

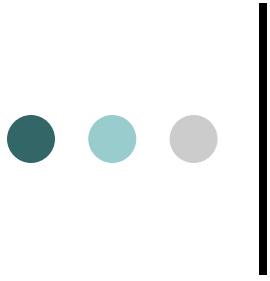


$C=(FS)^2$: Cubing by Composition of Faceted Search

Ronny Lempel

Dafna Sheinwald

IBM Haifa Research Lab



Introduction to Multifaceted Search and to On-Line Analytical Processing (OLAP)

Intro

Multifaceted Search With Aggregated Measures

IBM Research Faceted Search with BI - Results - Mozilla Firefox

http://lnx-trec.haifa.ibm.com:2407/search?query=travel+europe

travel europe Search settings

Found 6612 matching book records in 1.08 seconds. Showing results 1 - 10.

Narrow search by:

Publication Date	Author	Min Price	Average Price	Max Price	Publisher	Cent per page	Cost effectiveness	More Facets
2001 (538)								Category
2003 (697)								Subject
2004 (768)	Fodor's (126)	9.95	17.47	27.5	William Morrow & Co (3)	19.53	5.01	Rating(worst 1-10 best)
2005 (887)	DK Publishing (81)	7	18.69	50	Dover Publications (28)	5.47	1.64	Binding
2006 (689)	Rick Steves (69)	4.95	18.17	69.95	Carter Bowman (1)	6.23	1.15	Salesrank
Other (3033)	Karen Brown (56)	12.95	19.45	19.95	DK Travel (115)	13.92	1.13	Price
	Other (7197)	1.25	23.23	300	Other (6465)	16.51	1.1	Weight (grams)

Traveling the Eurail Express (Traveling Europe's Trains)

By **Jay Brunhouse**, April 2004 (Pelican Publishing Company).
Paperback. \$18.95.

Traveling the Eurail Express (Traveling Europe's Trains). . Jay Brunhouse

Travellers Krakow, 2nd (Travellers - Thomas Cook)

By **Scott Simpson**, April 1 2006 (Thomas Cook Publishing).

query

3 featured dimensions

other dimensions

Drill-down options

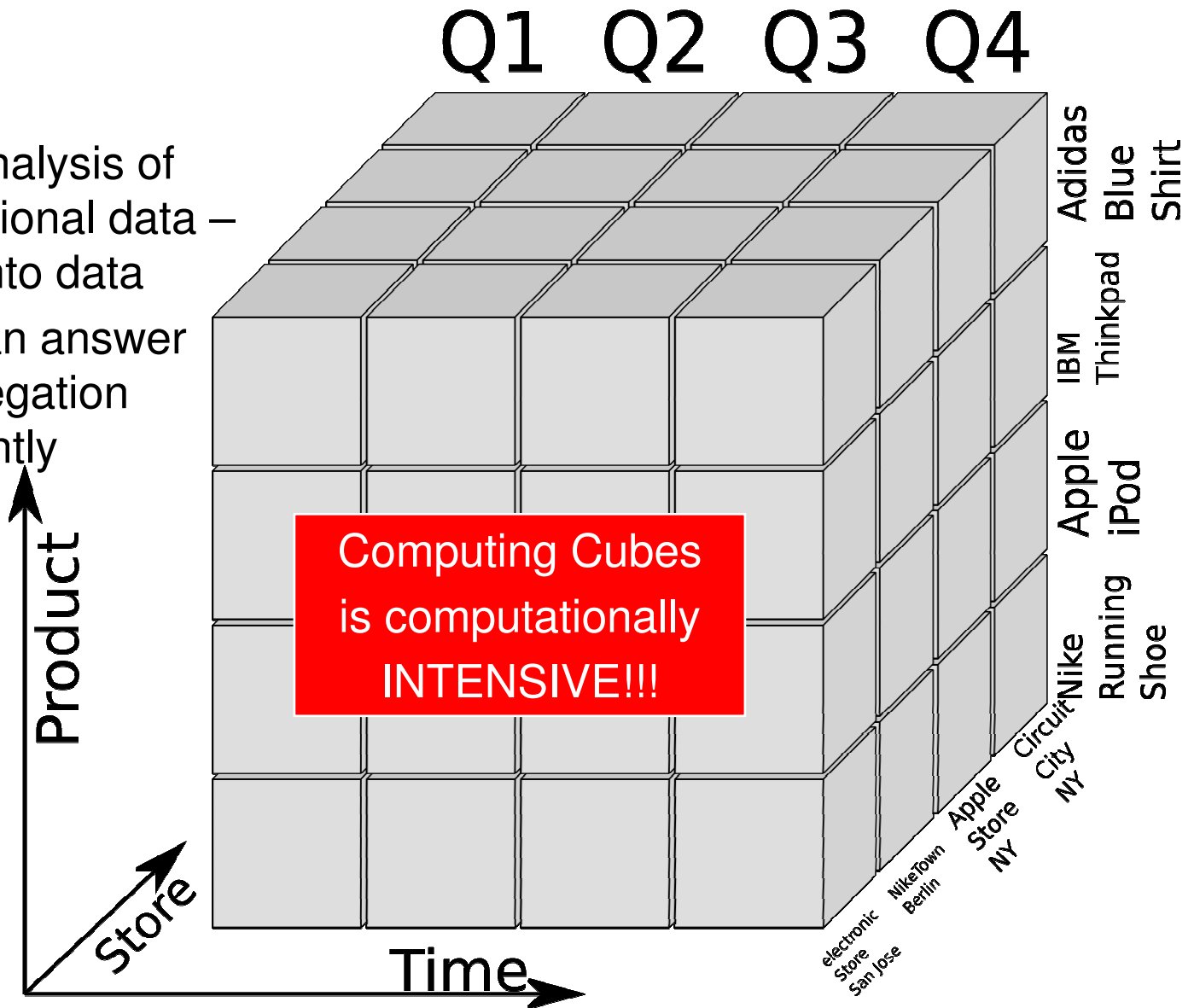
search results

Result counts and aggregated measures

Intro

OLAP Cubes

- Primary tool for multidimensional analysis of structured and relational data – qualitative insight into data
- Once computed, can answer selection and aggregation queries very efficiently
 - Drill down
 - Roll up
 - Measures



