

Column-Oriented Database Systems

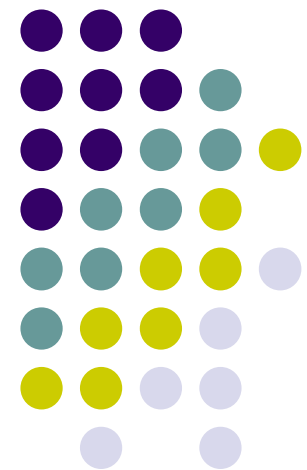
VLDB
2009
Tutorial



Part 1: Stavros Harizopoulos (HP Labs)

Part 2: Daniel Abadi (Yale)

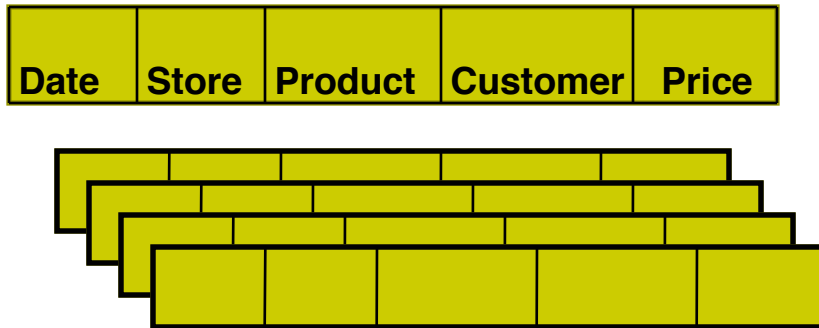
Part 3: Peter Boncz (CWI)





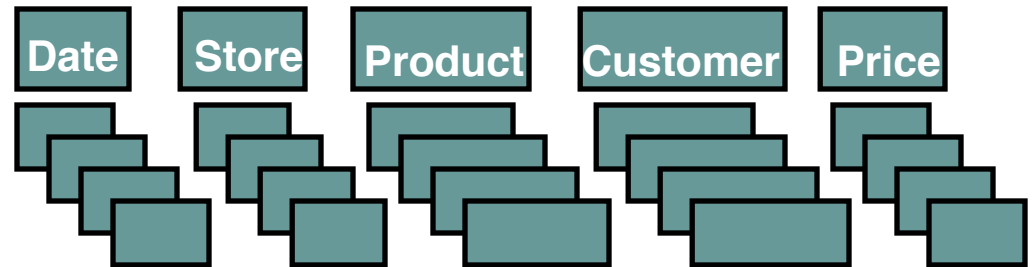
What is a column-store?

row-store



- + easy to add/modify a record
- might read in unnecessary data

column-store



- + only need to read in relevant data
- tuple writes require multiple accesses

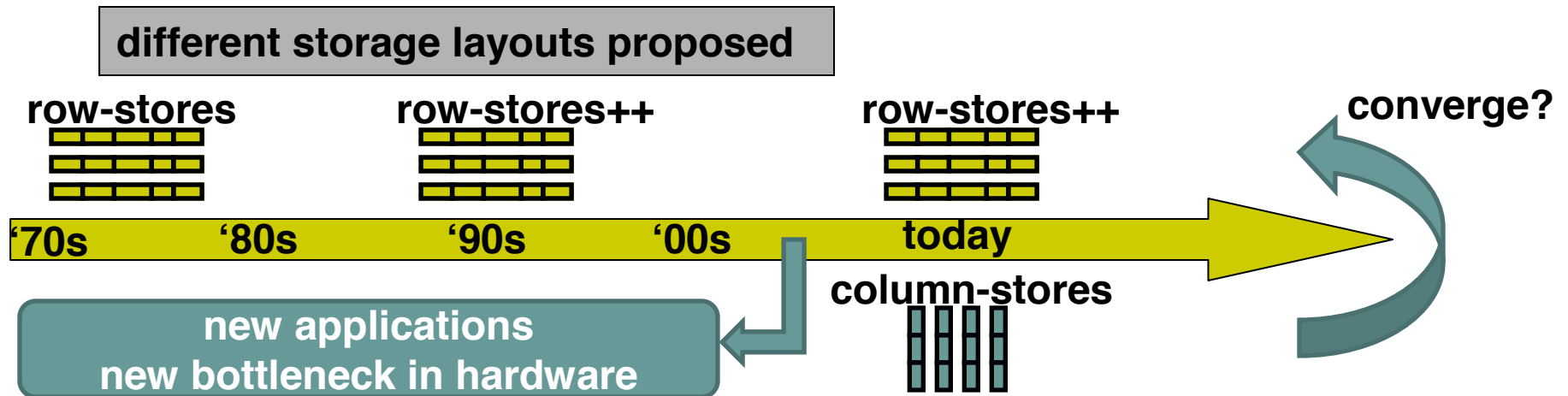
=> suitable for read-mostly, read-intensive, large data repositories





Are these two fundamentally different?

- 1 The only fundamental difference is the storage layout
- 1 However: we need to look at the big picture



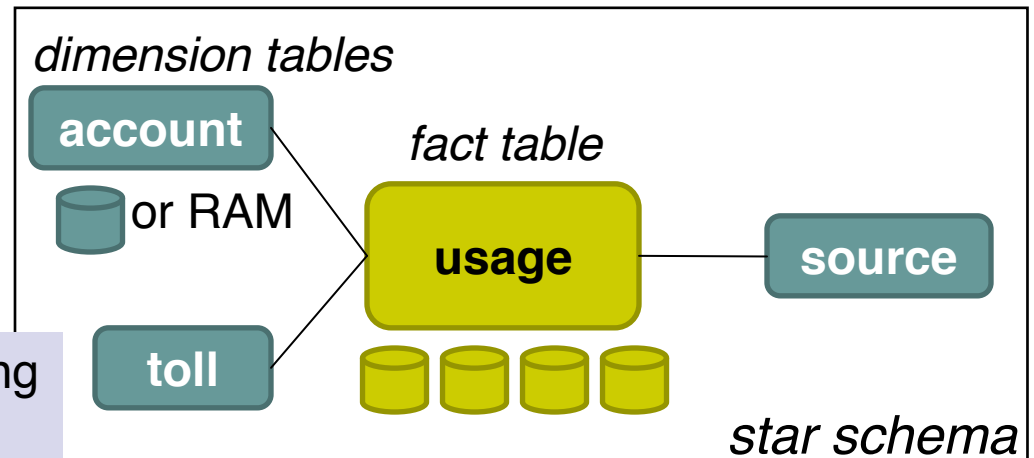


Telco Data Warehousing example

1 Typical DW installation

1 Real-world example

“One Size Fits All? - Part 2: Benchmarking Results” Stonebraker et al. CIDR 2007



QUERY 2

```
SELECT account.account_number,  
sum (usage.toll_airtime),  
sum (usage.toll_price)  
FROM usage, toll, source, account  
WHERE usage.toll_id = toll.toll_id  
AND usage.source_id = source.source_id  
AND usage.account_id = account.account_id  
AND toll.type_ind in ('AE', 'AA')  
AND usage.toll_price > 0  
AND source.type != 'CIBER'  
AND toll.rating_method = 'IS'  
AND usage.invoice_date = 20051013  
GROUP BY account.account_number
```

