Speicherstrukturen* zur Ablage von XML-Fragmenten

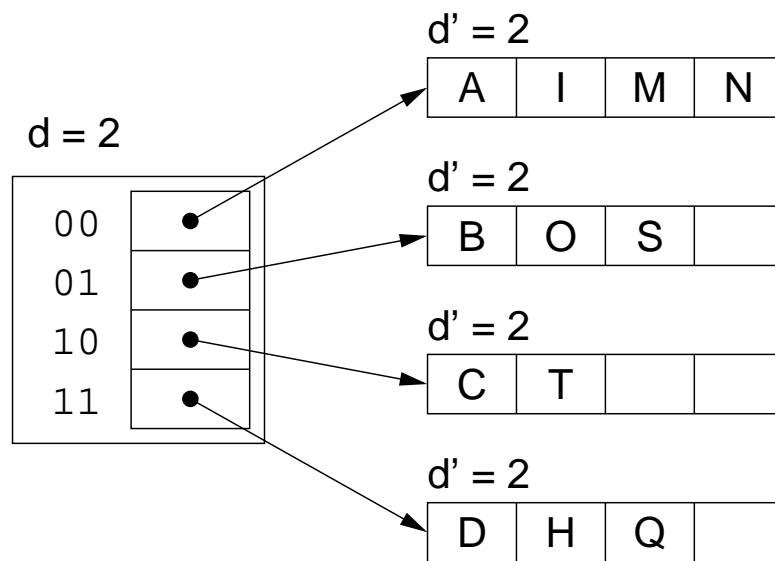
1.1 Kriterien

- ⑥ Dateiorganisationsform oder Zugriffspfad
- ⑥ Primärindex oder Sekundärindizes
- ⑥ Statische oder dynamische Strukturen
- ⑥ Ein- oder Mehrattributindex, ein- oder mehrdimensionaler Index
- ⑥ Erhaltung einer Ordnung
- ⑥ Form der unterstützten Anfragen: Punktanfragen, Bereichsanfragen, partielle Anfragen, vage Anfragen
- ⑥ Pfadindex: Schemainformation vorhanden oder nicht
- ⑥ Pfadindex: akkurat, eindeutig, vollständig

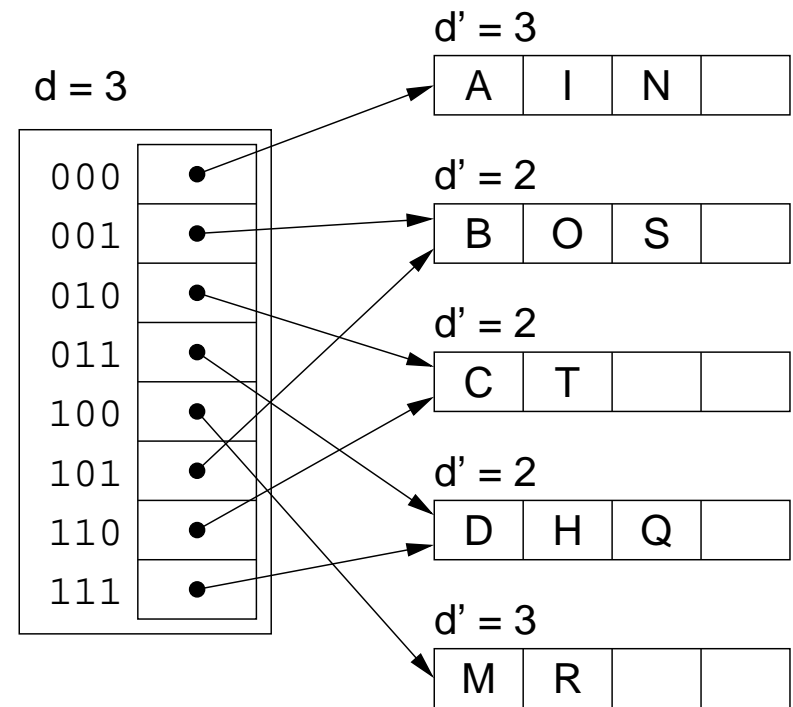
2 Werteindex

- ⑥ Hash-Strukturen am Beispiel Erweiterbares Hashing
 - △ Punktanfragen, $O(c)$
- ⑥ dynamische Index-Bäume am Beispiel B⁺-Baum
 - △ Bereichsanfragen, Präfixsuche
- ⑥ Signatur-Bäume
 - △ n -dimensionaler Zugriff, Ähnlichkeitssuche

2.1 Erweiterbares Hashing



(a) vor



(b) nach dem Einfügen
von „R“

2.2 B⁺-Baum

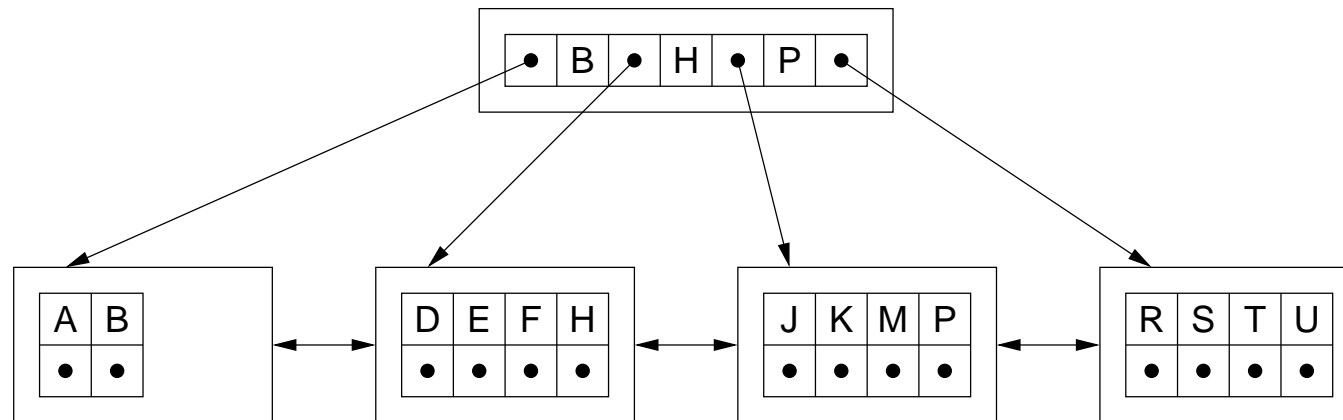


Abbildung 2: B⁺-Baum, Ordnung $M = 4$

B⁺-Baum (cont'd)

