

A large, thin, gold-colored circle is positioned behind the title bar, partially obscured by the title text.

Enhancing Bitmap Indices

Guadalupe Canahuate
The Ohio State University

Advisor: Hakan Ferhatosmanoglu

Introduction

- n Bitmap indices:
 - i Widely used in data warehouses and read-only domains
 - i Implemented in commercial DBMS such as Oracle, Informix and Sybase
 - i Extremely efficient for execution of point and range queries

Applications

- n Data warehouses
- n Scientific data
 - i **High-energy physics:** Simulations are continuously run, and notable events are stored with all the details.
 - i **Climate modeling:** sensor data.
 - i **Astro-physics:** telescopes devoted for observations.
- n Visualization applications



Bitmap Indices

- n Data is quantized into b categories using b bits
- n Each tuple is encoded based on which category its attribute belongs to:
 - i Bitmap Equality Encoding (BEE)
(Projection Index, Binning):
suited for point queries
 - i Bitmap Range Encoding (BRE):
suited for range queries
- n Fast Bitwise Operations over bit vectors (AND, OR, NOT, XOR)

Value	EE			RE		
	1	2	3	1	2	3
1	1	0	0	1	1	1
2	0	1	0	0	1	1
3	0	0	1	0	0	1

Bitmap Compression

- n Run-length encoders
 - i Runs of 0s (0, run-count)
 - i Byte-aligned Bitmap Code (BBC) [Oracle '94]
 - i Word-Aligned Hybrid Code (WAH) [LBNL '02]

 - i No need to decompress for query processing
 - i Full scan of the bitmap is needed

Our Goal

- n Enhance bitmap indices:
 - i Reduce Index Size (ICDE'05, TKDE sub)
 - n Improve compression
 - i Direct access over compressed bitmaps (VLDB'06)
 - n Bitmap Hashing
 - i Handle Missing Data (EDBT'06)
 - n Query support / semantics
 - i Handle data updates (SSDBM'07)
 - n Appends of new data

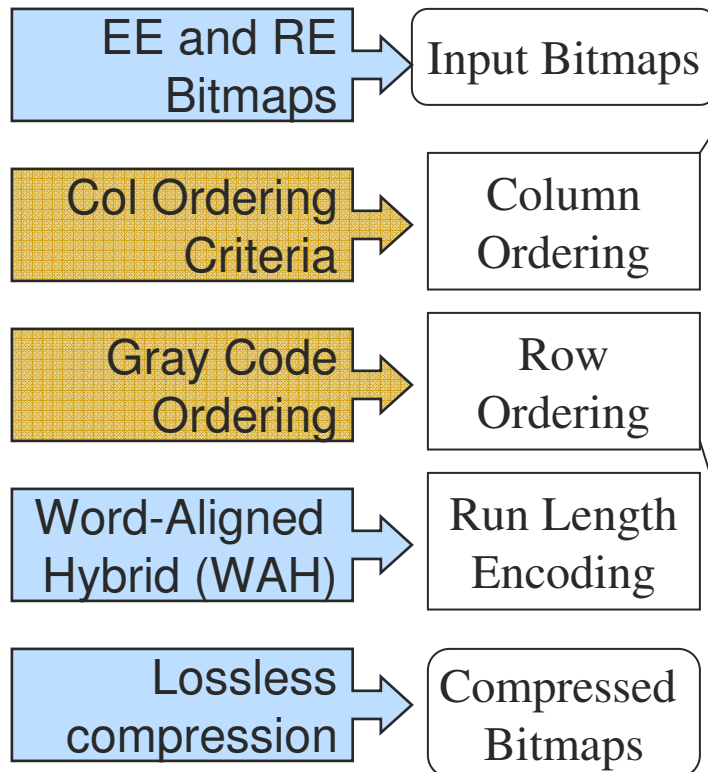
Our Goal

- n Enhance bitmap indices:
 - i Reduce Index Size (ICDE'05, TKDE sub)
 - n Improve compression
 - i Direct access over compressed bitmaps (VLDB'06)
 - n Bitmap Hashing
 - i Handle Missing Data (EDBT'06)
 - n Query support / semantics
 - i Handle data updates (SSDBM'07)
 - n Appends of new data