	<i>Column-store</i>	<i>Row-store</i>
Query 1	2.06	300
Query 2	2.20	300
Query 3	0.09	300
Query 4	5.24	300
Query 5	2.88	300

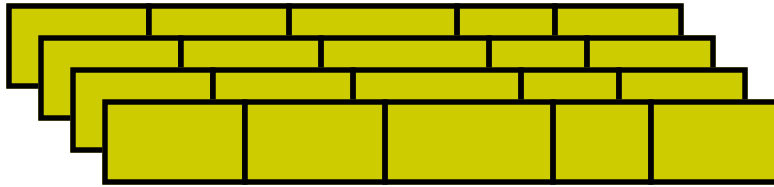
Why? Three main factors (next slides)



Telco example explained (1/3): *read efficiency*



row store



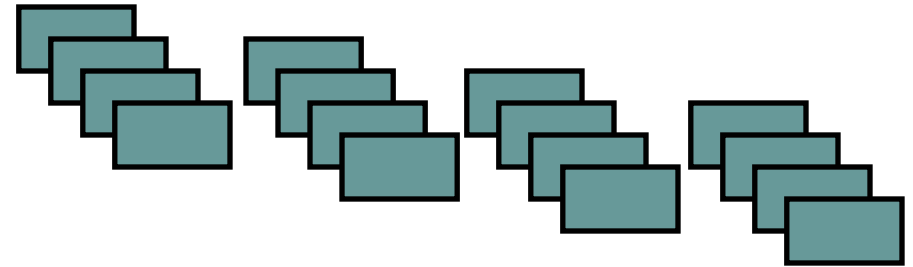
read pages containing entire rows

one row = 212 columns!

is this typical? (it depends)

What about vertical partitioning?
(it does not work with ad-hoc queries)

column store



read only columns needed

in this example: 7 columns

caveats:

- “select * ” not any faster
- clever disk prefetching
- clever tuple reconstruction





Telco example explained (2/3): *compression efficiency*

- 1 Columns compress better than rows
 - 1 Typical row-store compression ratio 1 : 3
 - 1 Column-store 1 : 10
- 1 Why?
 - 1 Rows contain values from different domains
=> more entropy, difficult to dense-pack
 - 1 Columns exhibit significantly less entropy
 - 1 Examples:

Male, Female, Female, Female, Male
1998, 1998, 1999, 1999, 1999, 2000
 - 1 Caveat: CPU cost (use lightweight compression)





Telco example explained (3/3): *sorting & indexing efficiency*

- 1 Compression and dense-packing free up space
 - 1 Use multiple overlapping column collections
 - 1 Sorted columns compress better
 - 1 Range queries are faster
 - 1 Use sparse clustered indexes

**What about heavily-indexed row-stores?
(works well for single column access,
cross-column joins become increasingly expensive)**





Additional opportunities for column-stores

- 1 Block-tuple / vectorized processing
 - 1 Easier to build block-tuple operators
 - 1 Amortizes function-call cost, improves CPU cache performance
 - 1 Easier to apply vectorized primitives
 - 1 Software-based: bitwise operations
 - 1 Hardware-based: SIMD
- 1 Opportunities with compressed columns
 - 1 *Avoid* decompression: operate directly on compressed
 - 1 *Delay* decompression (and tuple reconstruction)
 - 1 Also known as: *late materialization*
- 1 Exploit columnar storage in other DBMS components
 - 1 Physical design (both static and dynamic)

Part 3

**more
in Part 2**

See: *Database Cracking*, from CWI





MonetDB (more in Part 3)

- 1 Late 1990s, CWI: Boncz, Manegold, and Kersten
- 1 Motivation:
 - 1 Main-memory
 - 1 Improve computational efficiency by avoiding expression interpreter
 - 1 DSM with virtual IDs natural choice
 - 1 Developed new query execution algebra
- 1 Initial contributions:
 - 1 Pointed out memory-wall in DBMSs
 - 1 Cache-conscious projections and joins
 - 1 ...





2005: the (re)birth of column-stores

- 1 New hardware and application realities
 - 1 Faster CPUs, larger memories, disk bandwidth limit
 - 1 Multi-terabyte Data Warehouses
- 1 New approach: combine several techniques
 - 1 Read-optimized, fast multi-column access, disk/CPU efficiency, light-weight compression
- 1 C-store paper:
 - 1 First comprehensive design description of a column-store
- 1 MonetDB/X100
 - 1 “proper” disk-based column store
- 1 Explosion of new products





Applications for column-stores

- 1 Data Warehousing
 - 1 High end (clustering)
 - 1 Mid end/Mass Market
 - 1 Personal Analytics
- 1 Data Mining
 - 1 E.g. Proximity
- 1 Google BigTable
- 1 RDF
 - 1 Semantic web data management
- 1 Information retrieval
 - 1 Terabyte TREC
- 1 Scientific datasets
 - 1 SciDB initiative
 - 1 SLOAN Digital Sky Survey on MonetDB





List of column-store systems

- 1 Cantor (history)
- 1 Sybase IQ
- 1 SenSage (former Addamark Technologies)
- 1 Kdb
- 1 1010data
- 1 MonetDB
- 1 C-Store/Vertica
- 1 X100/VectorWise
- 1 KickFire
- 1 SAP Business Accelerator
- 1 Infobright
- 1 ParAccel
- 1 Exasol





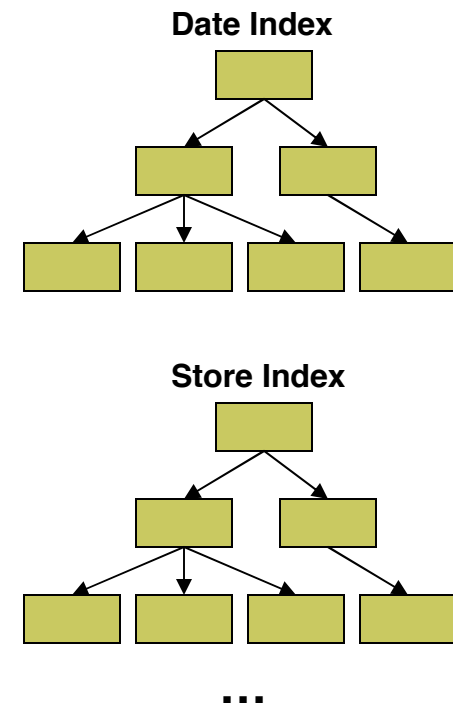
Simulate a Column-Store inside a Row-Store