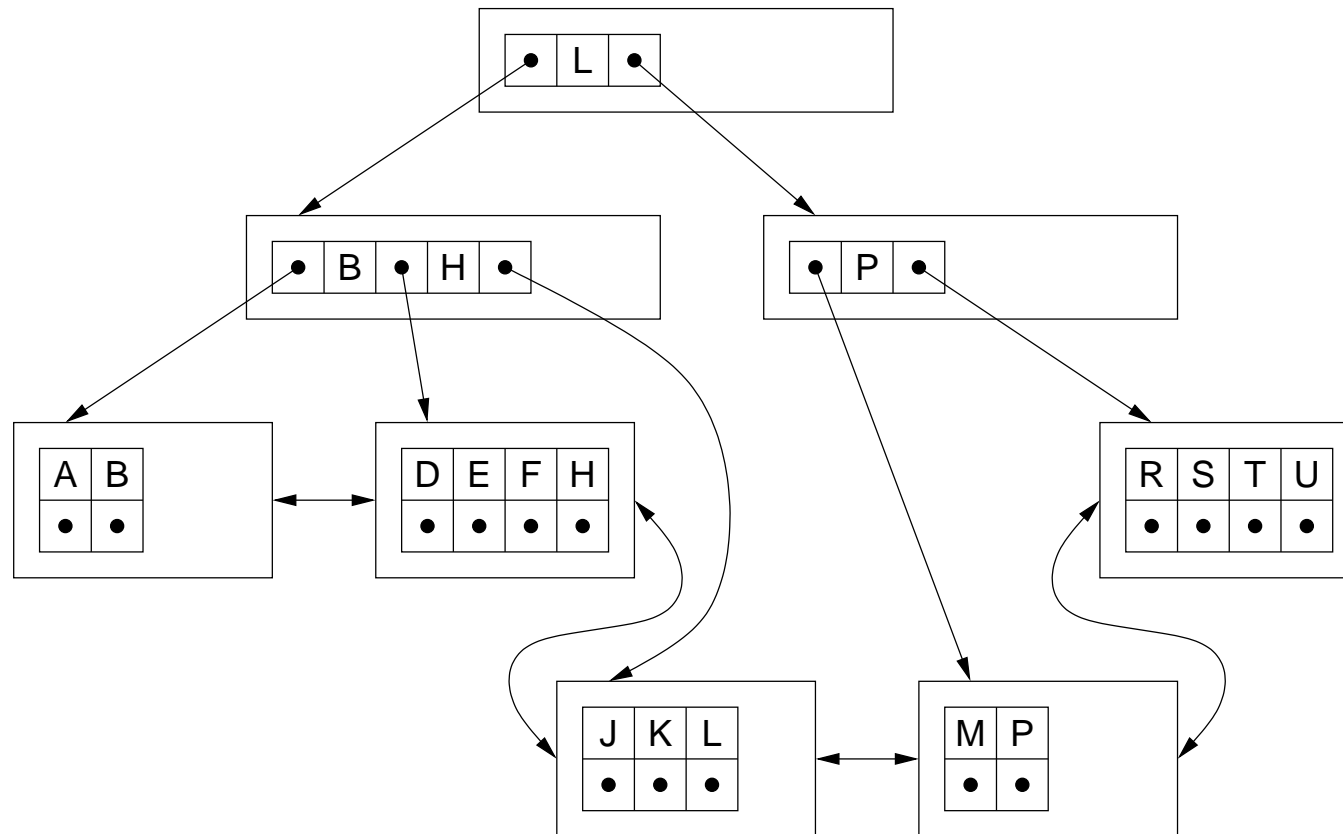


Abbildung 3: nach dem Einfügen von „L“

2.3 Signaturbaum

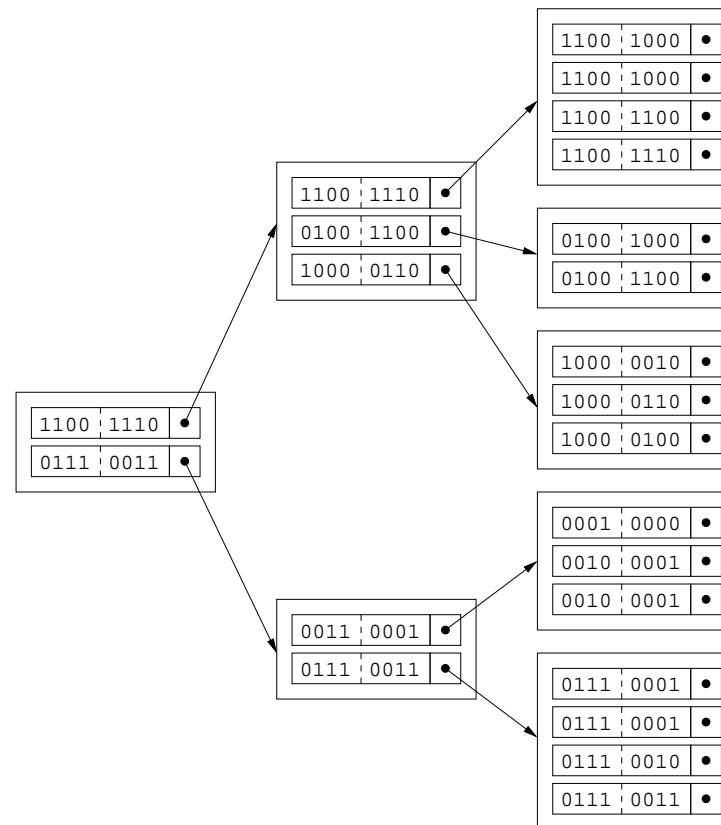


Abbildung 4: Baum, Ordnung $M = 4$ und $m = 2$

Signaturbaum (cont'd)

