

Column-Oriented Database Systems

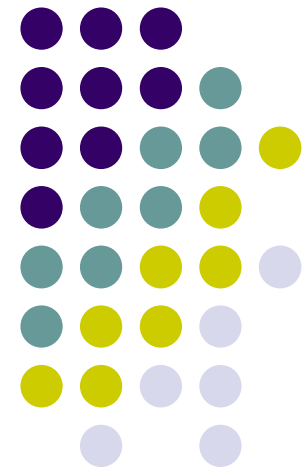
VLDB
2009
Tutorial



Part 1: Stavros Harizopoulos (HP Labs)

~~Part 2: Daniel Abadi (Yale)~~

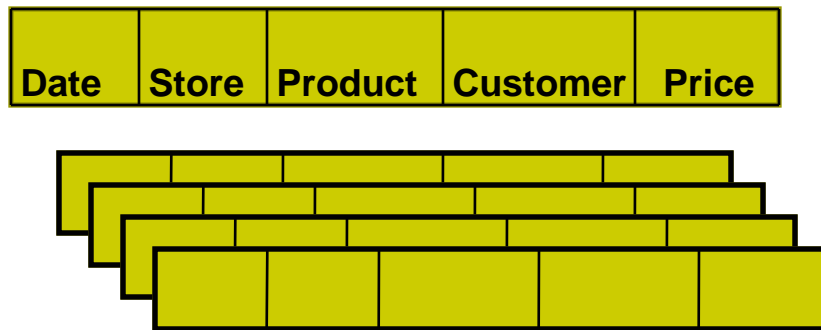
~~Part 3: Peter Boncz (CWI)~~





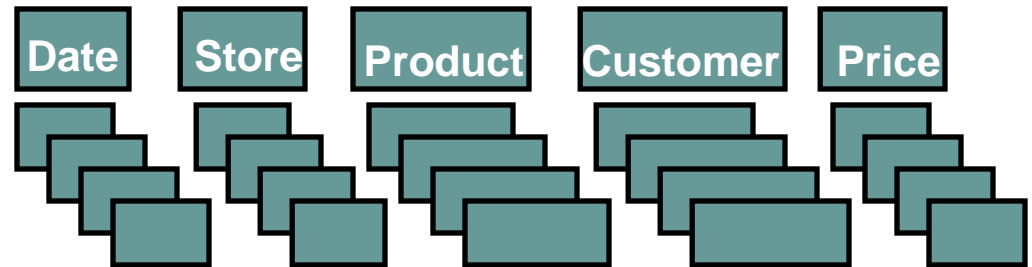
What is a column-store?

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

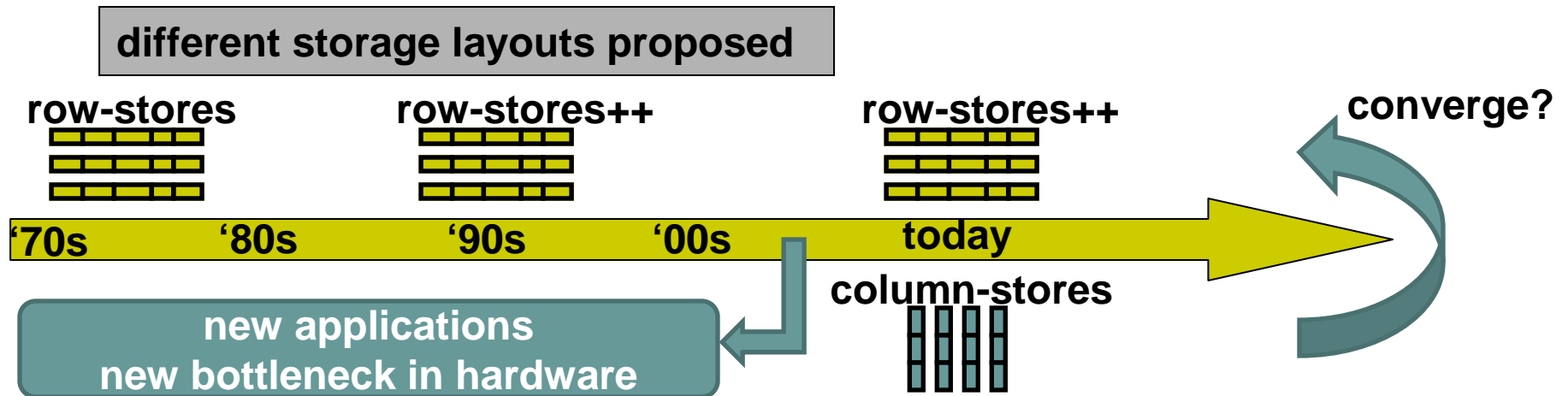
=> *suitable for read-mostly, read-intensive, large data repositories*