Improving Compression by Data Reorganization

- n NP-Complete
- n Most TSP heuristics are ineffective
- n Minimize the hamming distance of adjacent numbers
- n Gray-code: A space filling-curve for hamming space
- n A scalable, in-place algorithm to gray-sort the tuples

Reordering Framework



Reorder ($A, start, end, b$)

1: $i \leftarrow start$

2: $j \leftarrow end$

3: **while** $i < j$ **do**

4: Decrement j until $S(j, Cols[b]) = 0$

5: Increment i until $S(i, Cols[b]) = 1$

6: **if** $i < j$ **then**

7: Swap the i^{th} and j^{th} tuples on the table

8: **end if**

9: **end while**

10: **if** $b < no_of_bits$ **then**

11: $Cols[b+1] = nextCol();$

12: Reorder ($A, start, j, b+1$)

13: Reorder ($A, j+1, end, b+1$)

14: Reverse ($j+1, end$)

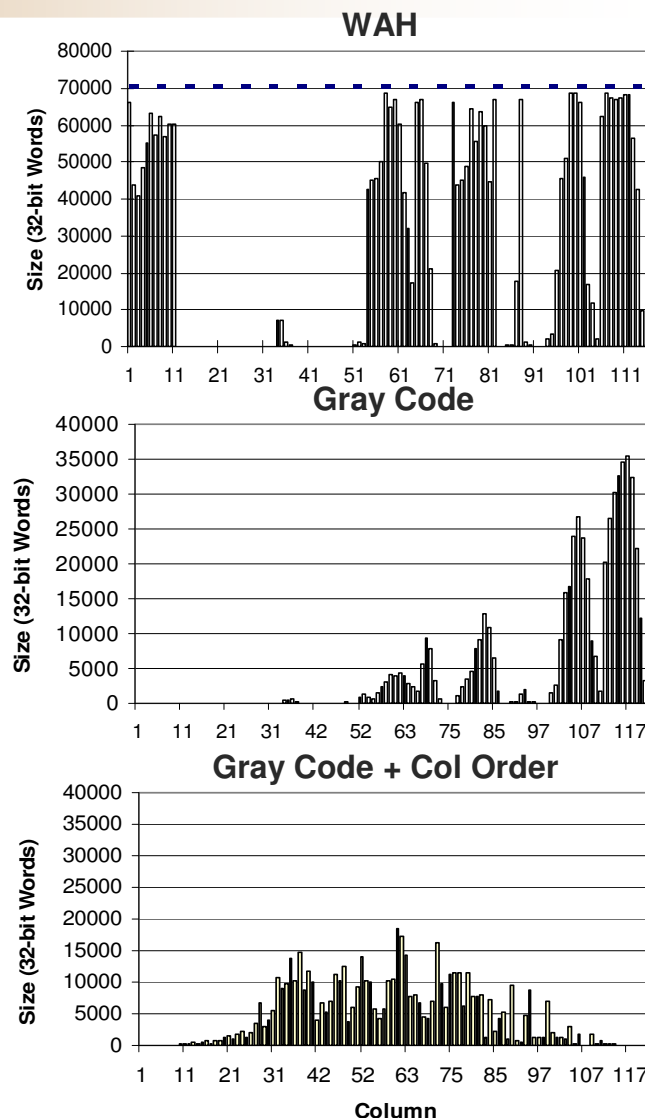
15: **end if**

Column Ordering Criteria

- n Using Bitmap Data
 - i Number of Set Bits
 - n Static: Most 1s Asc/Desc
 - n Dynamic: Most 1s in the Longest Run
 - i Compressibility
 - n Static: $c_i = \left| \frac{n}{2} - s_i \right|$
 - n Dynamic: Weighted Sum of Compressibility over the Gray code segments
- n Using Query Workloads
 - i Most frequent accessed columns first

Experimental Results

- n 2-10x improvement in compression over already compressed bitmaps
- n 4-7x improvement in query execution time