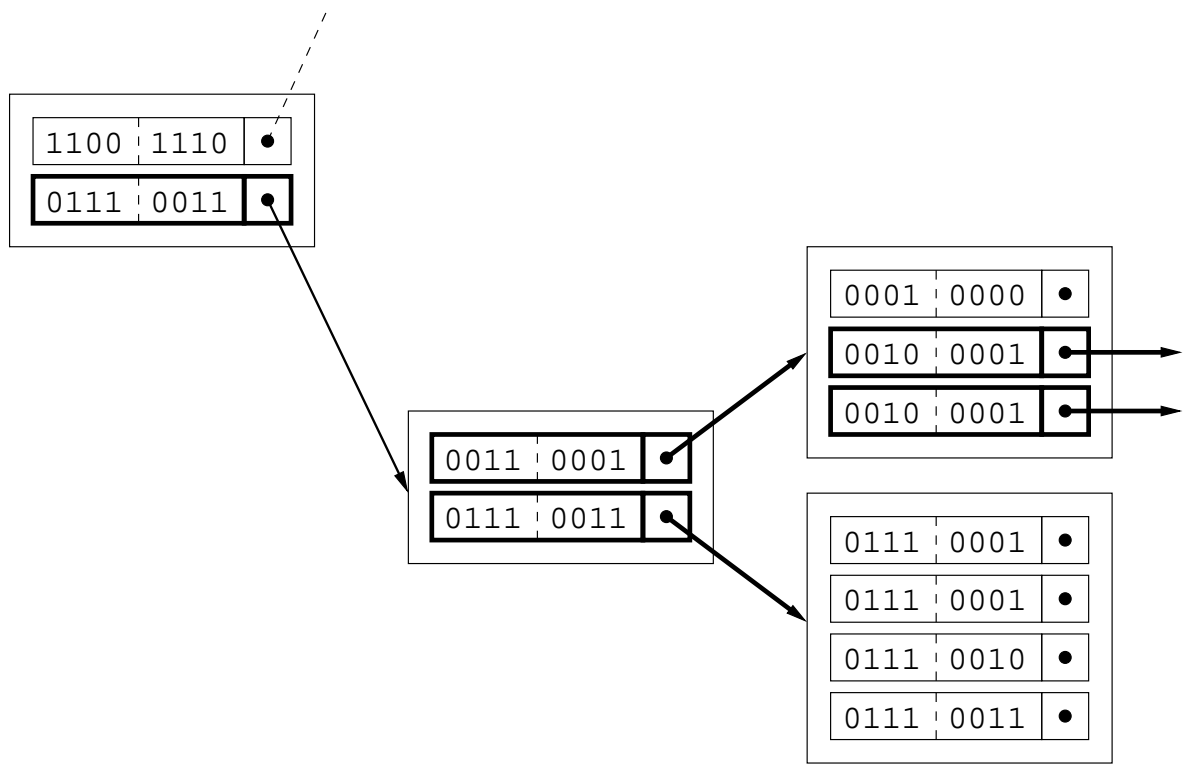


Abbildung 5: Suche nach der Signatur 00100001

3 Volltextindex

Anfragen an Volltext kann man grob in drei Formen unterteilen:

- ⑥ Stichwortbasierte Suche,
- ⑥ Mustersuche und
- ⑥ Strukturierte Anfragen.

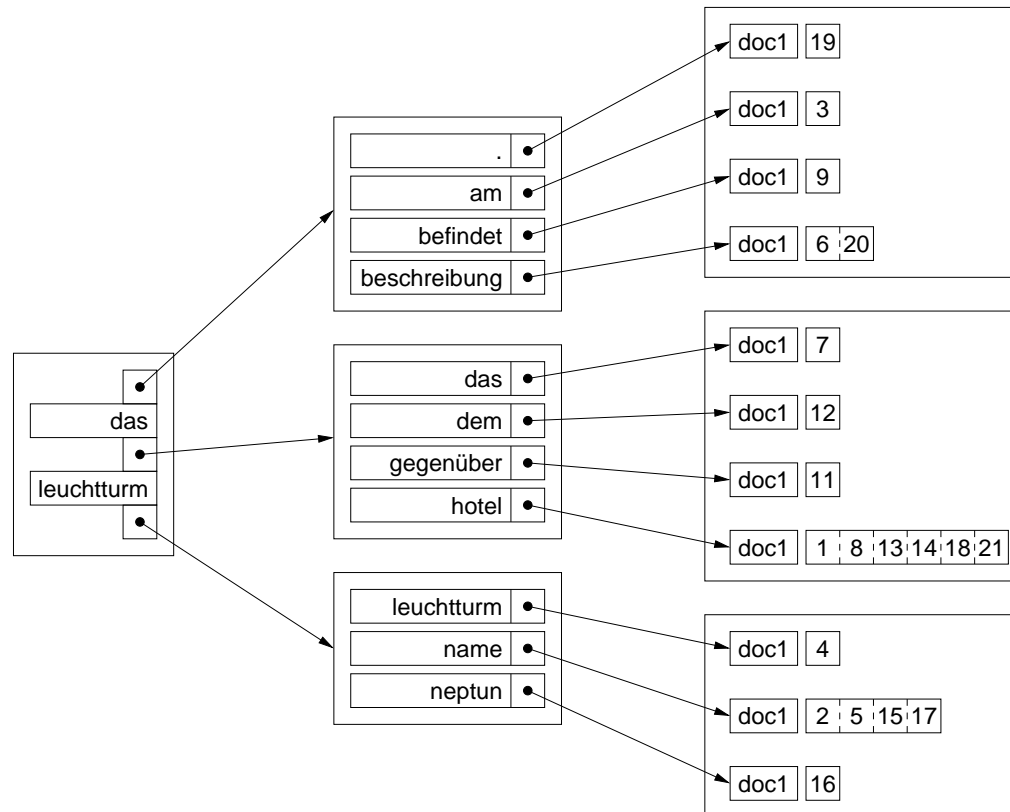
3.1 Invertierte Datei

```
<hotel1>  
  <name2>Am3 Leuchtturm4</name5>  
  <beschreibung6>  
    Das7 Hotel8 befindet9 sich10 gegenüber11 dem12 Hotel13  
    <hotel14><name15>Neptun16</name17></hotel18>.19  
  </beschreibung20>  
</hotel21>
```