Goal: General OLAP Cubing over Search Results

- Observation: faceted search with BI expressions is practically equivalent to having several “one-dimensional cubes” with measures
 - Equivalently, cubes can be seen as aggregations over the cross-product (Cartesian product) of several facets
- One cannot compute cubes accounting for free-text queries, since the space of possible queries is practically infinite
- Two natural approaches for combining free-text search and cubing:
 - Pipe the search engine’s results into a cubing engine
 - Extend a faceted search engine to aggregate by cross-products of facets
- We will actually do something different...



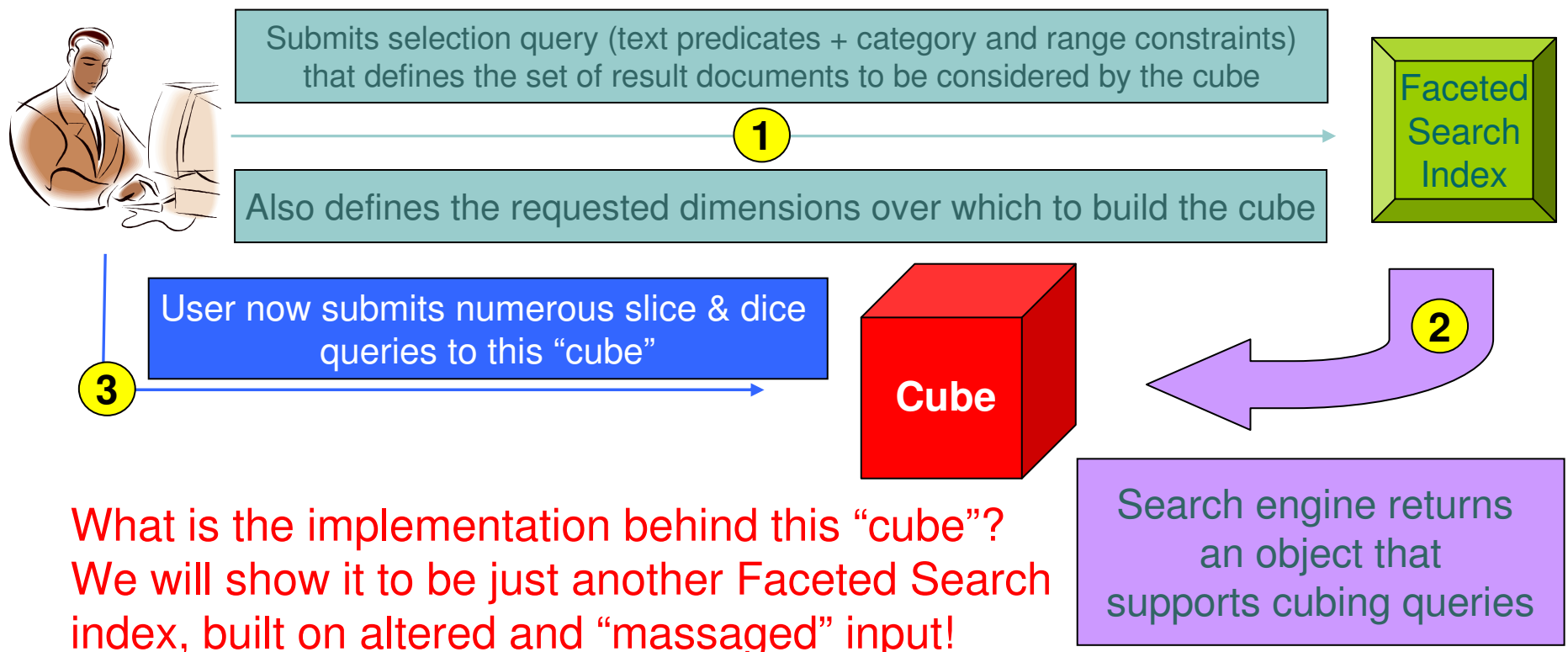
Interaction Paradigms

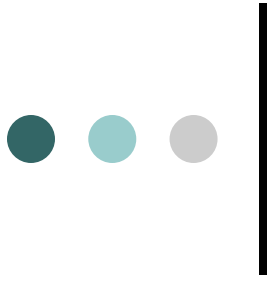
- Faceted search: ad-hoc discovery sessions of a multidimensional corpus using a mixture of keyword search and zoom-in/zoom-out browse operations
 - Typically few browsing operations per keyword query
 - Used by everyday users, in sessions lasting minutes
- OLAP Cubes: deep analysis of a fixed multidimensional corpus
 - A dataset is fixed, the cube is defined and built, and many cubing queries are executed
 - Used by professional analysts in sessions lasting days
- Mission: support cubing operations over subsets of a corpus defined by search-engine result sets

IDEA

Cubing Queries: Usage Scenario and Flow

Usage scenario: on the cube that was built over the documents and dimensions defined by some initial query, many subsequent slice-and-dice operations will be performed.