Our Goal

- n Enhance bitmap indices:
 - i Reduce Index Size (ICDE'05, TKDE sub)
 - n Improve compression
 - i **Direct access over compressed bitmaps (VLDB'06)**
 - n Bitmap Hashing
 - i Handle Missing Data (EDBT'06)
 - n Query support / semantics
 - i Handle data updates (SSDBM'07)
 - n Appends of new data

Compression and Direct Access

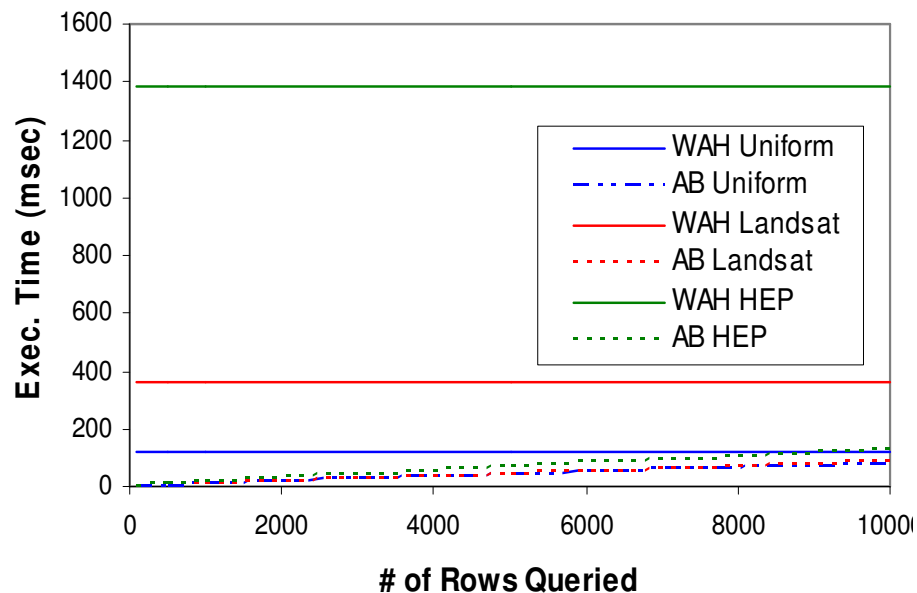
- n Bitmaps need to be compressed:
 - i Row numbers do not longer correspond to the bit position in the bitmap
- n Cannot pin-point a row
 - i Need to scan the bitmap
- n Enable direct access over the compressed bitmaps
- n Trade-off accuracy vs. efficiency
 - i No false negatives

Direct Access Over Compressed Bitmaps

- n Hashing/Bloom Filters
 - i A 2^m bit array indexed using k independent hash functions
 - i False positives can happen, but false negatives cannot
 - i Only the set bits are inserted into the AB
 - i Three levels of encoding:
 - n Per table, per attribute, per bitmap column
 - i Precision estimation

Direct Access Over Compressed Bitmaps

- n Always set the array size to be smaller than WAH bitmap
- n Execution time depends on the number of **rows queried**, not in the size of the bitmaps
- n For queries over less than ~15% of the rows, execution time is up to 3 orders of magnitude faster than WAH



Our Goal

- n Enhance bitmap indices:
 - i Reduce Index Size (ICDE'05, TKDE sub)
 - n Improve compression
 - i Direct access over compressed bitmaps (VLDB'06)
 - n Bitmap Hashing
 - i **Handle Missing Data (EDBT'06)**
 - n **Query support / semantics**
 - i Handle data updates (SSDBM'07)
 - n Appends of new data

Handling missing data

- n Incomplete databases are common in all research and industry domains
- n Index performance degrades in the presence of missing data
- n Missing data map to a distinguished category
- n Query Semantics:
 - i Missing is a match
 - n In a patient database retrieve the records whose first name is “John” and middle name “Paul”.
 - i Missing is not a match
 - n In a census database retrieve the interviewee that answered “C” in question 4 and “A” in question 8