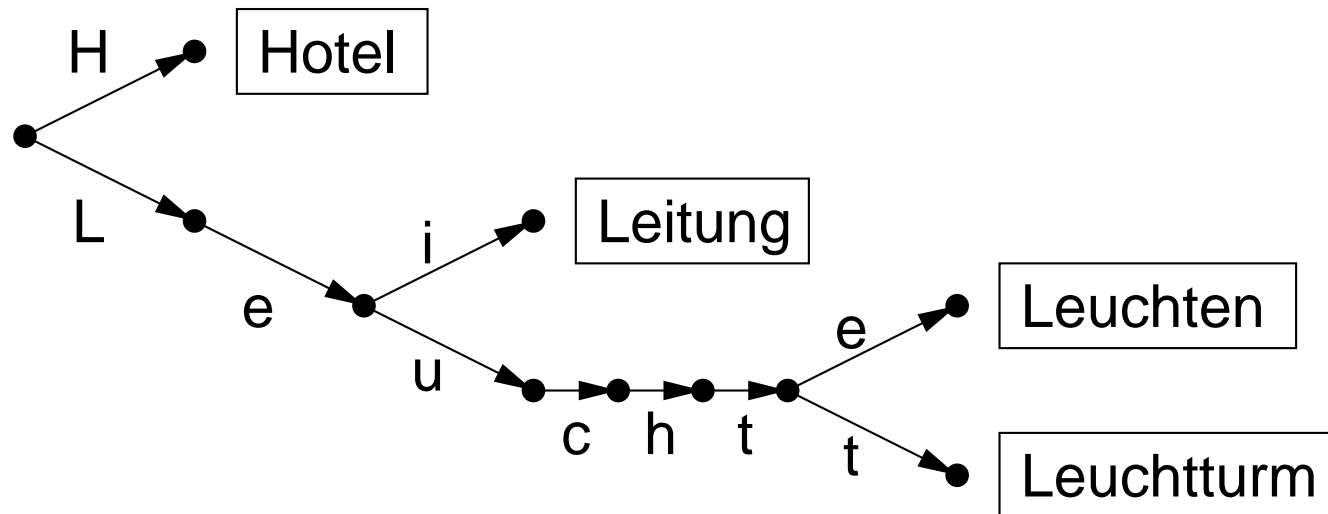
Invertierte Listen — Beispiel

Word	Vorkommen
.	{ 19 }
am	{ 3 }
befindet	{ 9 }
beschreibung	{ 6, 20 }
das	{ 7 }
dem	{ 12 }
gegenber	{ 11 }
hotel	{ 1, 8, 13, 14, 18, 21 }
leuchtturm	{ 4 }
name	{ 2, 5, 15, 17 }
neptun	{ 16 }