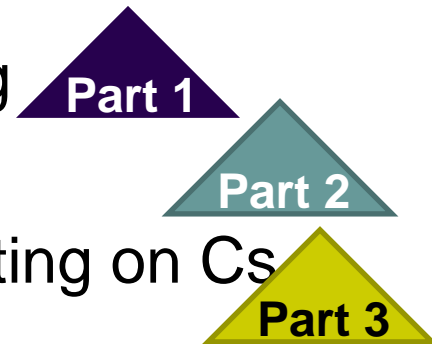


Are these two fundamentally different?

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



- 1 How did we get here, and where we are heading
- 1 What are the column-specific optimizations?
- 1 How do we improve CPU efficiency when operating on Cs





Outline

- 1 Part 1: Basic concepts — *Stavros*
 - 1 Introduction to key features
 - 1 From DSM to column-stores and performance tradeoffs
 - 1 Column-store architecture overview
 - 1 Will rows and columns ever converge?

- 1 Part 2: Column-oriented execution — *Daniel*

- 1 Part 3: MonetDB/X100 and CPU efficiency — *Peter*

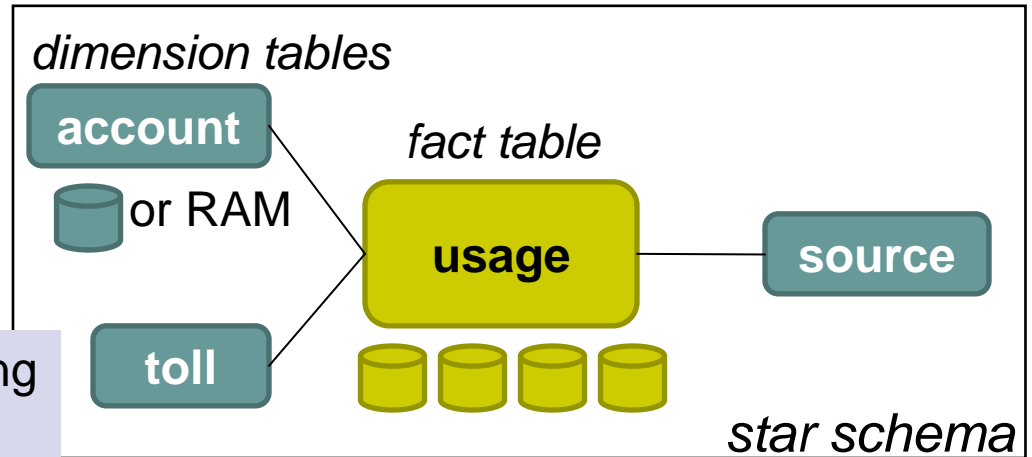




Telco Data Warehousing example

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

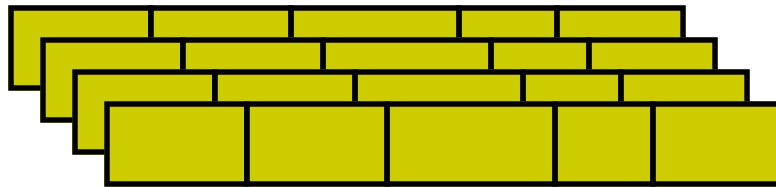
Why? Three main factors (next slides)





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

row store



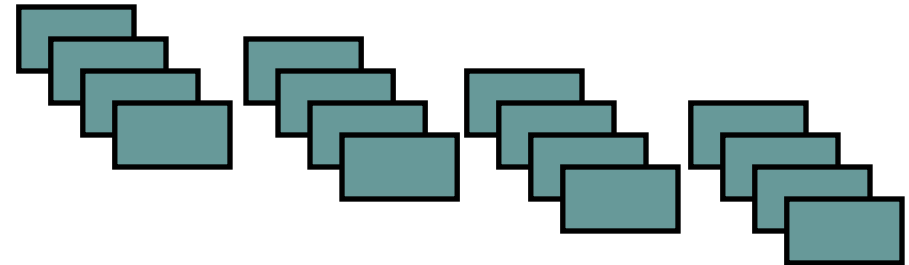
read pages containing entire rows

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

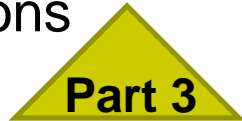
**What about heavily-indexed row-stores?
(works well for single column access,
cross-column joins become increasingly expensive)**





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)