Date	Store	Product	Customer	Price
01/01	BOS	Table	Mesa	\$20
01/01	NYC	Chair	Lutz	\$13
01/01	BOS	Bed	Mudd	\$79

**Option A:
Vertical Partitioning**

Date		Store		Product		Customer		Price	
TID	Value	TID	Value	TID	Value	TID	Value	TID	Value
1	01/01	1	BOS	1	Table	1	Mesa	1	\$20
2	01/01	2	NYC	2	Chair	2	Lutz	2	\$13
3	01/01	3	BOS	3	Bed	3	Mudd	3	\$79

**Option B:
Index Every Column**





Simulate a Column-Store inside a Row-Store

Date	Store	Product	Customer	Price
01/01	BOS	Table	Mesa	\$20
01/01	NYC	Chair	Lutz	\$13
01/01	BOS	Bed	Mudd	\$79

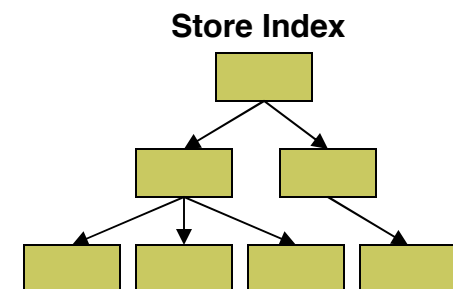
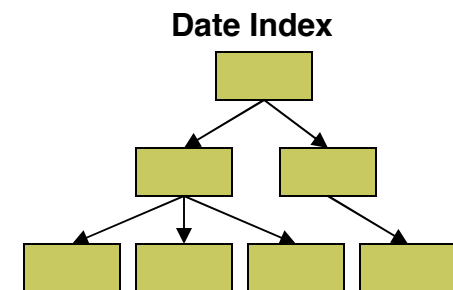
Option A: Vertical Partitioning

Date		
Value	StartPos	Length
01/01	1	3

Can explicitly run-length encode date

“Teaching an Old Elephant New Tricks.”
Bruno, CIDR 2009.

Option B: Index Every Column



...





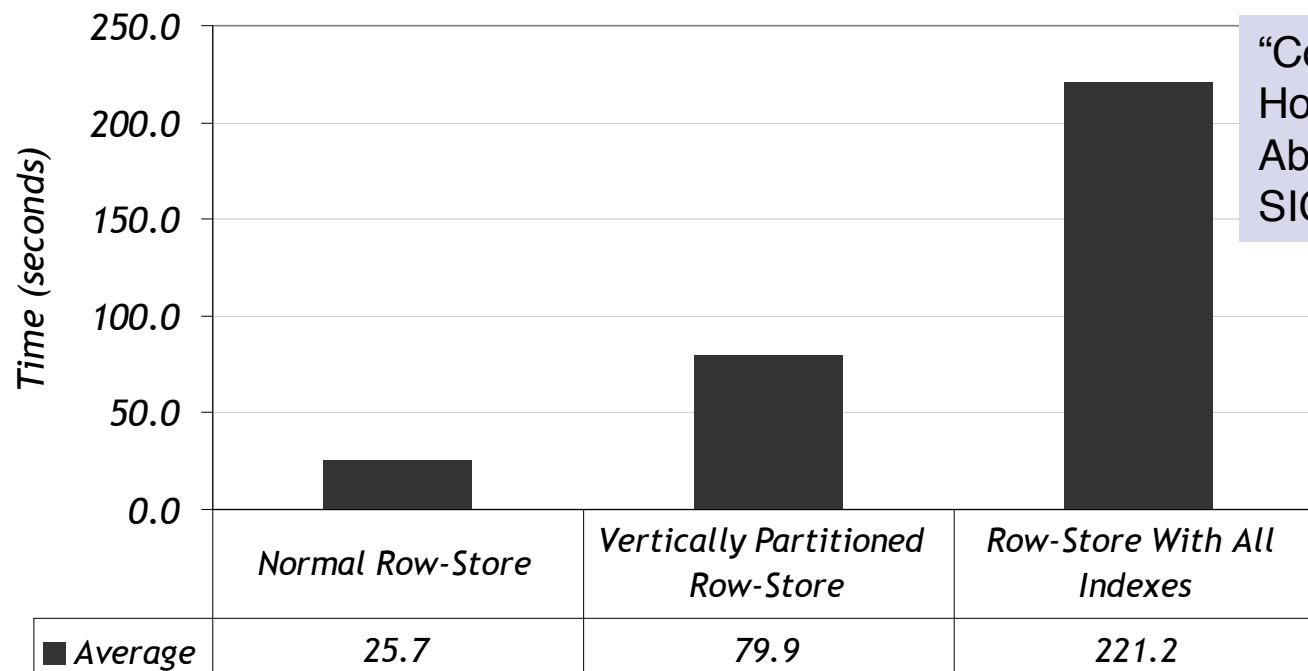
Experiments

1 Star Schema Benchmark (SSBM)

Adjoined Dimension Column Index (ADC Index) to Improve Star Schema Query Performance”. O’Neil et. al. ICDE 2008.

1 Implemented by professional DBA

1 Original row-store plus 2 column-store simulations on same row-store product



“Column-Stores vs Row-Stores: How Different are They Really?” Abadi, Hachem, and Madden. SIGMOD 2008.



Column-Oriented Database Systems

VLDB
2009
Tutorial

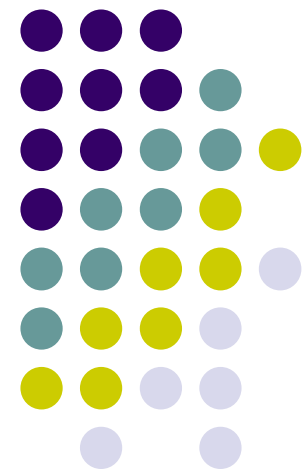


Compression

“Super-Scalar RAM-CPU Cache Compression”
Zukowski, Heman, Nes, Boncz, ICDE’06

“Integrating Compression and Execution in Column-Oriented Database Systems” Abadi, Madden, and Ferreira, SIGMOD ’06

•Query optimization in compressed database systems” Chen, Gehrke, Korn, SIGMOD’01