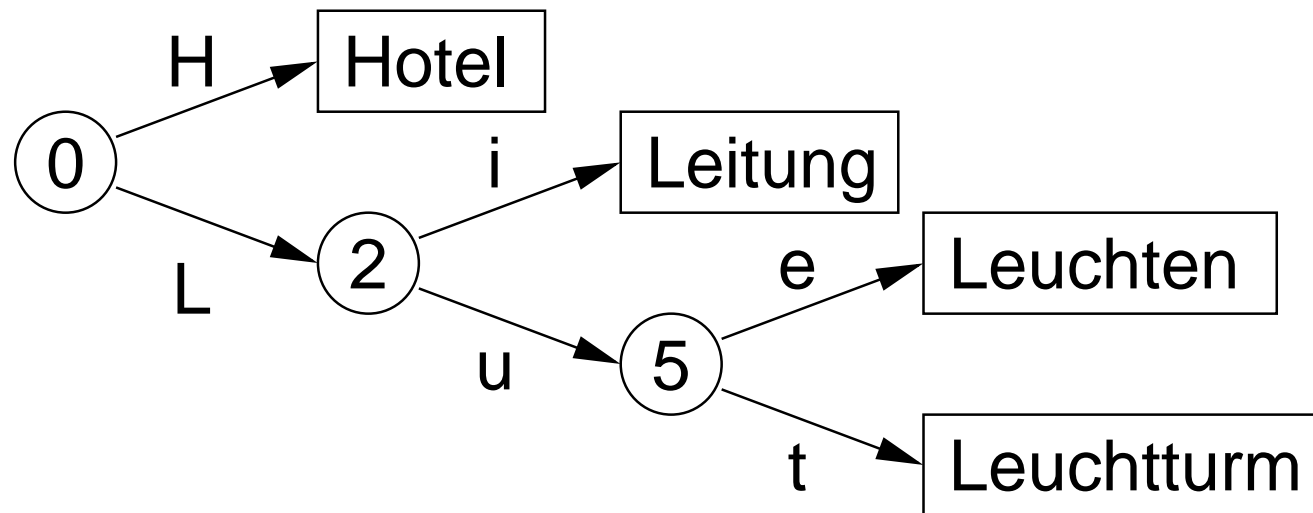
B⁺-Baum, Ordnung $M = 4$



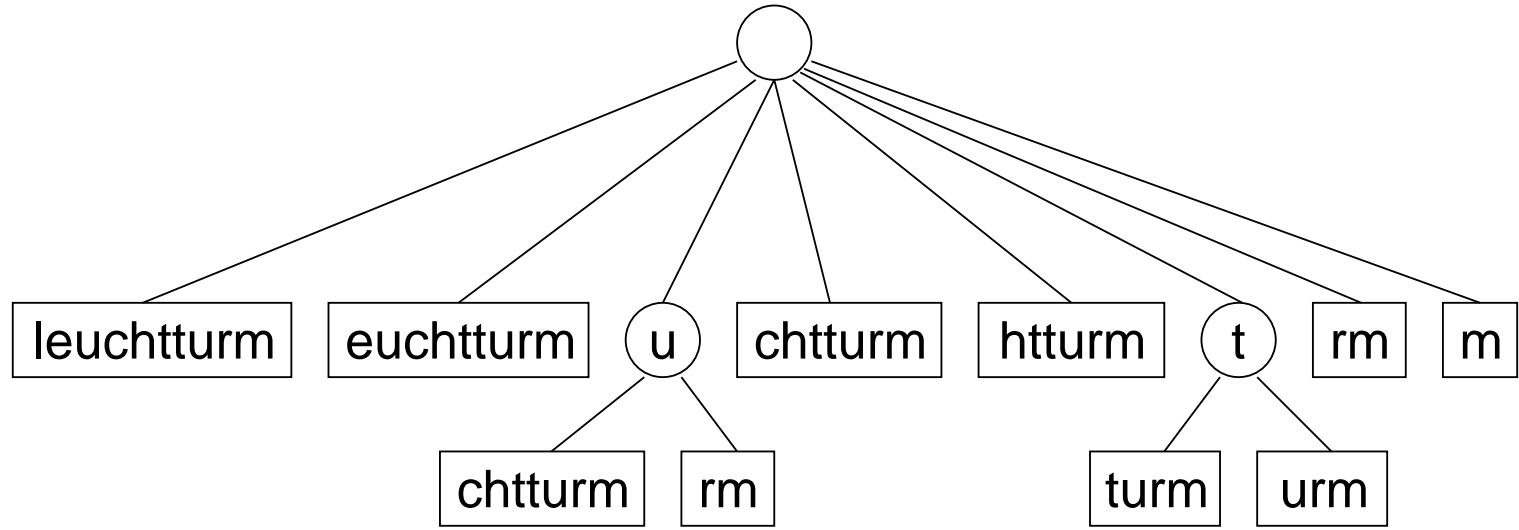
3.2 Tries



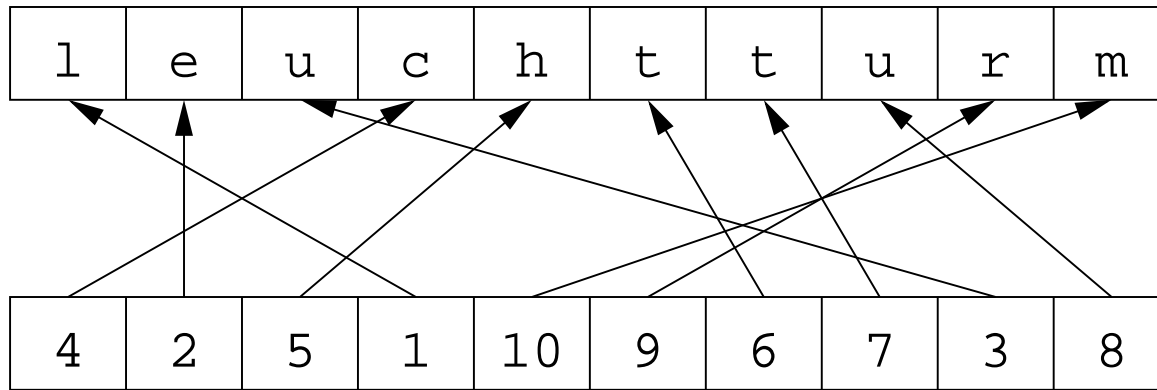
3.3 Patricia-Baum



3.4 Suffixbaum



3.5 Suffix Arrays



Vergleich Volltextindizes

Kriterium	Invert. Datei	Suffix Array	Signatur- baum
exakte Suche	✓	✓	✓
Bereichsanfragen	✓	✓	—
Präfixsuche	✓	✓	—
Teilzeichenketten	—	✓	—
partielle Anfragen	—	✓	✓

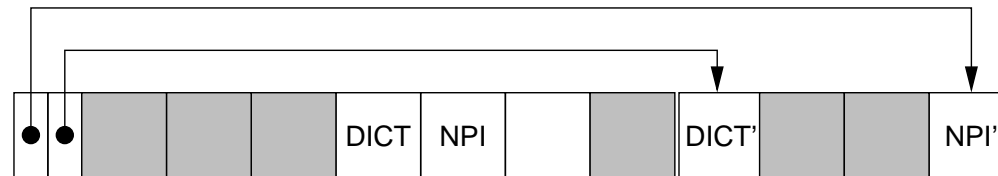
4 Dokumentenspeicherung

- ⑥ komplette physische Speicherung von XML-Dokumenten
- ⑥ → Speicherung der (DOM-)Knoten, Clustering
- ⑥ Indizierung der Baum- oder Graphenstruktur eines XML-Dokumentes
- ⑥ → Pfadindizes, XPath sind nicht General Path Expression (GPE)

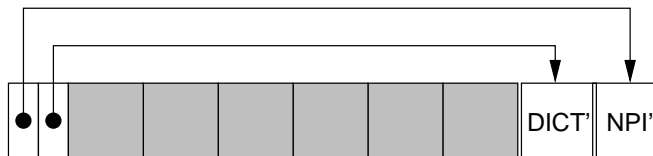
4.1 Persistentes DOM in Infonyte-DB

- ⑥ persistente Speicherung der durch DOM implizierten Struktur
- ⑥ pro XML-Dokument eine Datei
- ⑥ seitenorganisiert: Knotenseite aus 128 DOM-Knoten
- ⑥ Knotenindex (Node page index, NPI), Elementindex (Dictionary, DICT)
- ⑥ Strukturindex auf Basis Knotensignatur
- ⑥ Volltextindex auf Basis invertierter Dateien
- ⑥ Xpath, XSLT, XQL, XQuery