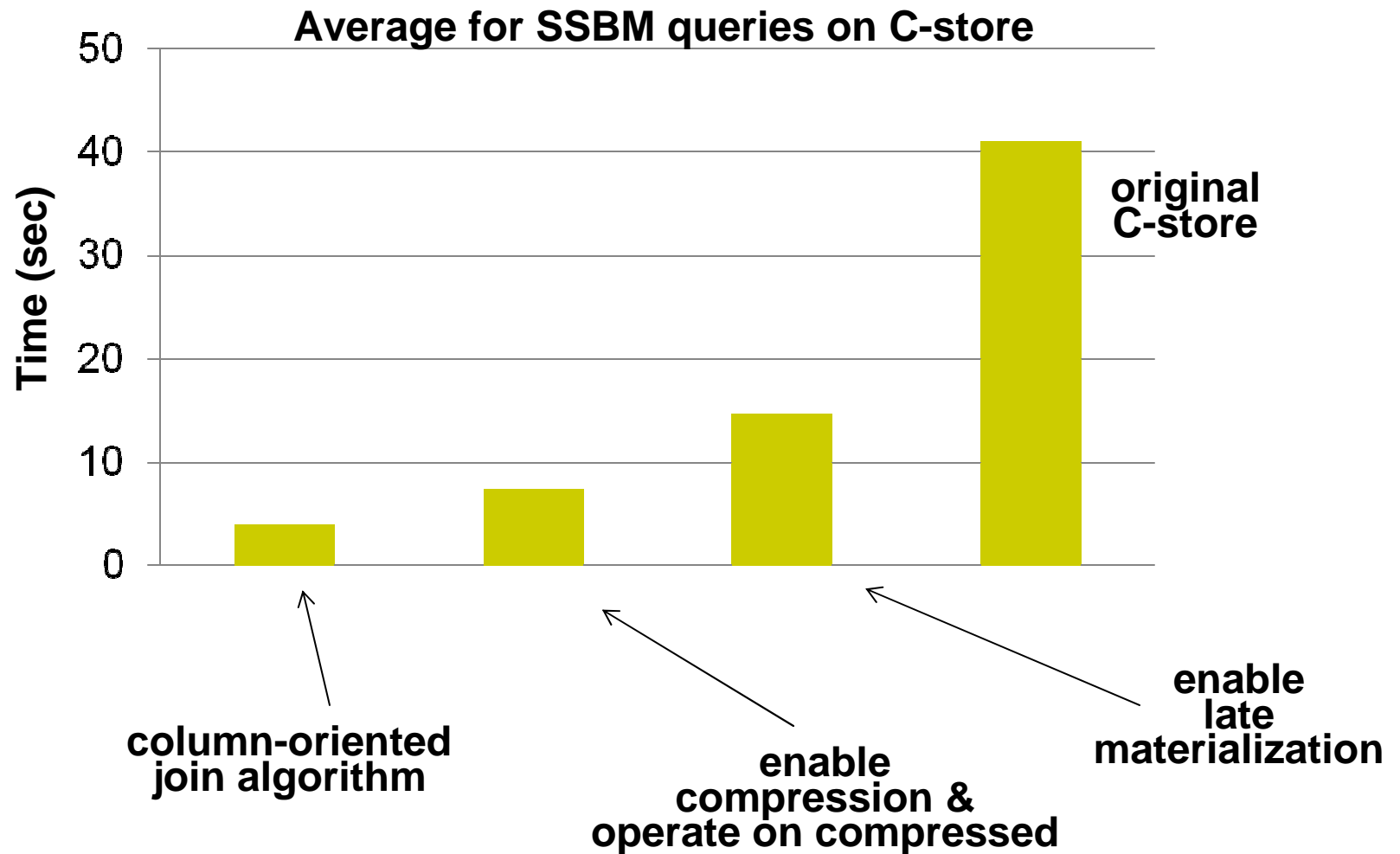


See: *Database Cracking*, from CWI



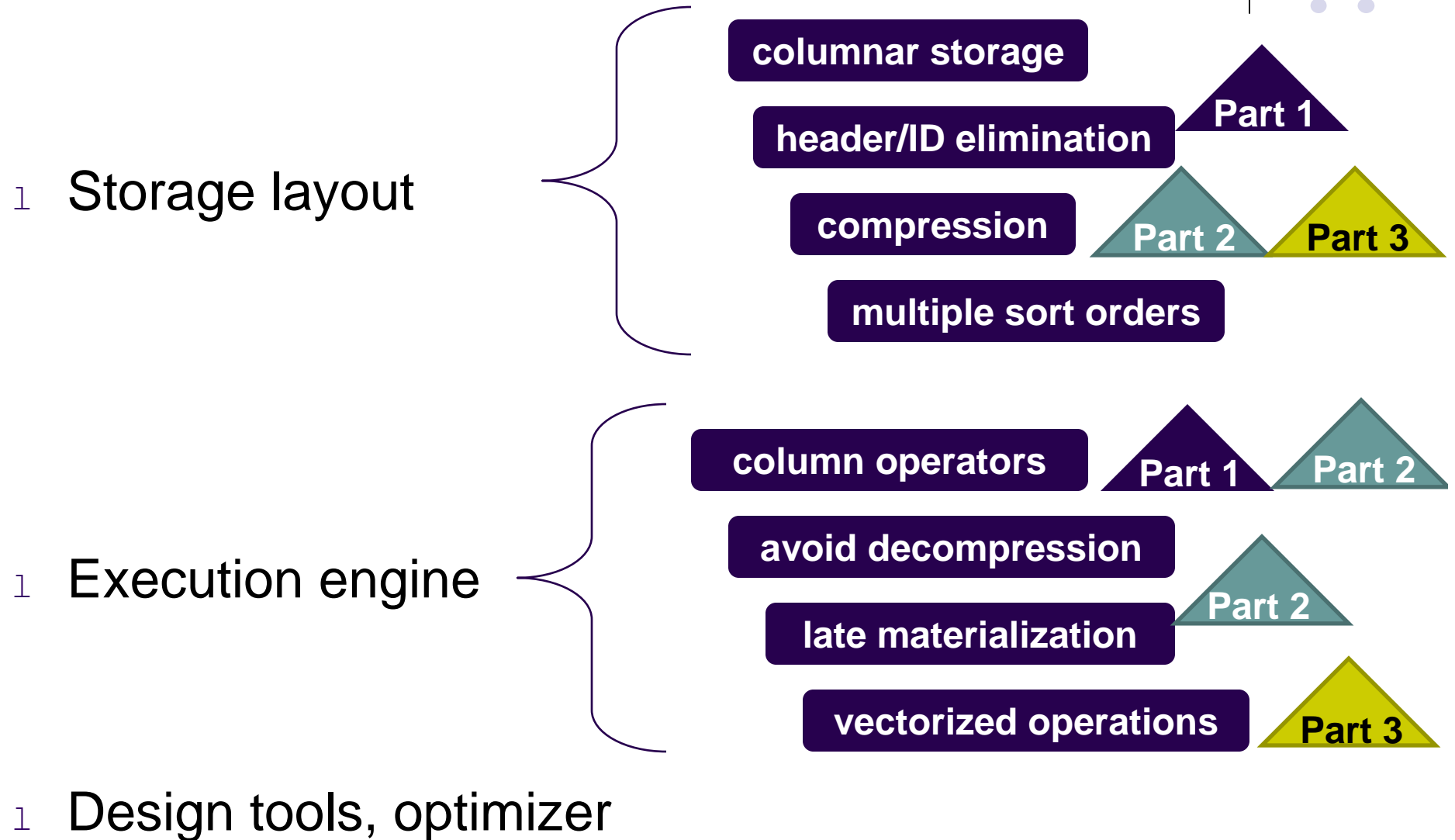
“Column-Stores vs Row-Stores: How Different are They Really?” Abadi, Hachem, and Madden. SIGMOD 2008.

Effect on C-Store performance





Summary of column-store key features





Outline

- 1 Part 1: Basic concepts — *Stavros*
 - 1 Introduction to key features
 - 1 From DSM to column-stores and performance tradeoffs
 - 1 Column-store architecture overview
 - 1 Will rows and columns ever converge?
- 1 Part 2: Column-oriented execution — *Daniel*
- 1 Part 3: MonetDB/X100 and CPU efficiency — *Peter*





From DSM to Column-stores

70s -1985:

TOD: Time Oriented Database – Wiederhold et al.
"A Modular, Self-Describing Clinical Databank System," *Computers and Biomedical Research*, 1975
More 1970s: Transposed files, Lorie, Batory,

"An overview of cantor: a new system for data analysis"
Karasalo, Svensson, SSDBM 1983

1985: DSM paper

"A decomposition storage model"
Copeland and Khoshafian. SIGMOD 1985.