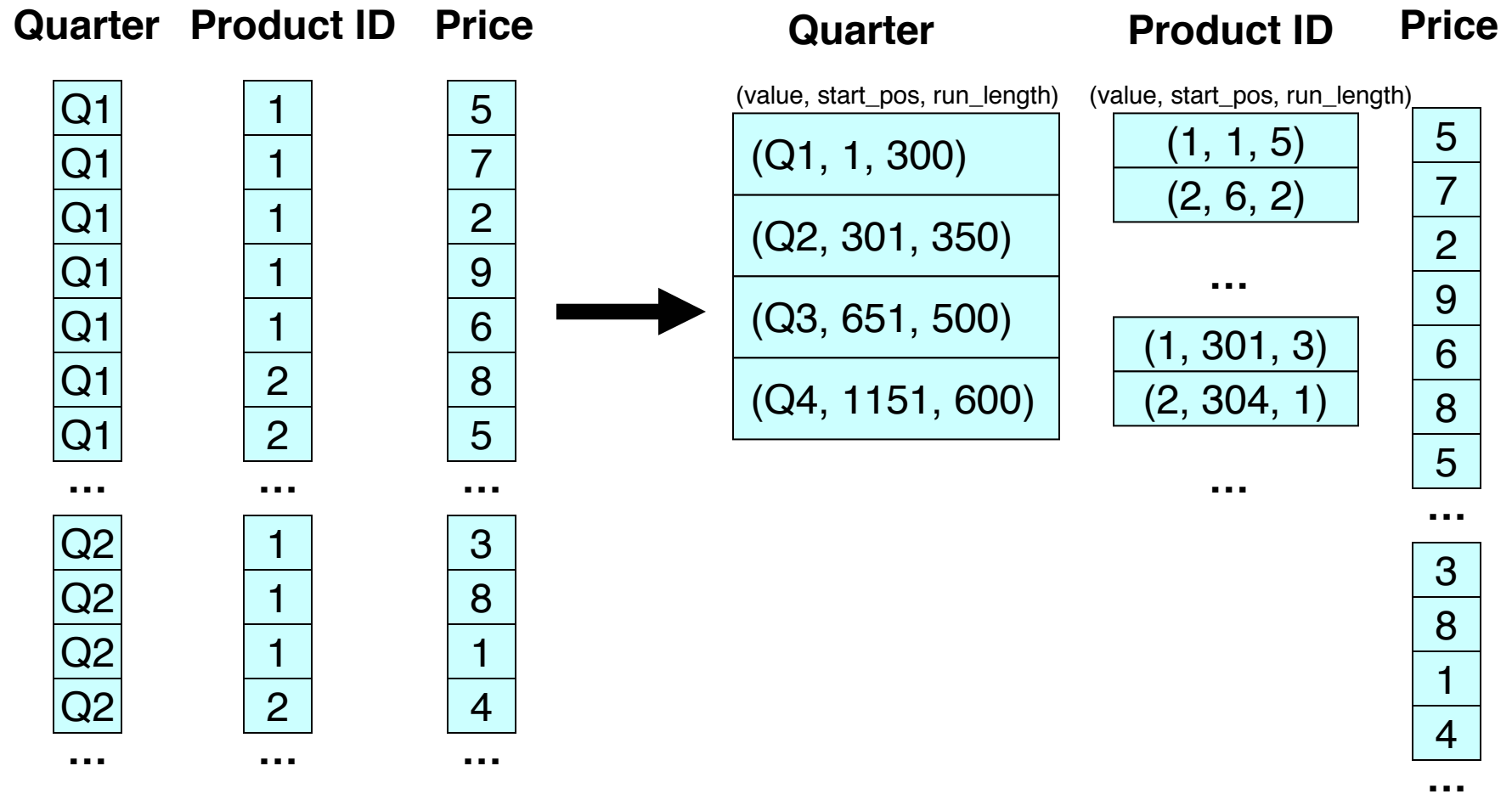
Compression

- 1 **Trades I/O for CPU**
- 1 **Increased column-store opportunities:**
 - 1 **Higher data value locality in column stores**
 - 1 **Techniques such as run length encoding far more useful**
 - 1 **Can use extra space to store multiple copies of data in different sort orders**