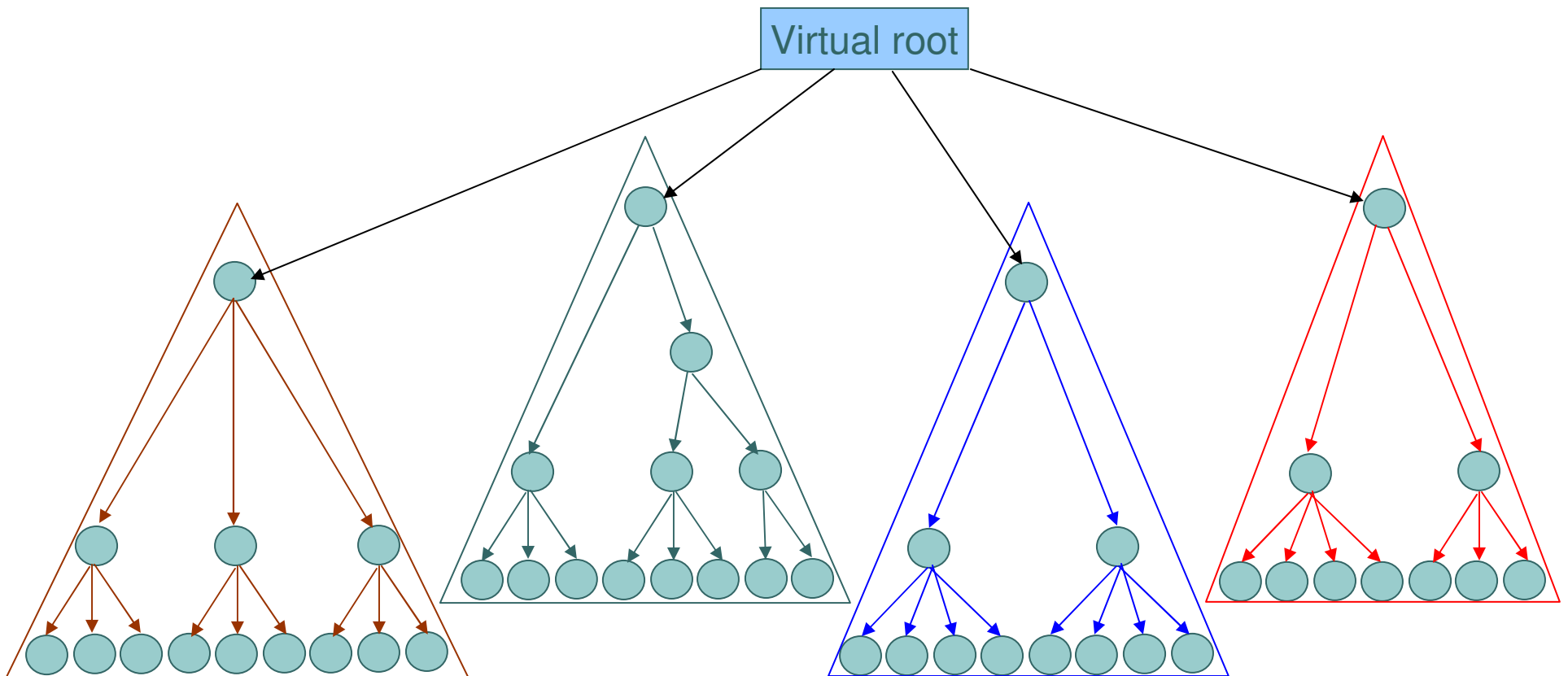
Multifaceted Search Implementation: Data Model, Indexing and Runtime

Data Model



Data Model (1)

- There's some background (latent) multidimensional hierarchy – four dimensions in this example

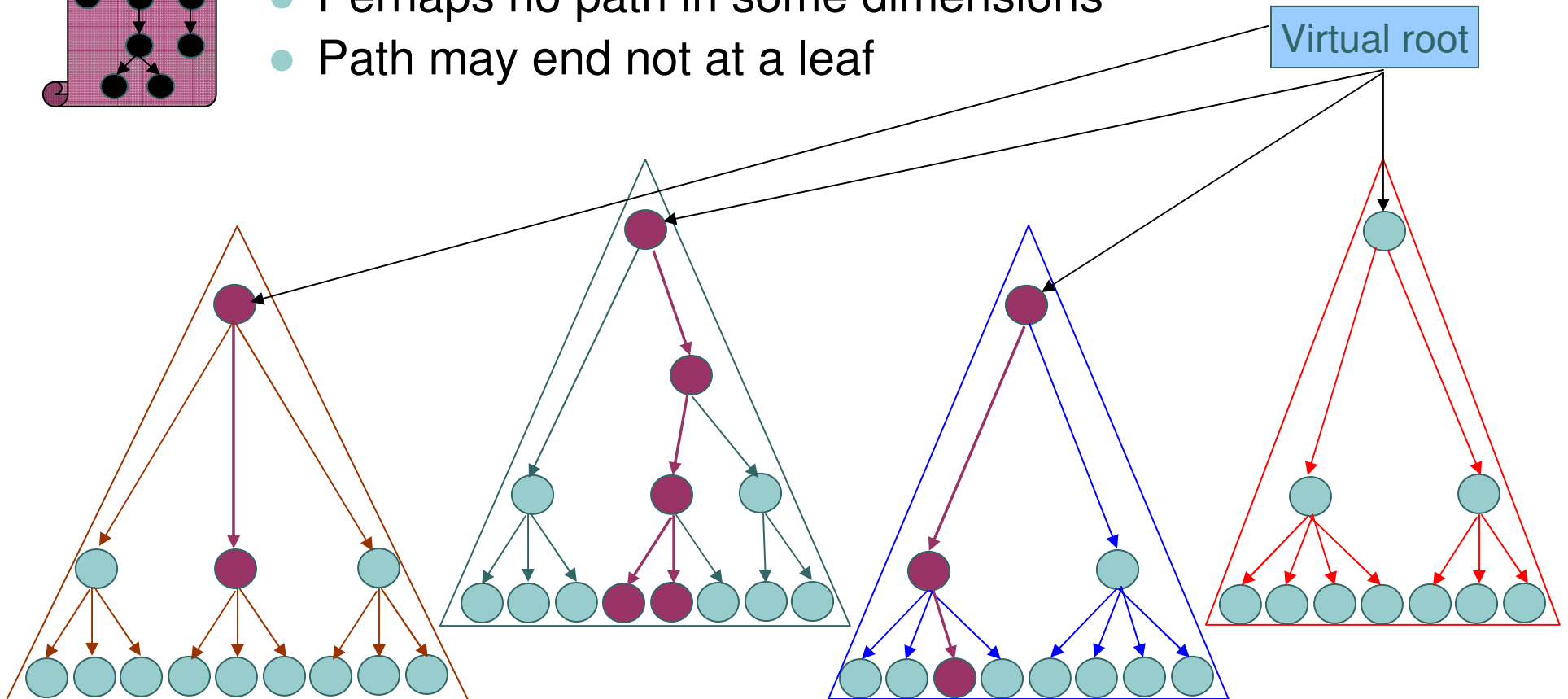
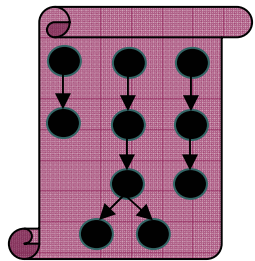