1990s: Commercialization through SybaseIQ

Late 90s – 2000s: Focus on main-memory performance

1 DSM "on steroids" [1997 – now] CWI: MonetDB

1 Hybrid DSM/NSM [2001 – 2004] Wisconsin: PAX, Fractured Mirrors

Michigan: Data Morphing

CMU: Clotho

2005 – : Re-birth of read-optimized DSM as "column-store"

MIT: C-Store

CWI: MonetDB/X100

10+ startups

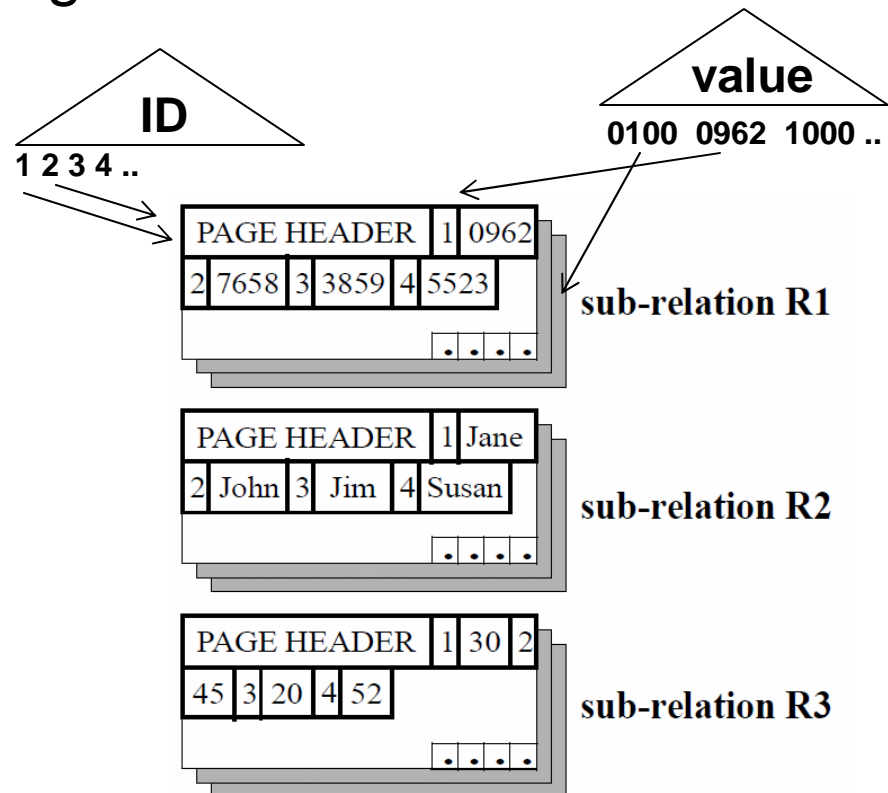
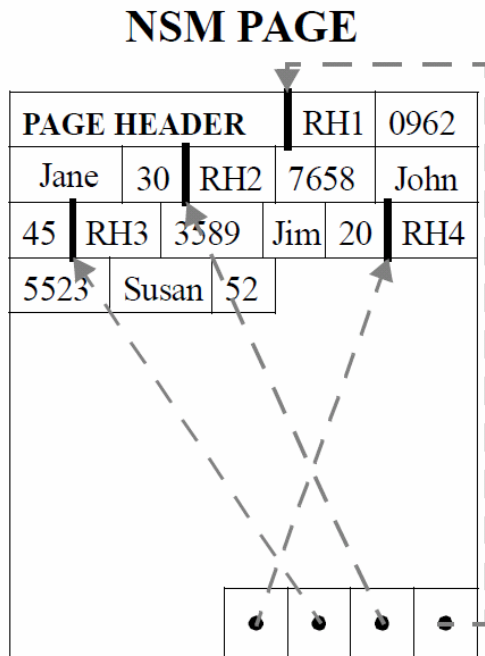




The original DSM paper

“A decomposition storage model” Copeland and Khoshafian. SIGMOD 1985.

- 1 Proposed as an alternative to NSM
- 1 2 indexes: clustered on ID, non-clustered on value
- 1 Speeds up queries projecting few columns
- 1 Requires more storage





Memory wall and PAX

1 90s: Cache-conscious research

from: “Cache Conscious Algorithms for Relational Query Processing.”
Shatdal, Kant, Naughton. VLDB 1994.

to: “Database Architecture Optimized for the New Bottleneck: Memory Access.”
Boncz, Manegold, Kersten. VLDB 1999.

and: “DBMSs on a modern processor: Where does time go?” Ailamaki, DeWitt, Hill, Wood. VLDB 1999.

1 PAX: Partition Attributes Across

1 Retains NSM I/O pattern

1 Optimizes cache-to-RAM communication

“Weaving Relations for Cache Performance.”
Ailamaki, DeWitt, Hill, Skounakis, VLDB 2001.