Run-length Encoding





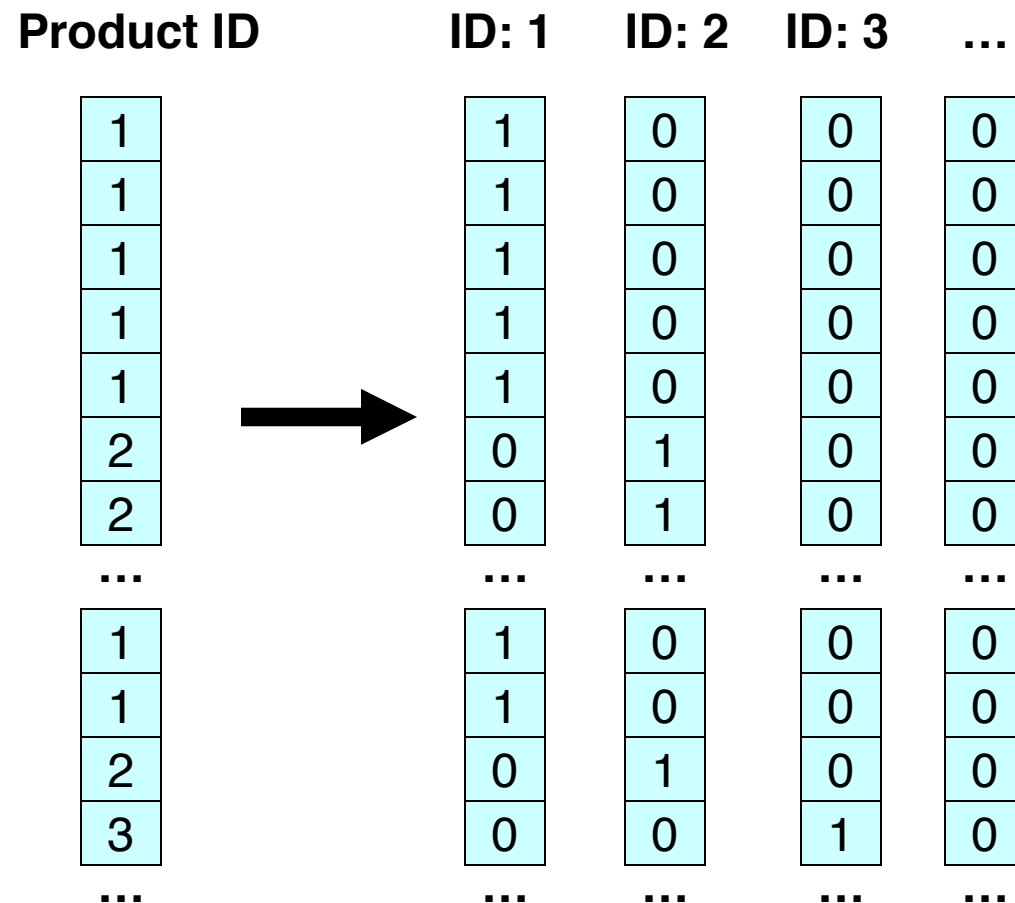
Bit-vector Encoding

- 1 For each unique value, v , in column c , create bit-vector b

- 1 $b[i] = 1$ if $c[i] = v$

- 1 Good for columns with few unique values

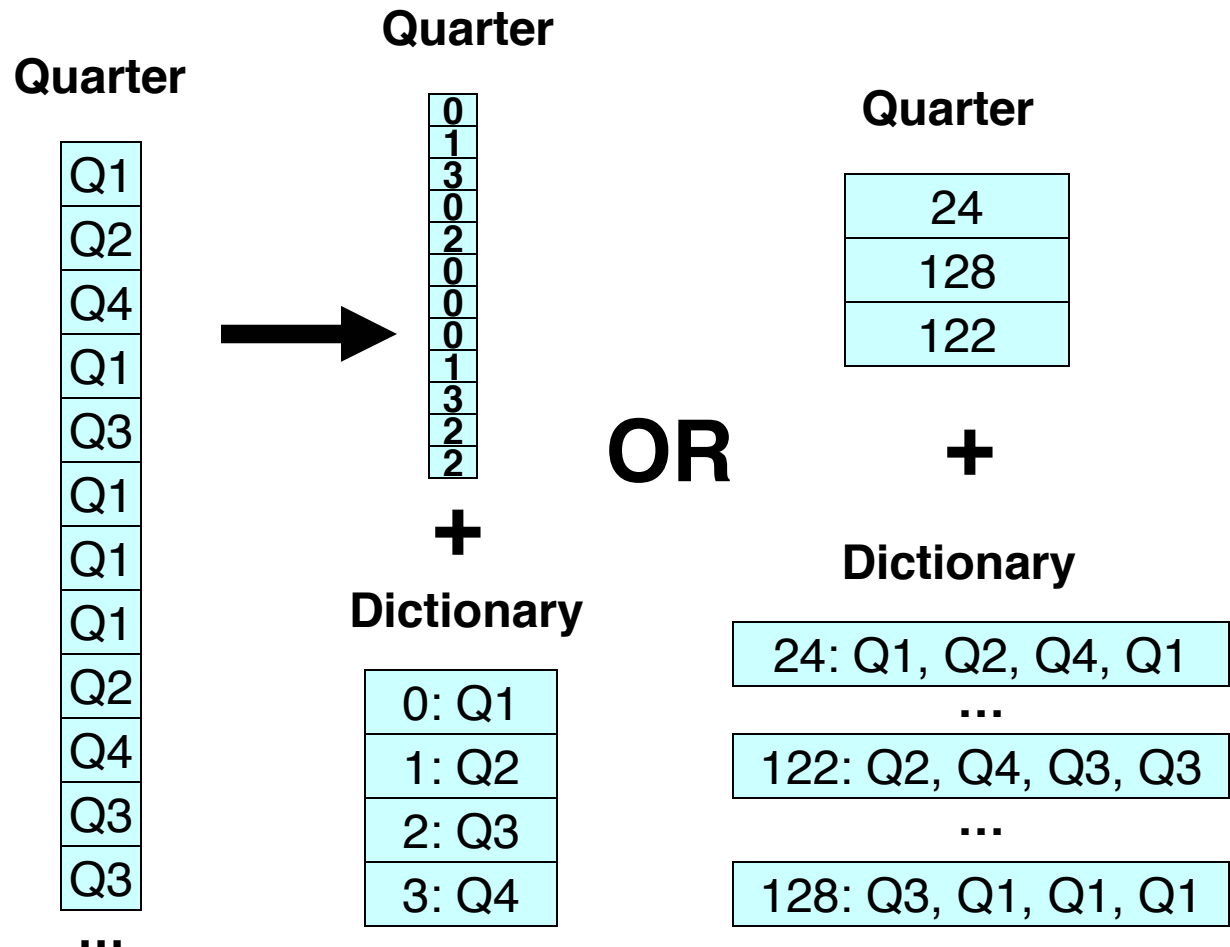
- 1 Each bit-vector can be further compressed if sparse





Dictionary Encoding

- 1 For each unique value create dictionary entry
- 1 Dictionary can be per-block or per-column
- 1 Column-stores have the advantage that dictionary entries may encode multiple values at once

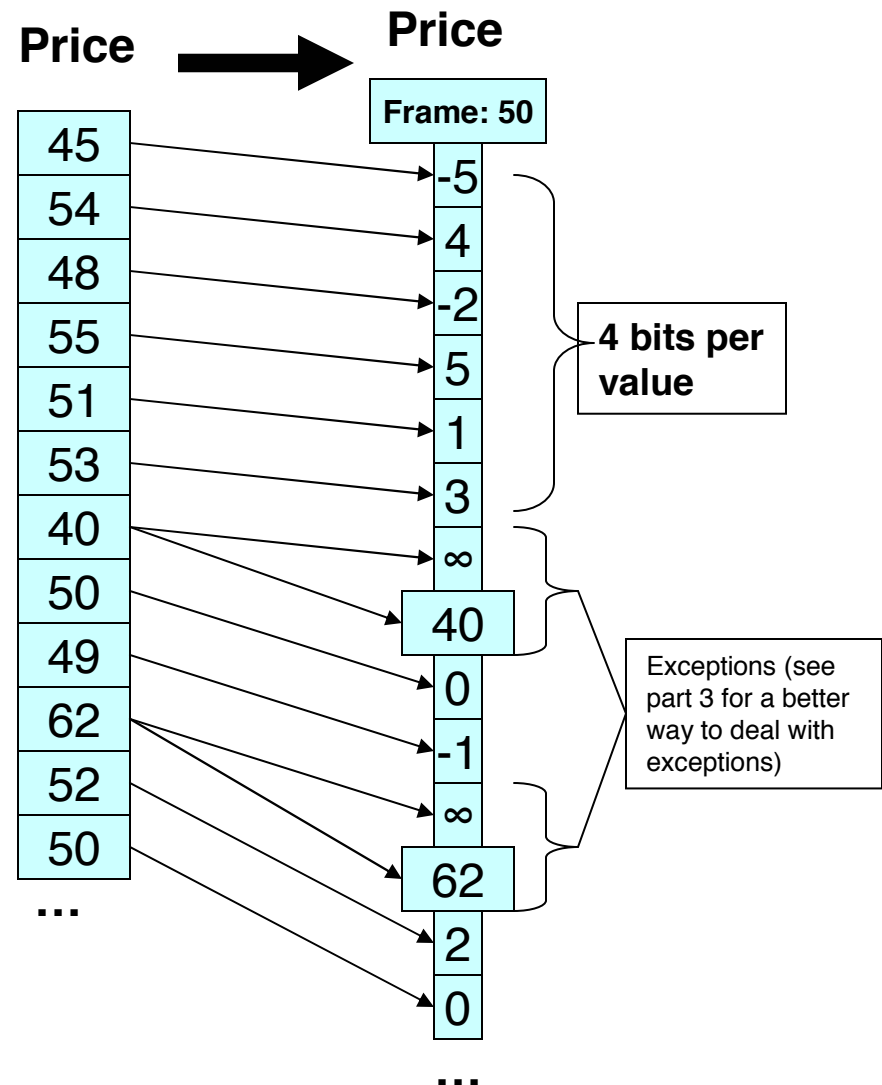




Frame Of Reference Encoding

- 1 Encodes values as b bit offset from chosen frame of reference
- 1 Special escape code (e.g. all bits set to 1) indicates a difference larger than can be stored in b bits
 - 1 After escape code, original (uncompressed) value is written