Data Model



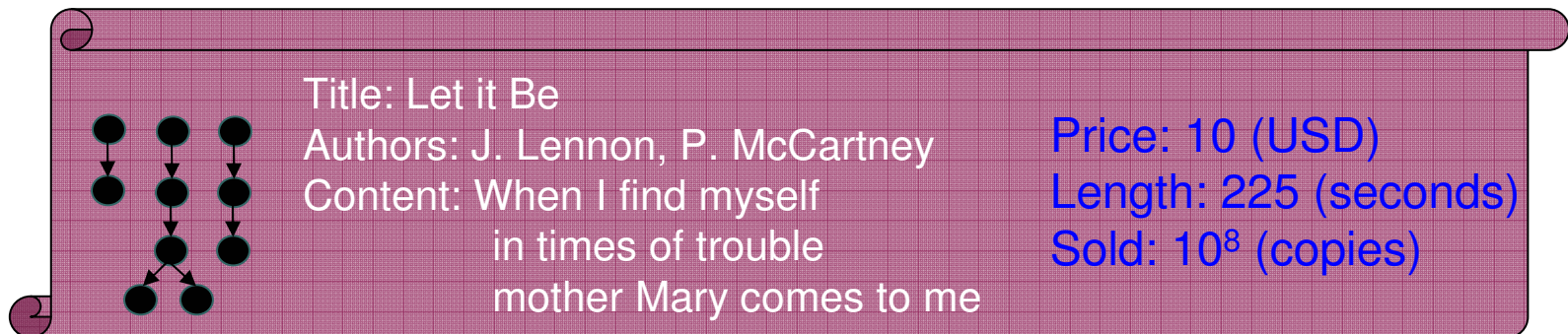
Data Model (2)

- Each document is associated with a set of paths
 - Perhaps several paths in some dimensions
 - Perhaps no path in some dimensions
 - Path may end not at a leaf



Data Model (3)

- Actually, each document can be thought of logically as composed of three parts:
 1. The collection of paths of categories/facets it is associated with
 2. Some free text, in one or more fields
 3. Several numeric attributes/fields that are associated with it





The Taxonomy Index

- Documents are ingested by the index one by one
- The **taxonomy index** component unions the set of paths associated with each document to form the overall taxonomy implied by the documents
- The taxonomy index provides the following services:
 - Maps human readable labels of each taxonomy node to some internal ID
 - Maps nodes to their ancestors
 - Maps nodes to their children



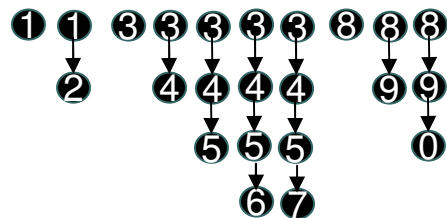
The Inverted Index (1)

- Keeps postings lists per each text token
- Keeps postings lists per each path prefix, listing all documents that are associated with the path prefix
- Keeps numeric postings lists to support range-constraint queries

1 3 8
2 4 9
5 0
6 7

Title: Let it Be
Authors: J. Lennon, P. McCartney
Content: When I find myself
in times of trouble
mother Mary comes to me

Price: 10 (USD)
Length: 225 (seconds)
Sold: 10^8 (copies)



← This document will be listed in the postings lists corresponding to these path prefixes



The Inverted Index (2)

- Keeps a special postings list that allows efficient access to the paths and numeric attributes associated with each document

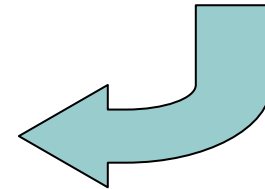
1 3 8
↓ ↓ ↓
2 4 9
↓ ↓
5 0
↓ ↓
6 7

Title: Let it Be
Authors: J. Lennon, P. McCartney
Content: When I find myself
in times of trouble
mother Mary comes to me

Price: 10 (USD)
Length: 225 (seconds)
Sold: 10^8 (copies)

The payload associated with the posting element of this document:

2	Price: 10
6	Length: 225
7	
0	Sold: 10^8



Note that only path tails are stored – higher-level nodes can be inferred using the taxonomy index