PAX PAGE

PAGE HEADER				0962	7658
3859	5523				
<hr/>					
Jane	John	Jim	Susan		
<hr/>					
				•	•
30	52	45	20		
<hr/>					





More hybrid NSM/DSM schemes

1 Dynamic PAX: Data Morphing

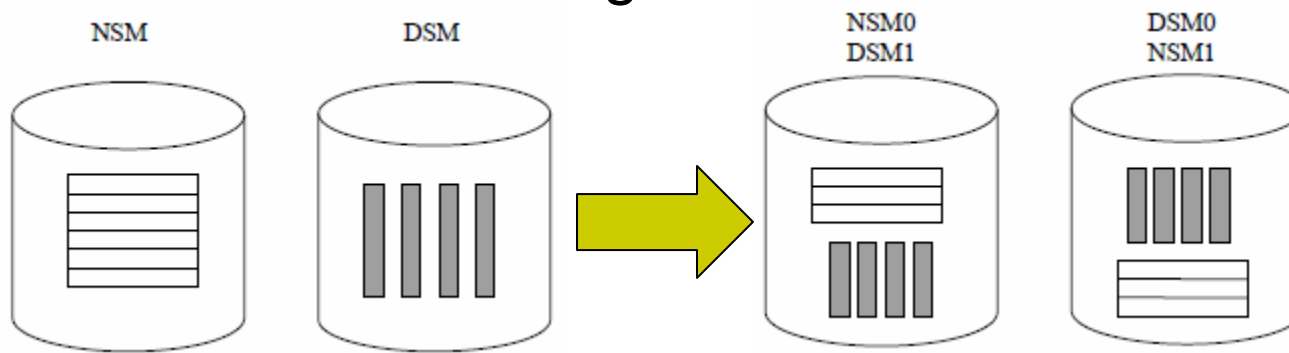
“Data morphing: an adaptive, cache-conscious storage technique.” Hankins, Patel, VLDB 2003.

1 Clotho: custom layout using scatter-gather I/O

“Clotho: Decoupling Memory Page Layout from Storage Organization.” Shao, Schindler, Schlosser, Ailamaki, and Ganger. VLDB 2004.

1 Fractured mirrors

1 Smart mirroring with both NSM/DSM copies



“A Case For Fractured Mirrors.” Ramamurthy, DeWitt, Su, VLDB 2002.





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





Performance tradeoffs: columns vs. rows

DSM traditionally was not favored by technology trends
How has this changed?

1 Optimized DSM in “Fractured Mirrors,” 2002

1 “Apples-to-apples” comparison

“Performance Tradeoffs in Read-Optimized Databases”
Harizopoulos, Liang, Abadi, Madden, VLDB’06

1 Follow-up study

“Read-Optimized Databases, In-Depth” Holloway, DeWitt, VLDB’08

1 Main-memory DSM vs. NSM

“DSM vs. NSM: CPU performance tradeoffs in block-oriented query processing” Boncz, Zukowski, Nes, DaMoN’08

1 Flash-disks: a come-back for PAX?

“Fast Scans and Joins Using Flash Drives” Shah, Harizopoulos, Wiener, Graefe. DaMoN’08

“Query Processing Techniques for Solid State Drives”
Tsirogiannis, Harizopoulos, Shah, Wiener, Graefe, SIGMOD’08