“Compressing Relations and Indexes” Goldstein, Ramakrishnan, Shaft, ICDE’98

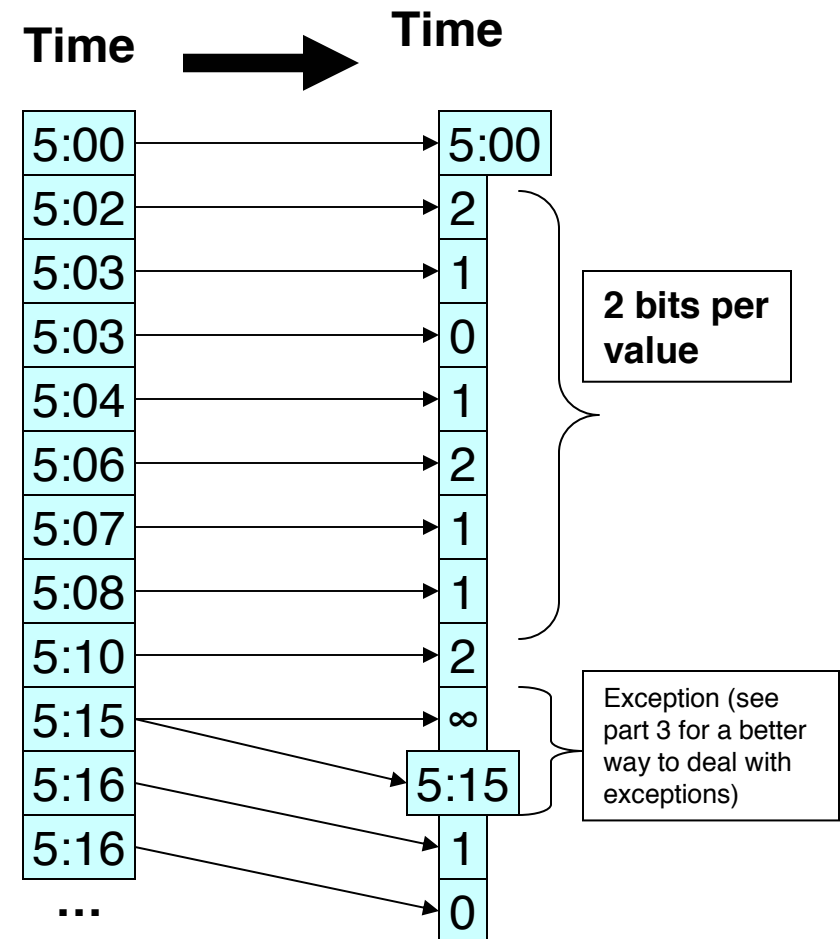




Differential Encoding

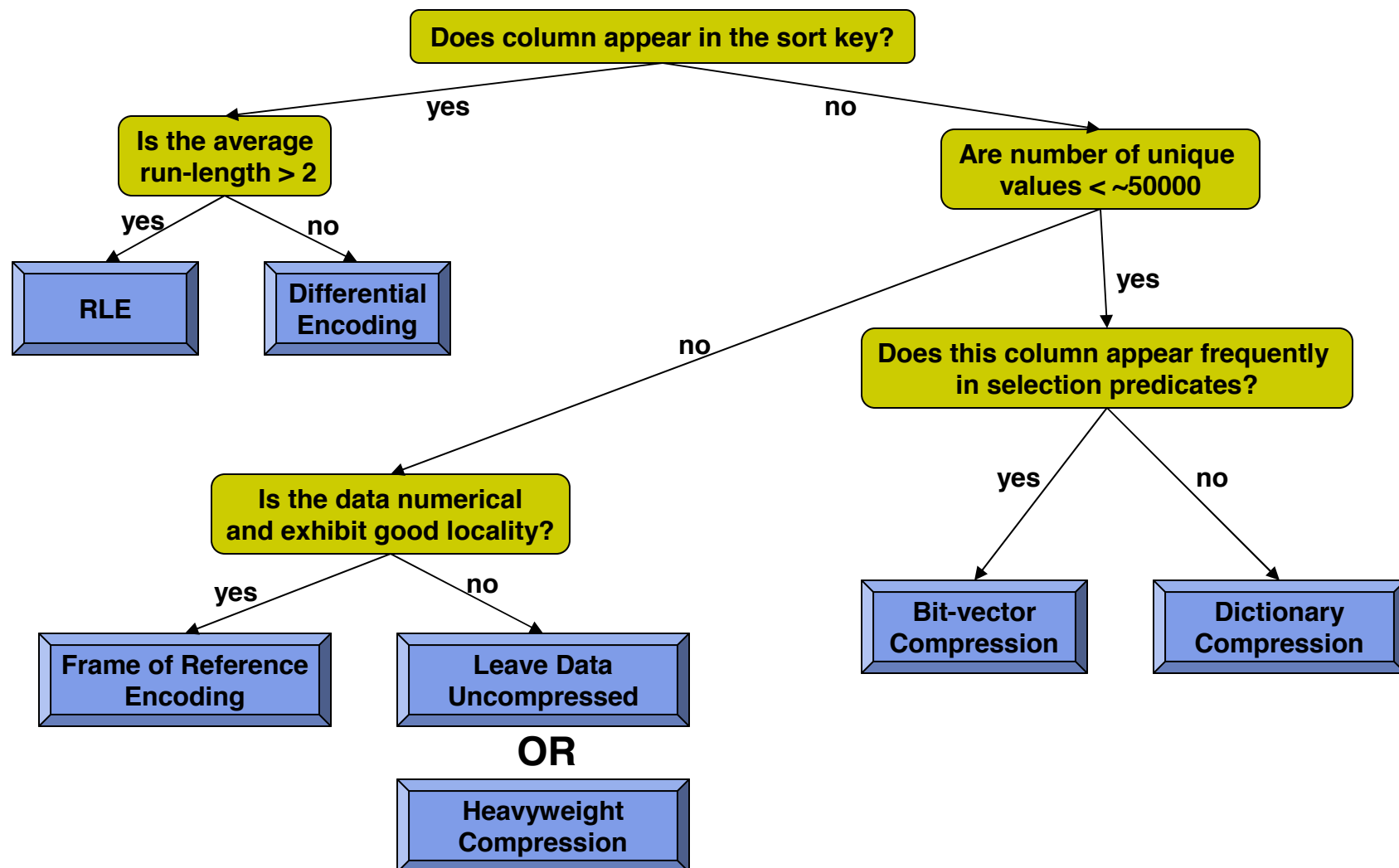
- 1 Encodes values as b bit offset from previous value
- 1 Special escape code (just like frame of reference encoding) indicates a difference larger than can be stored in b bits
 - 1 After escape code, original (uncompressed) value is written
- 1 Performs well on columns containing increasing/decreasing sequences
 - 1 inverted lists
 - 1 timestamps
 - 1 object IDs
 - 1 sorted / clustered columns

“Improved Word-Aligned Binary Compression for Text Indexing”
Ahn, Moffat, TKDE’06





What Compression Scheme To Use?