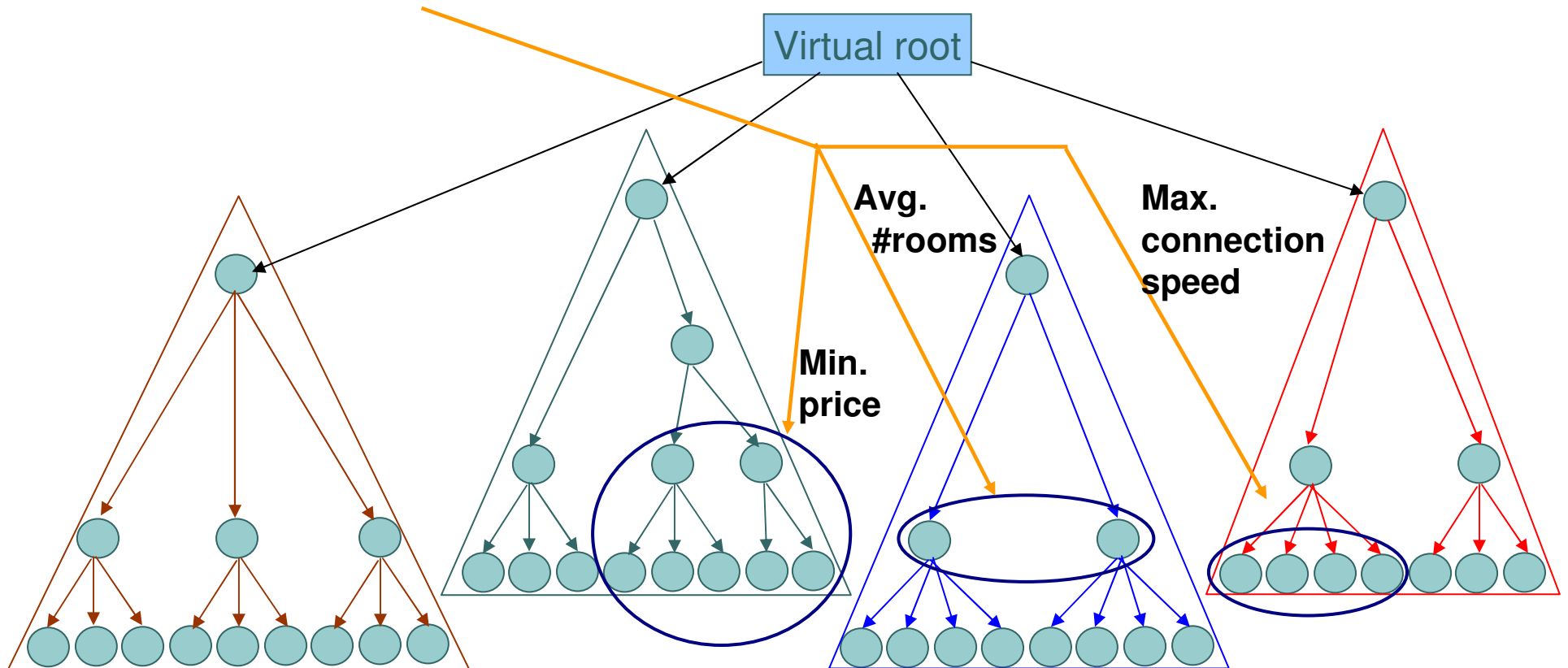
Faceted Search



Faceted Queries

Faceted queries are comprised of two parts:

1. Selection criteria: text(hotels with internet access), categories(geo:US, chain=Hilton), numerics(price<200)
2. Facet Context: sub-trees where counts/aggregations are required





Faceted Query Evaluation

1. Given the selection portion of the faceted query, identify and rank the matching documents using the index
2. Access the special per-document entry (payload) corresponding to each matching document
3. For any category included in the facet context, aggregate counts/expressions

2	Price: 10
6	Length: 225
7	Sold: 10^8

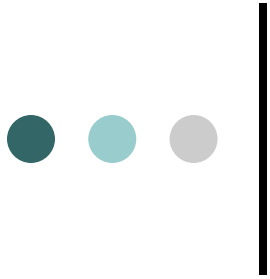
2	Price: 3
6	Length: 215
8	Sold: 10^7

4	Price: 18
5	Length: 185
7	Sold: 10^6

5	Price: 12
8	Length: 220
9	Sold: 10^8

Category:2	Count: 2	Average price: 6.5
Category:0	Count: 3	Max sold: 10^8

Note that per query, computation time scales linearly with size of result set (number of matching docs)