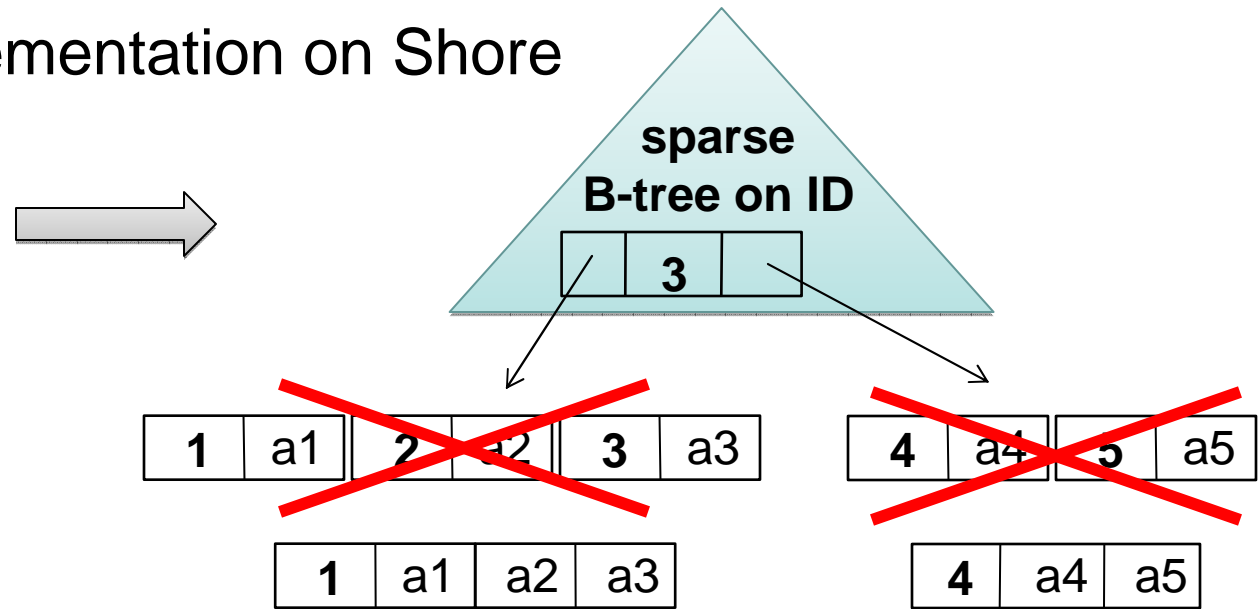


Fractured mirrors: a closer look

- 1 Store DSM relations inside a B-tree
 - 1 Leaf nodes contain values
 - 1 Eliminate IDs, amortize header overhead
 - 1 Custom implementation on Shore

“A Case For Fractured Mirrors” Ramamurthy, DeWitt, Su, VLDB 2002.

Tuple Header	TID	Column Data
	1	a1
	2	a2
	3	a3
	4	a4
	5	a5



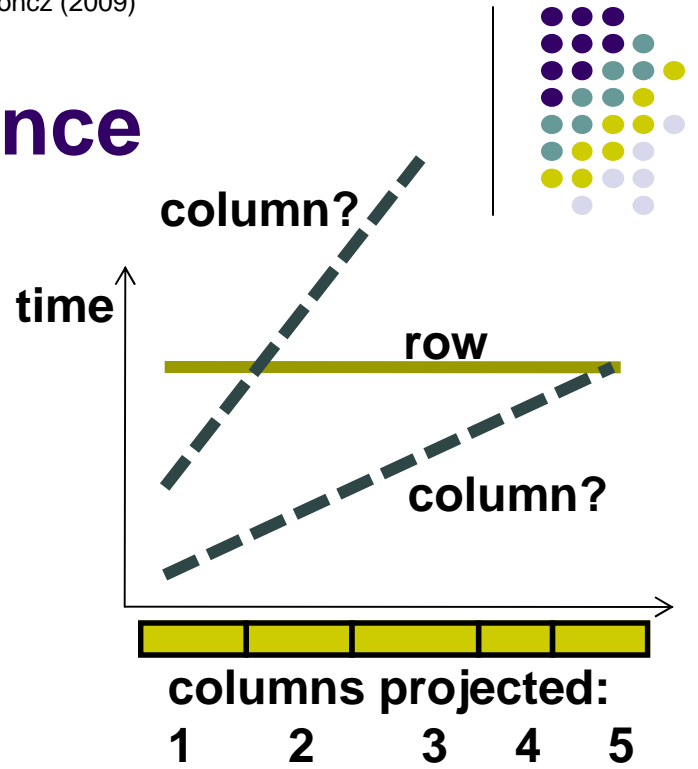
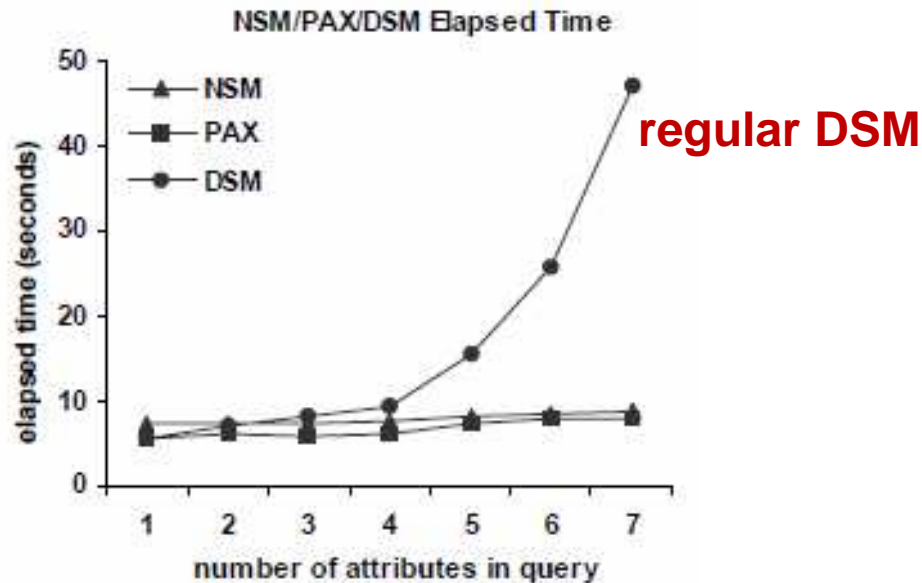
Similar: storage density comparable to column stores

“Efficient columnar storage in B-trees” Graefe. Sigmod Record 03/2007.