Aufbau einer *PDOM*-Datei



(a) vor

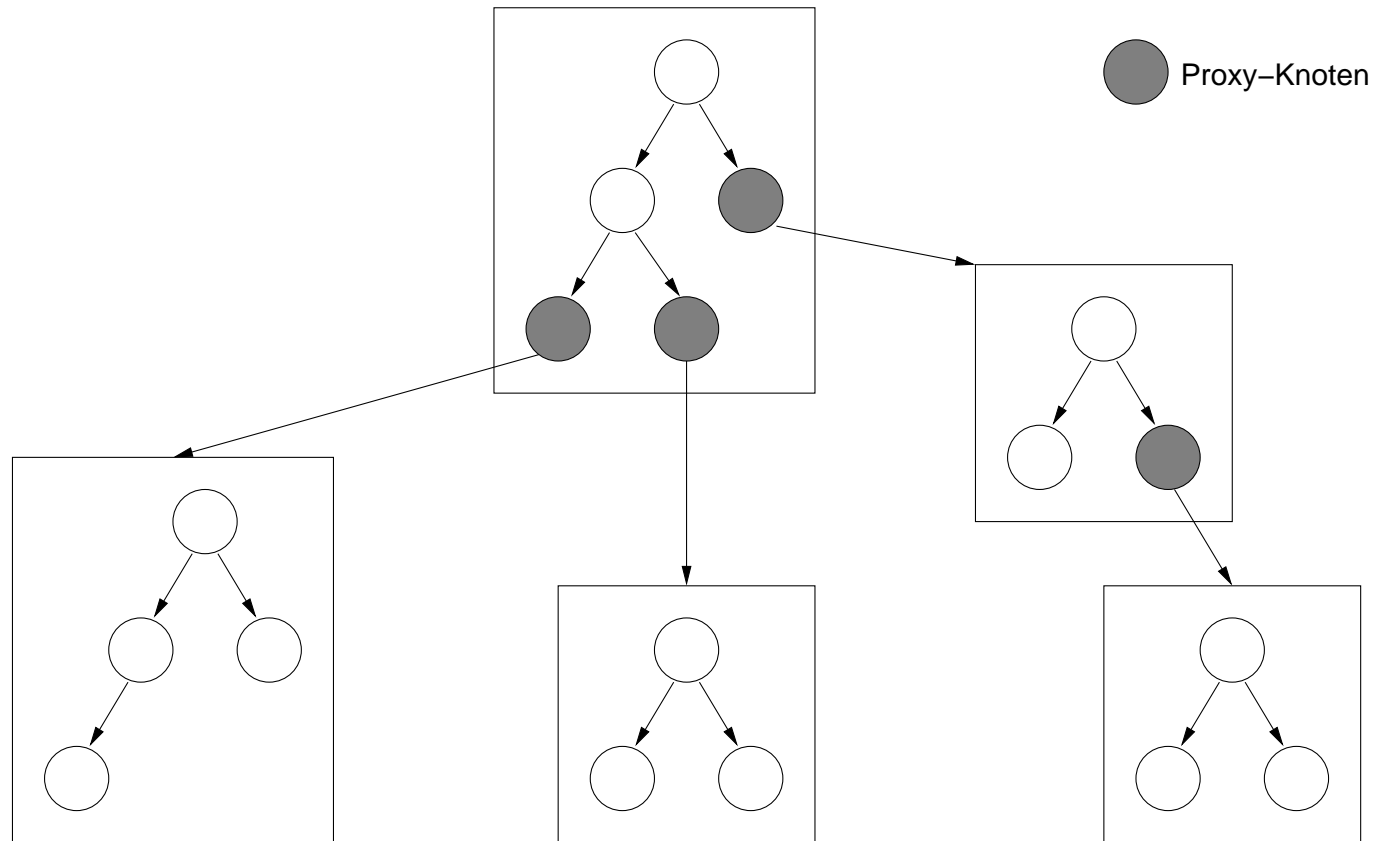


(b) nach Reorganisa-
tion

4.2 Hybride Speicherung in Nativ

- ⑥ Knotenstruktur analog DOM,
- ⑥ seitenbasierte Speicherung,
- ⑥ Clustering: Teilbäume, Proxy-Knoten,
- ⑥ B-Baum-ähnliche Struktur, Splitting, Merging

Dokument und Clustering in Nativ



5 Struktur-/Pfadindizes

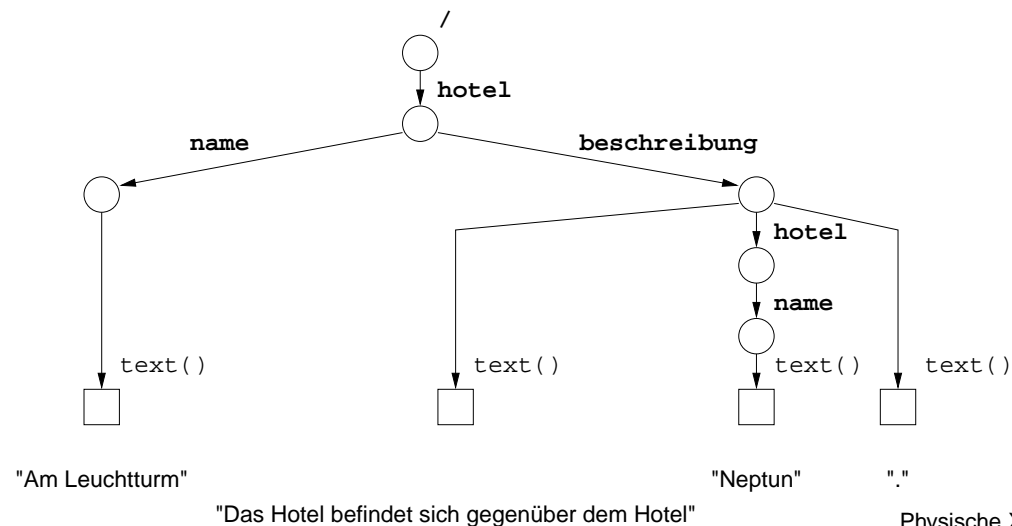
- ⑥ allgemeine Pfadausdrücke (GPE)
- ⑥ mit oder ohne Schema (DTD, XML Schema)
- ⑥ *eindeutig*, jeder Pfad genau einmal
- ⑥ *akkurat*, alle Pfade exakt abgebildet
- ⑥ *vollständig*, ganze Pfade, nicht nur Teilpfade

6 Pfadauswertung

- ⑥ einfache Pfade: $/a/b$, $a/b/c$
- ⑥ Achsennavigation: Vorgänger, Nachfolger, Geschwister, ...
- ⑥ einfache Muster: $a//b$, $a/*/b$
- ⑥ komplizierter: $//$ und $*$ in Kombination
- ⑥ allgemein: Pattern matching in Bäumen, reguläre Ausdrücke, NFA ($O(nm)$)

6.1 XML-sensitiver Volltextindex

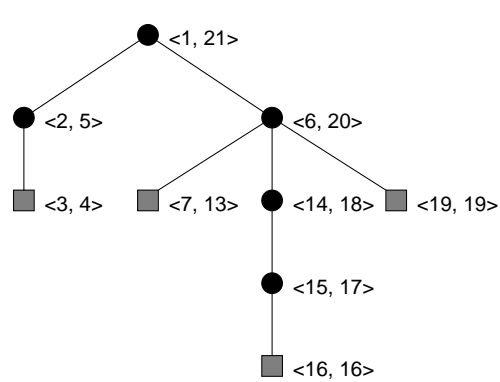
```
<hotel1>  
  <name2>Am3 Leuchtturm4</name5>  
  <beschreibung6>  
    Das7 Hotel8 befindet9 sich10 gegenüber11 dem12 Hotel13  
    <hotel14><name15>Neptun16</name17></hotel18>.19  
  </beschreibung20>  
</hotel21>
```



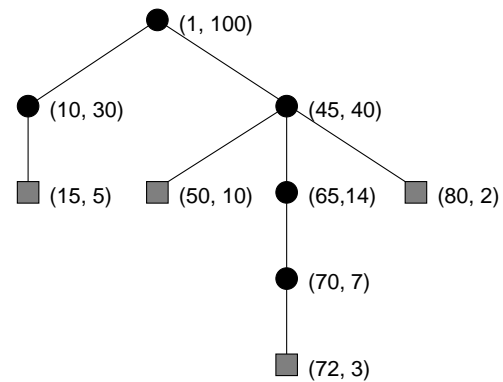
Numerierungsschemata

- ⑥ Kodierung der Knotenabfolge → Achsenavigation, Volltextoperationen
 - △ Dietz (pre-, post-order), Element-Subelement, Vorgänger/Nachfolger, Phrasen,
 - △ Position-Size, Updates,
 - △ Bitvektor, Element-Subelement, Vorgänger/Nachfolger, Subskript, Hierarchie
- ⑥ Grundlagen: Bereichsalgebren, Suche in metrischen Räumen