When should columns be projected?

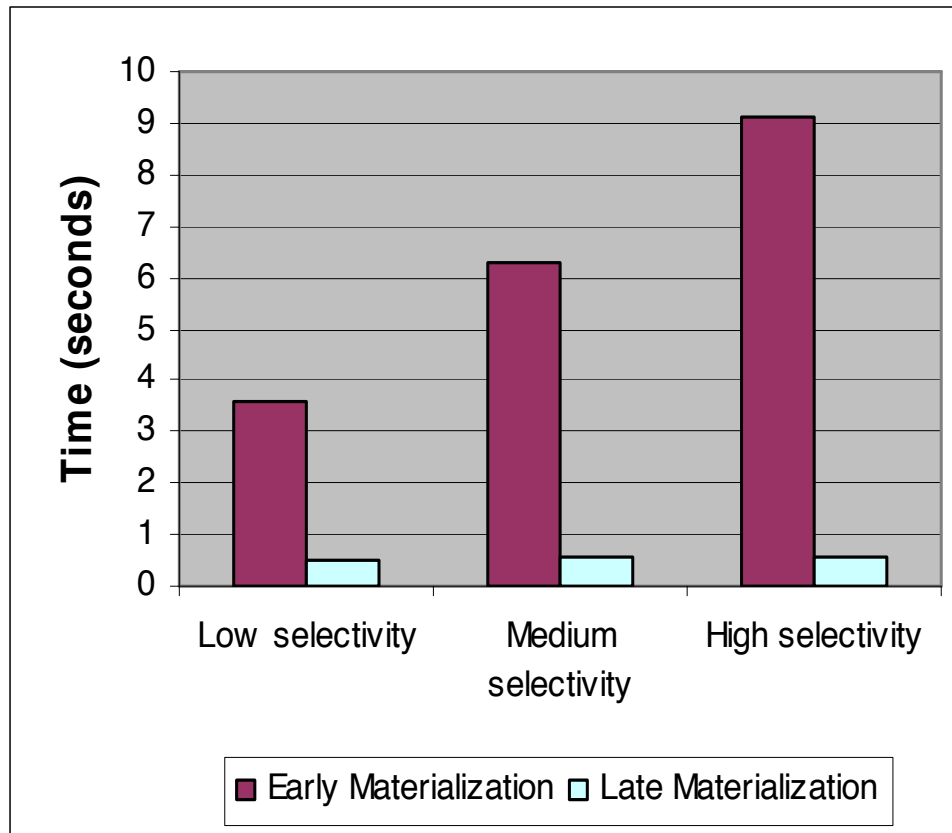
- 1 **Where should column projection operators be placed in a query plan?**
 - 1 **Row-store:**
 - 1 Column projection involves removing unneeded columns from tuples
 - 1 Generally done as early as possible
 - 1 **Column-store:**
 - 1 Operation is almost completely opposite from a row-store
 - 1 Column projection involves reading needed columns from storage and extracting values for a listed set of tuples
 - § This process is called “materialization”
 - 1 **Early materialization: project columns at beginning of query plan**
 - § Straightforward since there is a one-to-one mapping across columns
 - 1 **Late materialization: wait as long as possible for projecting columns**
 - § More complicated since selection and join operators on one column obfuscates mapping to other columns from same table
 - 1 **Most column-stores construct tuples and column projection time**
 - § Many database interfaces expect output in regular tuples (rows)
 - § Rest of discussion will focus on this case



“Materialization Strategies in a Column-Oriented DBMS”
Abadi, Myers, DeWitt, and Madden. ICDE 2007.



For plans without joins, late materialization is a win



QUERY:

```
SELECT C1, SUM(C2)  
FROM table  
WHERE (C1 < CONST) AND  
      (C2 < CONST)  
GROUP BY C1
```

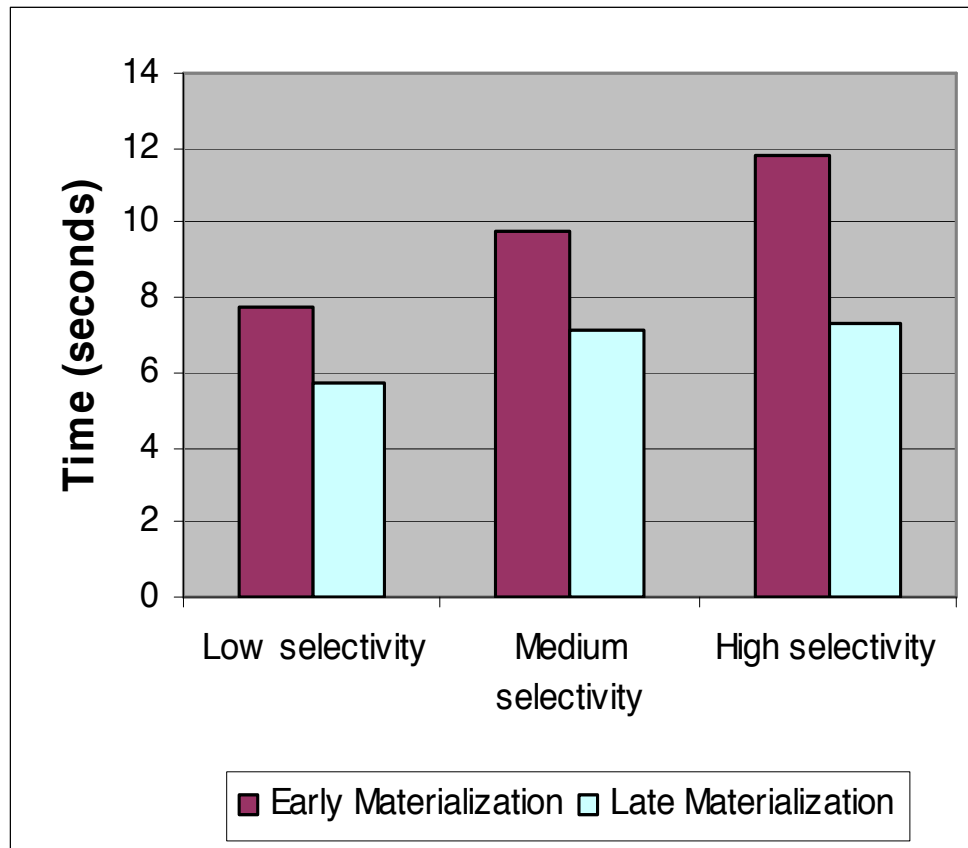
1 Ran on 2 compressed columns from TPC-H scale 10 data



“Materialization Strategies in a Column-Oriented DBMS”
Abadi, Myers, DeWitt, and Madden. ICDE 2007.