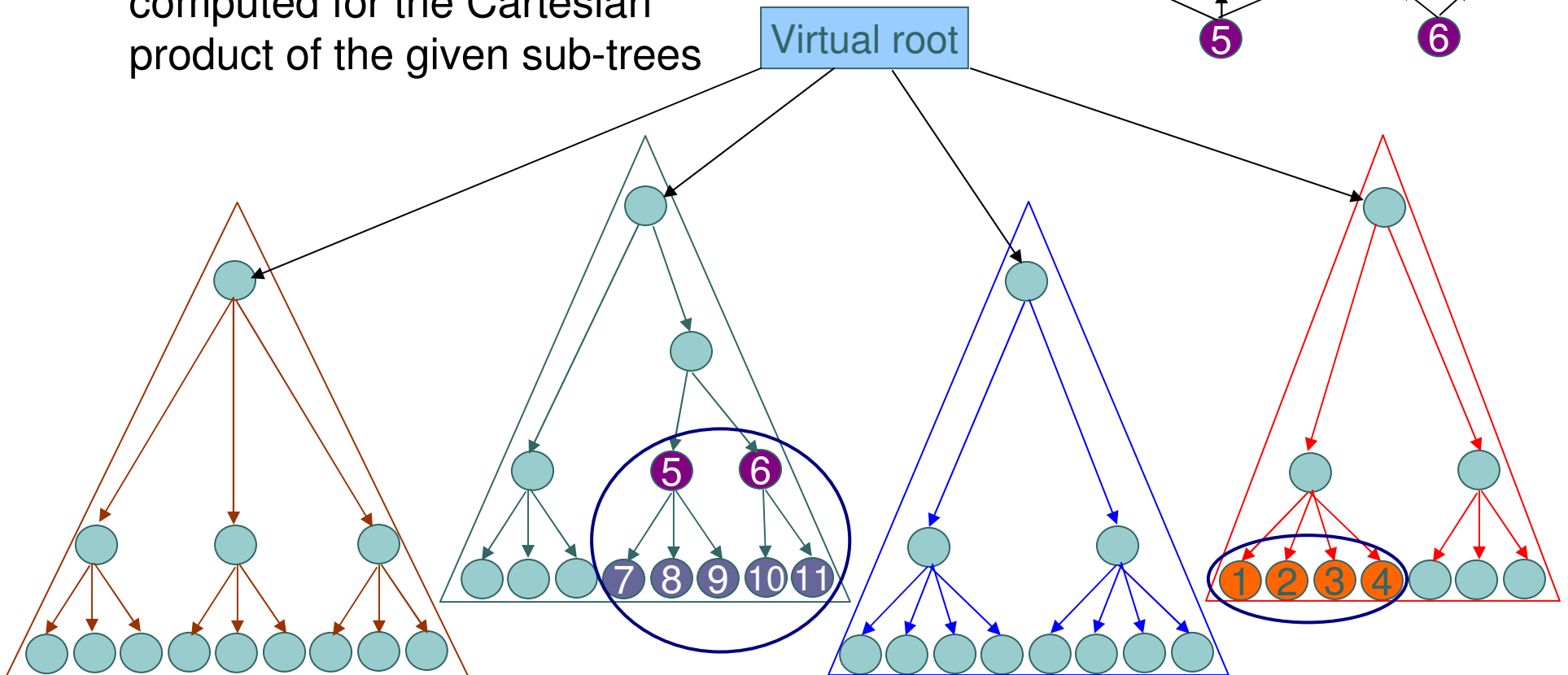
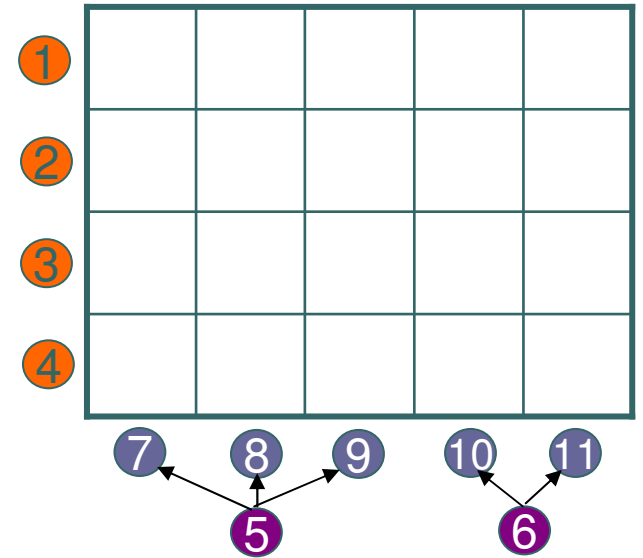
Building a Cubing Engine by a Recursive Application of Faceted Indexing

Cubing



Cubing Queries

Same as faceted queries, except that the output counts and aggregations are to be computed for the Cartesian product of the given sub-trees



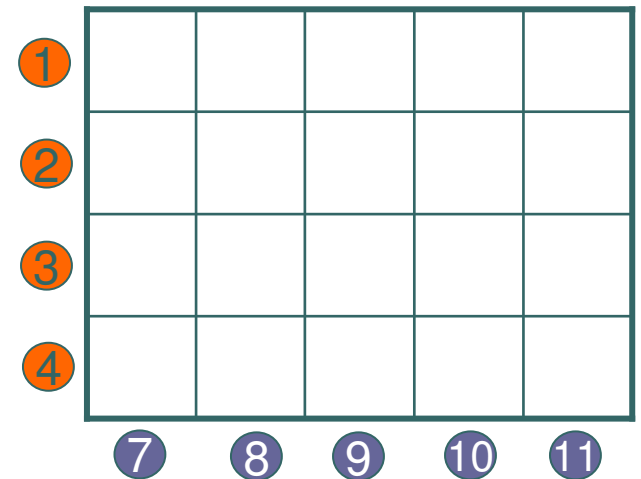
Cubing



Cubing Queries

Naïve solution:

- Every cell can be computed by a single regular query, with the appropriate category constraints
- Every line can be computed by a single faceted query, given constraints on the other dimensions
- **Problem: each such query runs in a time that is linear in the size of the result set, which is not acceptable for cubing queries**
- So let's look at a solution that trades-off some of the runtime cost with some preprocessing



Cubing

Cubing Query Evaluation

1. Given the selection portion of the initial query, identify the matching documents
2. Access the special per-document payload corresponding to each matching document
3. Extract the **Cube-DNA (CDNA)** of each document, i.e. the sequence of categories to which the document belongs and which are part of the cube
 - Example: if the cubed sub-trees are $\{1,2,3,4\} \times \{7,8,9,10,11\}$, the **CDNA** sequences of the sample documents are as follows:

2	Price: 10
6	Length: 225
7	Sold: 10^8

{2,7,10}

2	Price: 3
6	Length: 215
8	Sold: 10^7

{2,8,10}

4	Price: 18
6	Length: 185
11	Sold: 10^6

{4,11}

1	Price: 12
2	Length: 220
9	Sold: 10^8

{1,2,9,10}

4. Aggregate counts and expressions per **CDNA** sequence

CDNA Sequences - Notes

- Each CDNA sequence affects a certain set of cube cells
 - A cell may be affected by multiple CDNA sequences

2	Price: 10
6	Length: 225
7	
10	Sold: 10^8

{1,2,3,4}
x
{7,8,9,10,11}

2	Price: 3
6	Length: 215
8	
10	Sold: 10^7

CDNA {2,7,10} affects cells
{(2,7), (2,10)}

CDNA {2,8,10} affects cells
{(2,8), (2,10)}

- Documents that share the same CDNA affect exactly the same set of cube cells
- Assumption is that for large result sets, the plurality of CDNA sequences is much smaller than the size of the result set itself
 - For small result sets, it doesn't matter - **number of CDNA sequences is always bounded by number of results**

Cubing

CDNA Sequences – What Now?