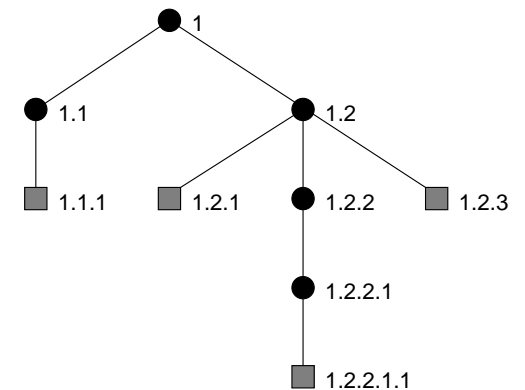
Numerierungsschemata (cont'd)



(c) Dietz

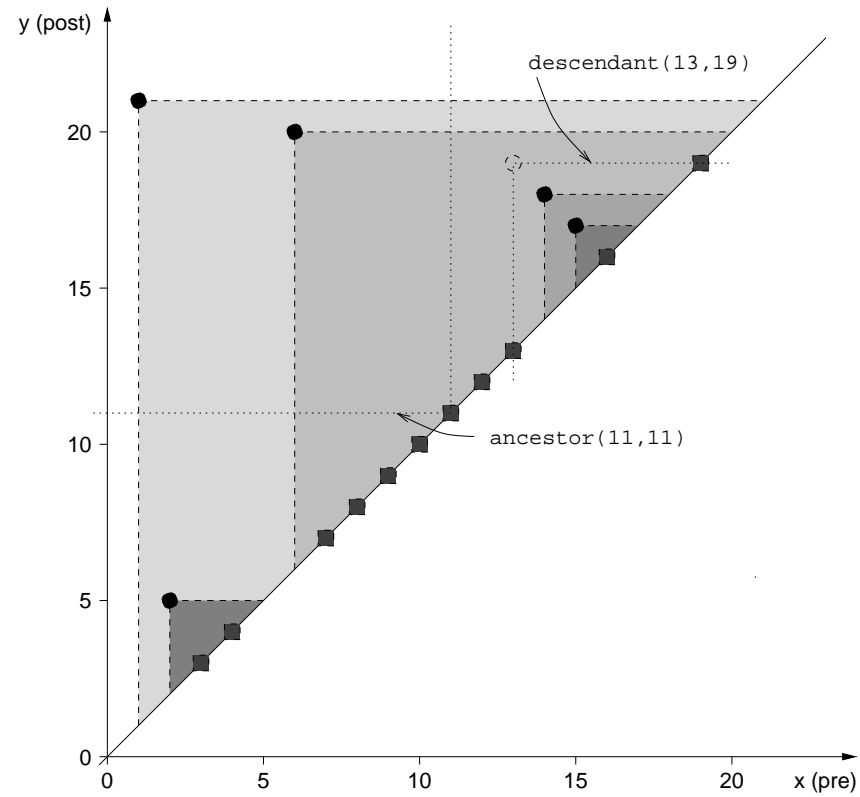


(d) Size



(e) Bitvektor

Numerierungsschemata (cont'd)



- ⑥ Wortposition (*Start-Tag*, *Ende-Tag*)
- ⑥ **CONTAINS**-Anfragen, SQL/MM Fulltext
- ⑥ einfache Pfadanfragen, Element-Subelement
- ⑥ Aufbau analog Volltextindex mit Wortpositionen

Element	Vorkommen
beschreibung	{(6, 20)}
hotel	{(1, 21), (14, 18)}
name	{(2, 5), (15, 17)}

Anfragetyp: *elem//subelem*

```
select sub.*  
from element e, element sub  
where e.begin < sub.begin  
       and sub.end < e.end  
       and e.name = elem  
       and sub.name = subelem
```


6.2 XASR

- ⑥ Pfadindizierung, Pre-Post-Order, Vaterknoten,
- ⑥ Element-Subelement-Anfragen → Verbundfolgen

docid	position		dparent	element
	dmin	dmax		
4711	1	21	0	hotel
4711	2	5	1	name
4711	6	20	1	beschreibung
4711	14	18	6	hotel
4711	17	17	14	name

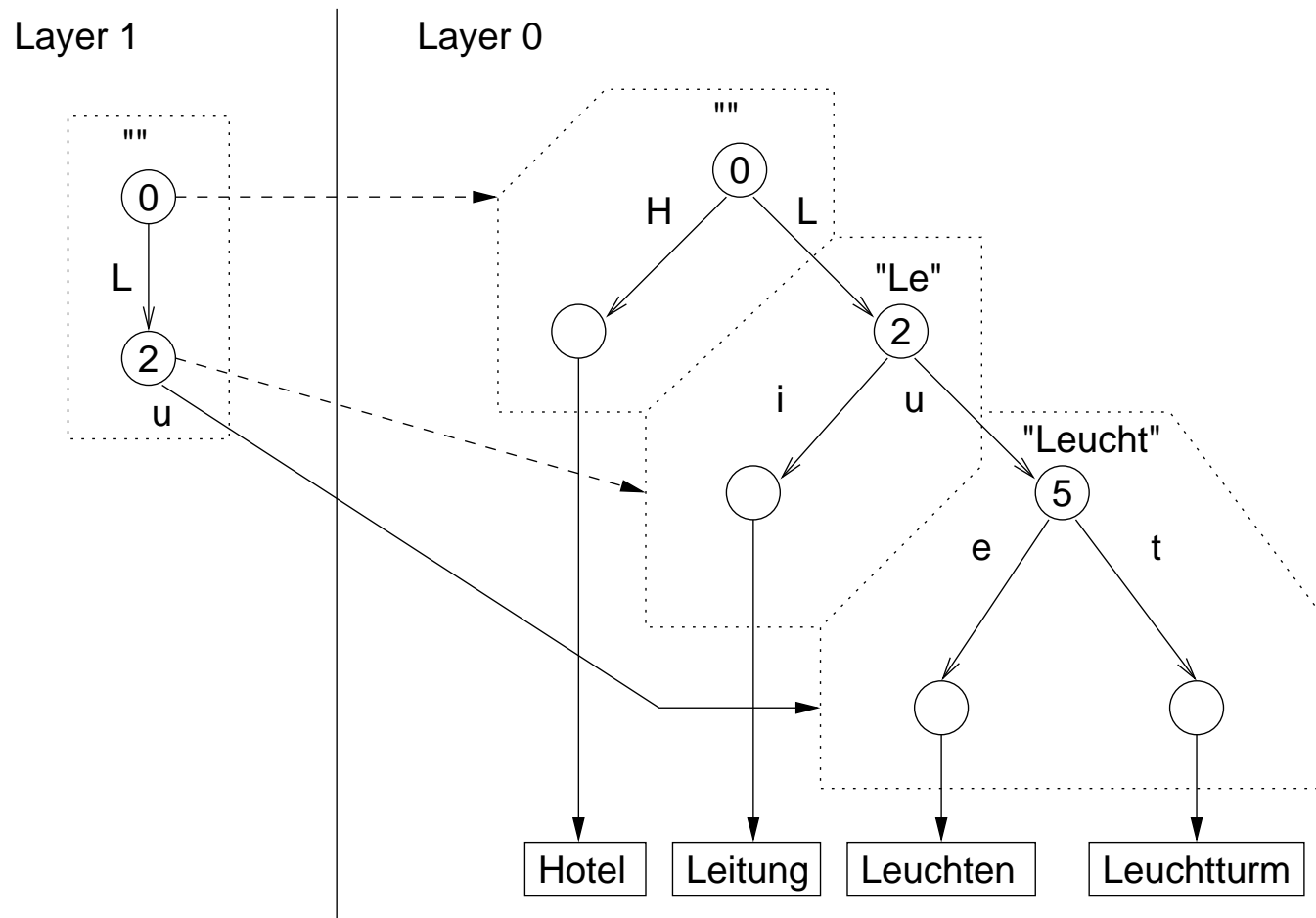
Anfragetyp: *elem/subelem*

```
select sub.*  
from element e, element sub  
where e.dmin = sub.dparent  
       and e.name = elem  
       and sub.name = subelem
```