Even on uncompressed data, late materialization is still a win



QUERY:

```
SELECT C1, SUM(C2)  
FROM table  
WHERE (C1 < CONST) AND  
      (C2 < CONST)  
GROUP BY C1
```

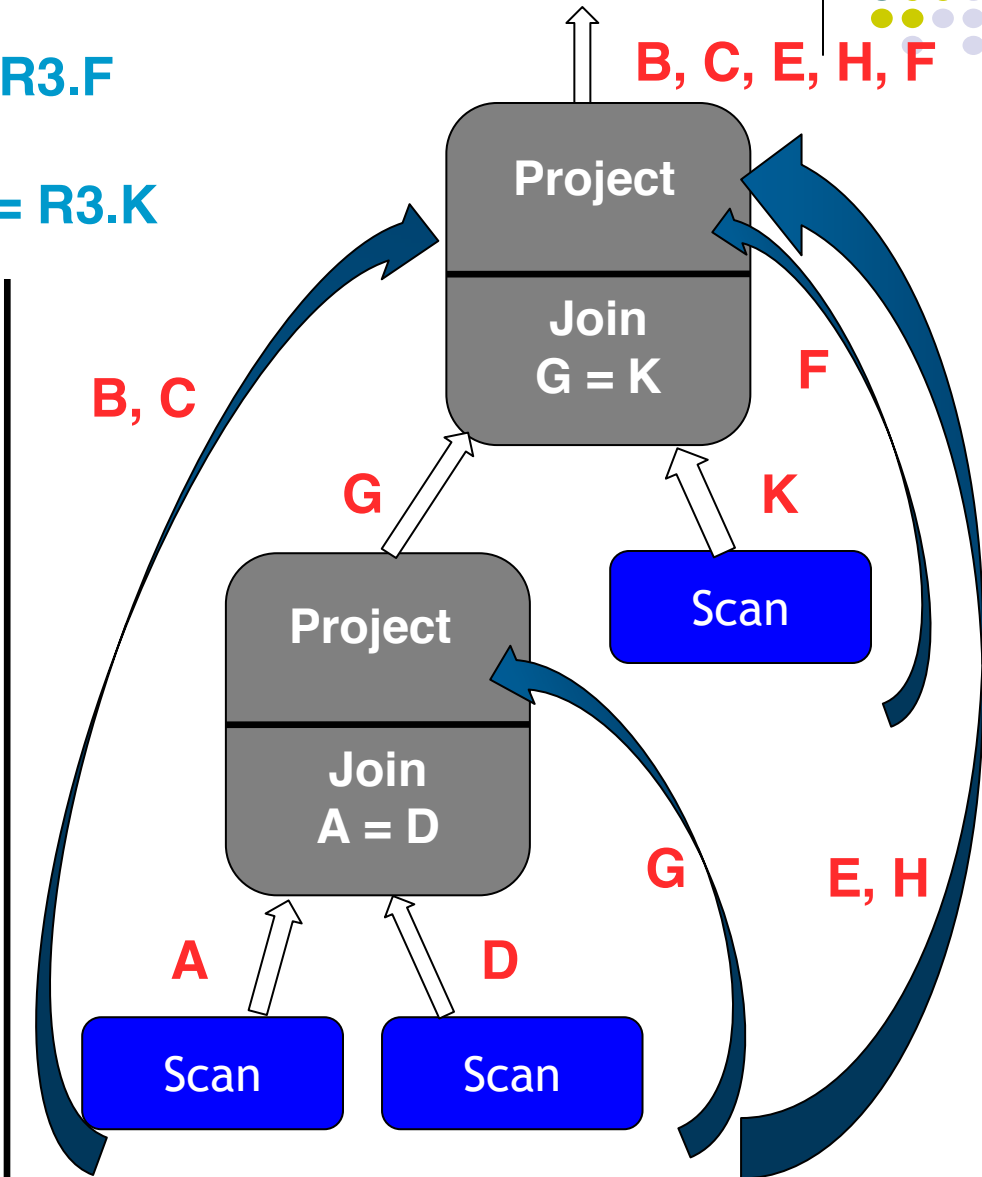
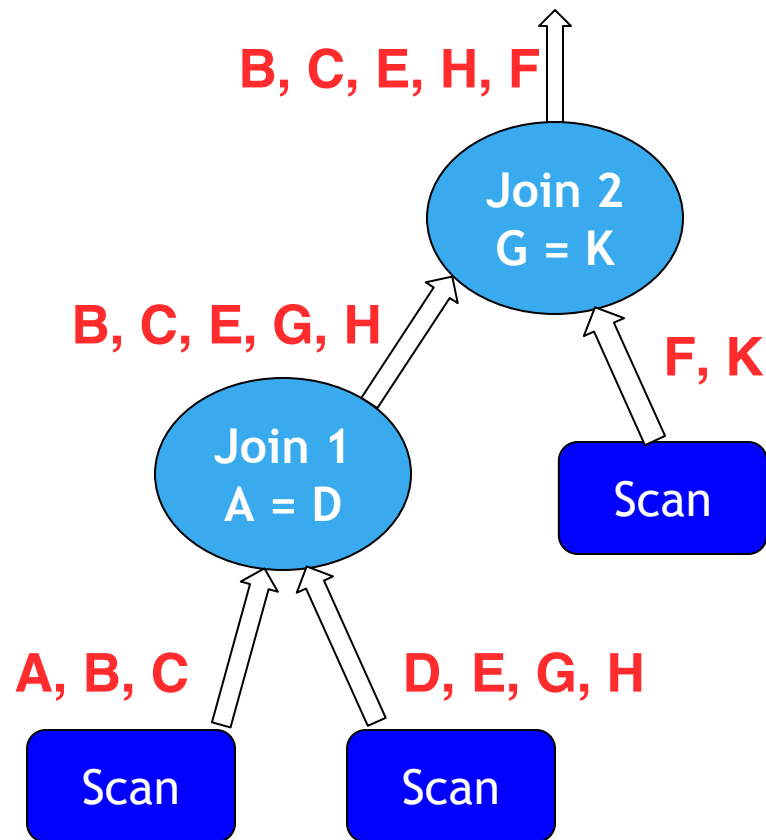
1 **Materializing late still works best**





What about for plans with joins?

Select R1.B, R1.C, R2.E, R2.H, R3.F
From R1, R2, R3
Where R1.A = R2.D AND R2.G = R3.K



What about for plans with joins?

Select R1.B, R1.C, R2.E, R2.H, R3.F
From R1, R2, R3
Where R1.A = R2.D AND R2.G = R3.K