- CDNA Sequences per-se aren't quite what we need, as they don't map 1:1 to cube cells
- To fill a cube cell, we need to aggregate the contributions of all CDNA sequences that affect the cell

Cell (2,7):

- Count = $20 + 61 = 81$
- Avg. Price = $(20 * 15 + 61 * 14) / 81$
- Max. Sold = $\text{Max}\{5 * 10^7, 4 * 10^7\} = 5 * 10^7$

- So what should we do?

Sequence	Count	Avg. Price	Max. Sold
{2,4,7}	20	15	$5 * 10^7$
{4,9}	100	18	$9 * 10^6$
{3,8,9}	84	21	$6 * 10^6$
{1,8,10}	123	19	$2 * 10^8$
{2,7}	61	14	$4 * 10^7$
{1,7,8}	34	25	$5 * 10^7$
{4,7,9,11}	78	12	$3 * 10^7$

Cubing

Building the Cube Structure: Indexing CDNA Sequences (1)

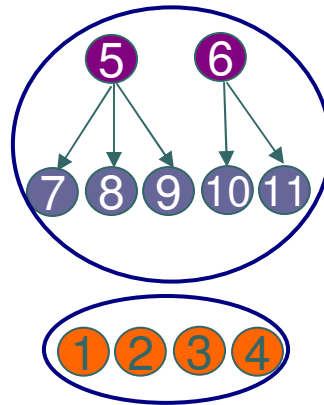
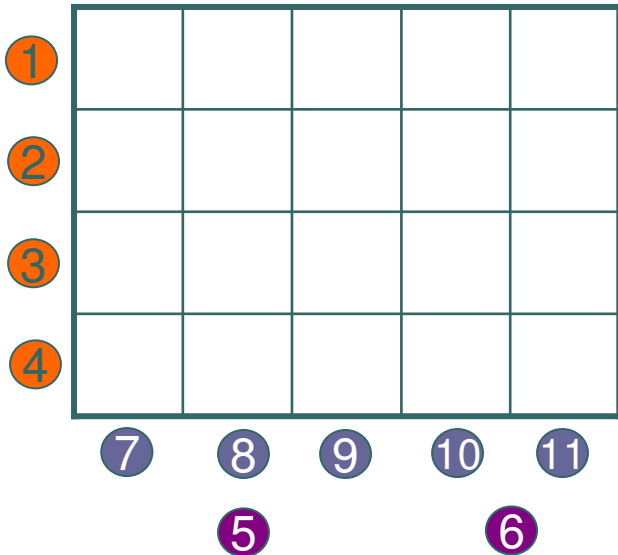
Idea:

1. Create a virtual document from each CDNA sequence:
 - No textual content
 - Categories are the DNA elements
 - Numeric attributes are the count, and the calculated expressions (measures)
2. Index the virtual documents by the same faceted search mechanism!
3. **The resulting faceted index will serve as the “cube”**

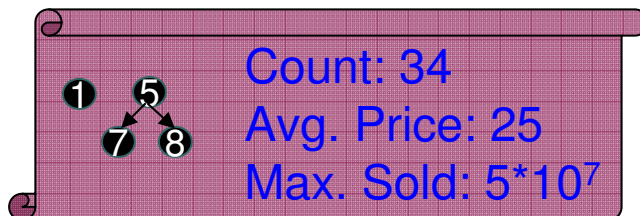
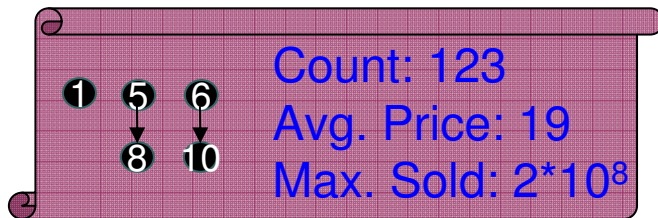
Sequence	Count	Avg. Price	Max. Sold
{2,4,7}	20	15	$5 \cdot 10^7$
{4,9}	100	18	$9 \cdot 10^6$
{3,8,9}	84	21	$6 \cdot 10^6$
{1,8,10}	123	19	$2 \cdot 10^8$
{2,7}	61	14	$4 \cdot 10^7$
{1,7,8}	34	25	$5 \cdot 10^7$
{4,7,9,11}	78	12	$3 \cdot 10^7$

Cubing

Building the Cube Structure: Indexing CDNA Sequences (2)