6.3 Index Fabric

- ⑥ Vollständige Pfade, ausgewählte Teilpfade für häufige Anfragen
- ⑥ akkurat, eindeutig
- ⑥ kein Schema
- ⑥ Mehr-Ebenen-Patricia-Baum
- ⑥ spezielle Elementkodierung
 - △ Elementnamen als Zeichen eines erweiterten Alphabets
 - △ Präfixnotation
 - △ Beispiel Hotelnamen: `hotel` → **H**, `name` → **N**
 - „**H N** Am Leuchtturm“
 - „**H N** Neptun“

Mehr-Ebenen-Patricia-Baum



6.4 Pfad-Indizierung in SphinX

- ⑥ Schema notwendig
- ⑥ Schemagraph und B^+ -Baum
- ⑥ akkurat, eindeutige Pfadindizierung
- ⑥ GPE

Sphinx Index-Struktur

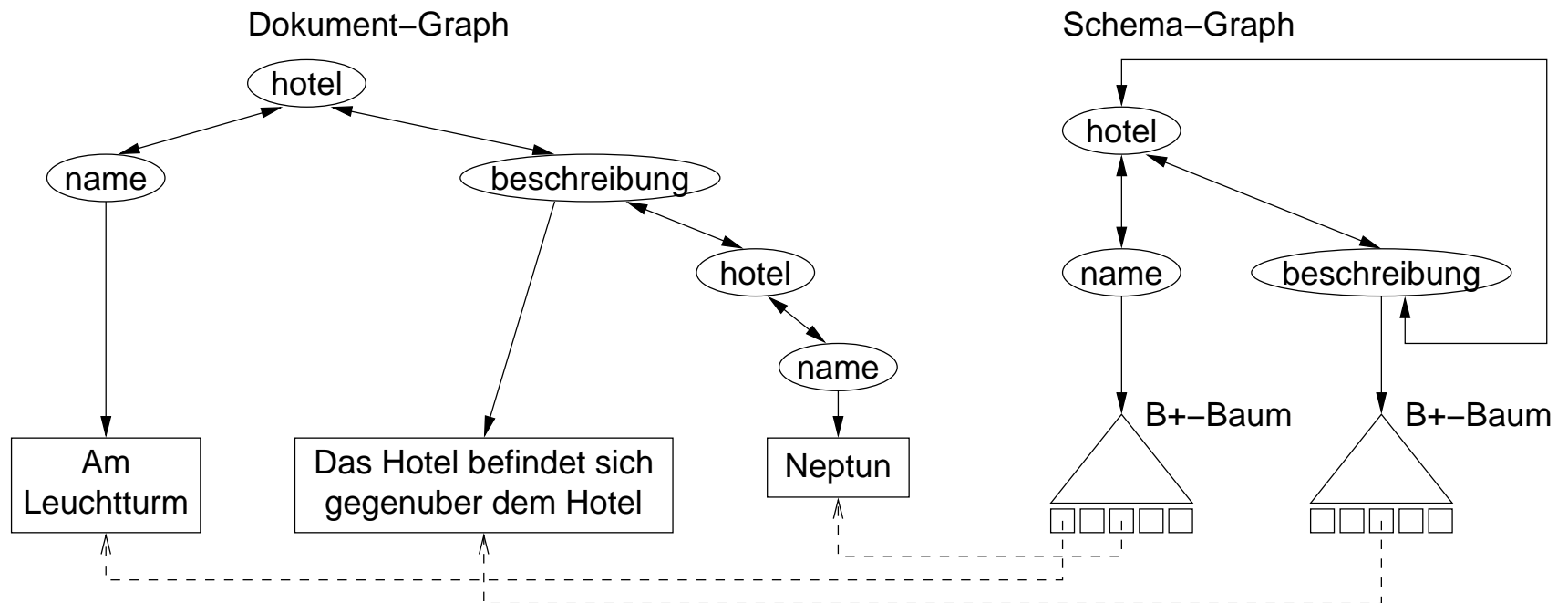


Abbildung 6: Dokument- und Schemagraph

7 Zusammenfassung

- ⑥ Vielzahl weiterer Techniken: Dataguides, 1-Index, 2-Index, T-Index, ToXin, ...
- ⑥ Unterstützung nicht nur für einfache Pfadausdrücke, sondern GPE
- ⑥ Aufwand für GPE: analog zu regulären Ausdrücken, NFA, DFA
- ⑥ Verbesserung wenn Schema vorhanden (Dokument- vs. Schemagraph)