Handling Missing Data

		Equality Encoded						Range Encoded					
Record	Value	B _{1,0}	B _{1,1}	B _{1,2}	B _{1,3}	B _{1,4}	B _{1,5}	B _{1,0}	B _{1,1}	B _{1,2}	B _{1,3}	B _{1,4}	B _{1,5}
1	5	0	0	0	0	0	1	0	0	0	0	0	1
2	2	0	0	1	0	0	0	0	0	1	1	1	1
3	3	0	0	0	1	0	0	0	0	0	1	1	1
4	missing	1	0	0	0	0	0	1	1	1	1	1	1
5	4	0	0	0	0	1	0	0	0	0	0	1	1
6	5	0	0	0	0	0	1	0	0	0	0	0	1
7	1	0	1	0	0	0	0	0	1	1	1	1	1
8	3	0	0	0	1	0	0	0	0	0	1	1	1
9	missing	1	0	0	0	0	0	1	1	1	1	1	1
10	2	0	0	1	0	0	0	0	0	1	1	1	1

Query Processing with Missing Data

Equality Encoded

$$v_1 \leq A_i \leq v_2 =$$

Missing is a Match

$$\begin{cases} \left(\bigcup_{j=v_1}^{v_2} B_{i,j} \right) \vee B_{i,0} & \text{if } v_2 - v_1 \leq \lfloor C_i/2 \rfloor \\ \bigcup_{j=1}^{v_1-1} B_{i,j} \vee \bigcup_{j=v_2+1}^{C_i} B_{i,j} & \text{otherwise} \end{cases}$$

Missing is not a Match

$$\begin{cases} \left(\bigcup_{j=v_1}^{v_2} B_{i,j} \right) & \text{if } v_2 - v_1 \leq \lfloor C_i/2 \rfloor \\ \bigcup_{j=1}^{v_1-1} B_{i,j} \vee \bigcup_{j=v_2+1}^{C_i} B_{i,j} \oplus B_{i,0} & \text{otherwise} \end{cases}$$

Range Encoded

$$\begin{cases} B_{i,1} & \text{if } v_2 = v_1 = 1 \\ \frac{B_{i,v_1} \oplus B_{i,v_1-1} \vee B_{i,0}}{B_{i,C_i-1} \vee B_{i,0}} & \text{if } 1 < v_1 = v_2 < C_i \\ \frac{B_{i,v_1-1} \vee B_{i,0}}{B_{i,v_2}} & \text{if } 1 < v_1 < C_i, v_2 = C_i \\ B_{i,v_2} & \text{if } v_1 = 1, 1 < v_2 < C_i \\ B_{i,v_2} \oplus B_{i,v_1-1} \vee B_{i,0} & \text{otherwise} \end{cases}$$

$$\begin{cases} B_{i,1} \oplus B_{i,0} & \text{if } v_2 = v_1 = 1 \\ \frac{B_{i,v_1} \oplus B_{i,v_1-1}}{B_{i,C_i-1}} & \text{if } 1 < v_1 = v_2 < C_i \\ \frac{B_{i,v_1-1}}{B_{i,v_2} \oplus B_{i,0}} & \text{if } 1 < v_1 < C_i, v_2 = C_i \\ B_{i,v_2} \oplus B_{i,0} & \text{if } v_1 = 1, 1 < v_2 < C_i \\ B_{i,v_2} \oplus B_{i,v_1-1} & \text{otherwise} \end{cases}$$

Our Goal

- n Enhance bitmap indices:
 - i Reduce Index Size (ICDE'05, TKDE sub)
 - n Improve compression
 - i Direct access over compressed bitmaps (VLDB'06)
 - n Bitmap Hashing
 - i Handle Missing Data (EDBT'06)
 - n Query support / semantics
 - i **Handle data updates (SSDBM'07)**
 - n **Appends of new data**

Handling Data Updates

- n Look for this project in the poster session... J

Conclusion

- n Enhance bitmap indices:
 - i Improve performance
 - n Better compression
 - n Direct access over compressed bitmap
 - i Overcome limitations
 - n Handle updates efficiently
 - n Support more types of queries
- n Extend the applicability of bitmaps to more domains that can benefit from bitmaps good performance

Ongoing Work

- n Support new types of queries
 - i Similarity Queries
- n Integrate bitmaps to the coastal monitoring system (NSF Cyber Infrastructure Project)
- n Bitmaps on emerging hardware technologies

Questions and Comments

n Thank you!



Email: canahuat@cse.ohio-state.edu