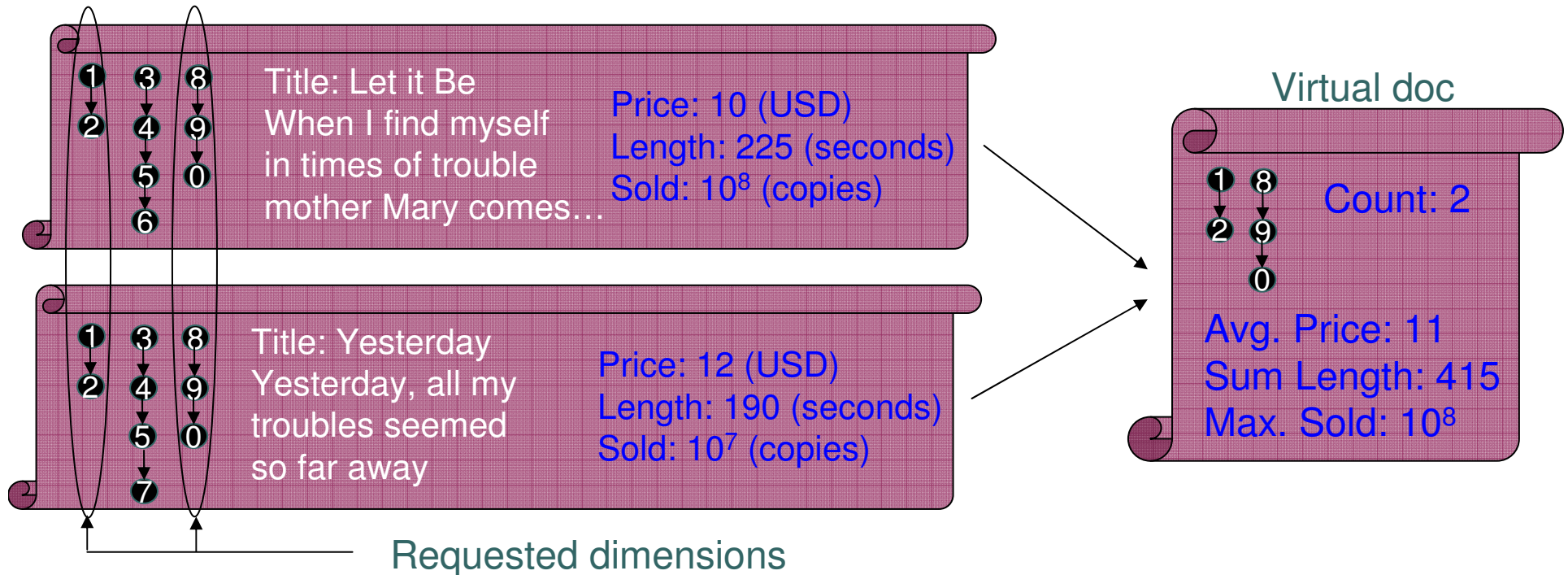
Sequence	Count	Avg. Price	Max. Sold
{2,4,7}	20	15	$5 \cdot 10^7$
{4,9}	100	18	$9 \cdot 10^6$
{3,8,9}	84	21	$6 \cdot 10^6$
{1,8,10}	123	19	$2 \cdot 10^8$
{2,7}	61	14	$4 \cdot 10^7$
{1,7,8}	34	25	$5 \cdot 10^7$
{4,7,9,11}	78	12	$3 \cdot 10^7$



Cubing

Building the Cube Structure: The Virtual CDNA Documents

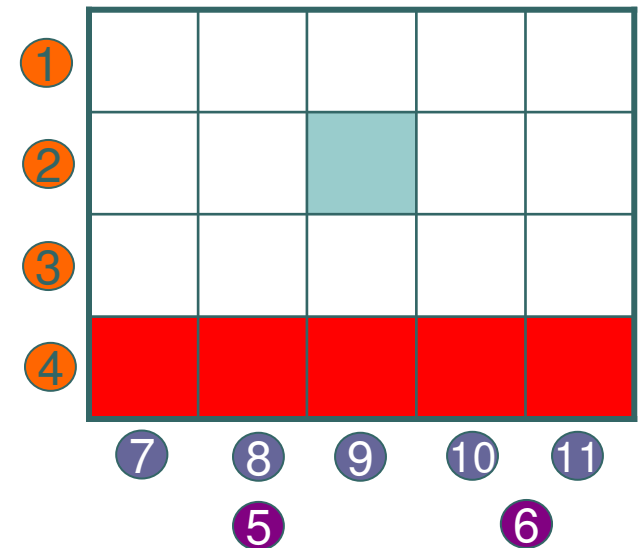
- The virtual documents represent aggregations of groups of original documents whose CDNA is identical, i.e. that are indistinguishable in the context of this cube
 - Many-to-one relationship between original documents (in OLAP terminology - facts) and CDNA virtual documents:



Cubing

The CDNA Index – Operations

- To fill the cell (2,9) – just query for documents containing **category:2 category:9**
- To fill the line (4,*), query for documents containing **category:4** and treat categories 7-11 as a one-dimensional facet
- Queries on the CDNA index run in $O(|\text{CDNA result set}|)$, i.e. the number of CDNA entries that affect the requested cell/hyper-plane



Summary



Recap of “Cubing” Scheme

- With one query over the original data (in the original faceted index), we can build a CDNA index in $O(|\text{result set}|)$
- Assumption: typically, the number of CDNA sequences per result set is not much greater than the number of populated cube cells in the requested hyper-plane
 - Number of virtual documents in CDNA index is $O(|\text{populated cells}|)$
 - Most cubes are sparse
 - Number of populated cube cells is typically much less than the number of (real) documents over which the cube is being built
- CDNA index CDNA index is yet another faceted index that supports cubing queries in $O(|\text{CDNA result set}|)$
 - It may reside on disk, in which case it is persistent and requires practically no RAM overhead – in particular, no need to materialize a cube in memory



Summary and Roadmap

- Presented a “recursive” invocation of faceted indexing that can support cubing operations
 - First invocation returns, as its result, a second index of the same type on which subsequent operations run
- Whether the performance tradeoffs implied by this approach make sense is subject to experiments
- But we like the algebraic idea nonetheless
- Use-case is real, and is yet another example of the merging between the structured, DB-oriented world and the semi-structured world of search engines

Thank You