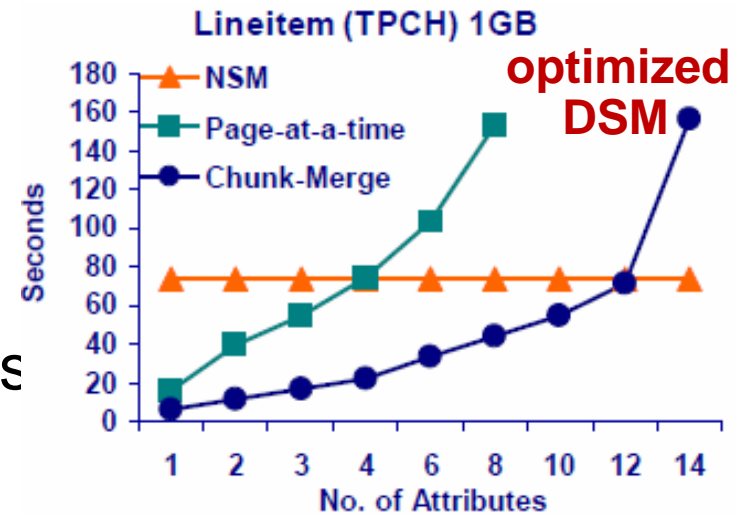


Fractured mirrors: performance

From PAX paper:



- 1 Chunk-based tuple merging
 - 1 Read in segments of M pages
 - 1 Merge segments in memory
 - 1 Becomes CPU-bound after 5 pages

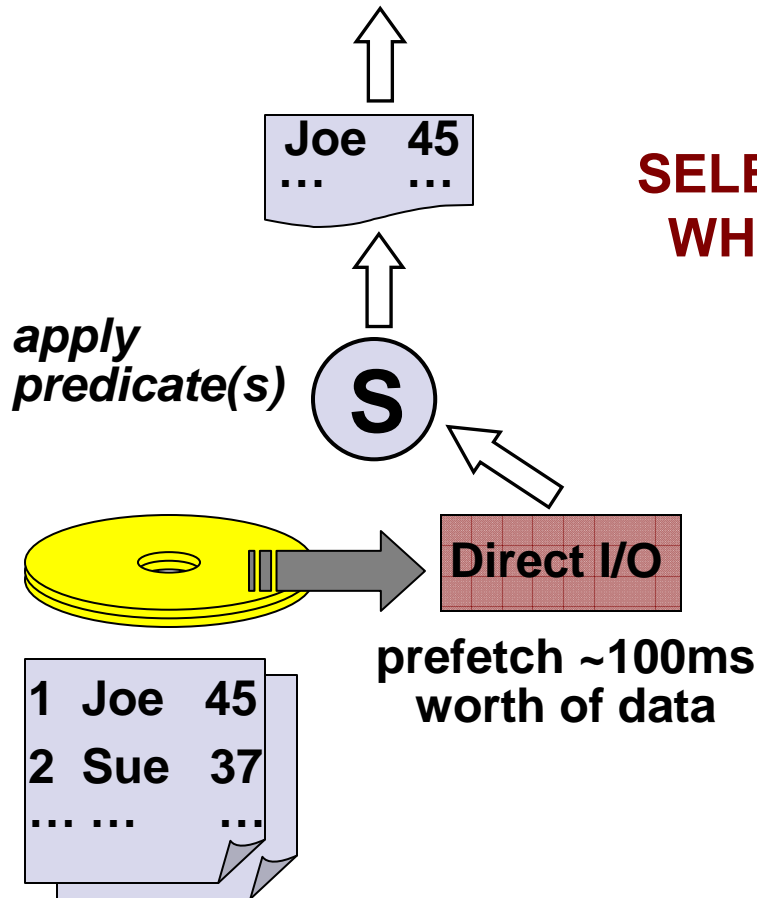


Column-scanner implementation

“Performance Tradeoffs in Read-Optimized Databases”
Harizopoulos, Liang, Abadi,
Madden, VLDB’06

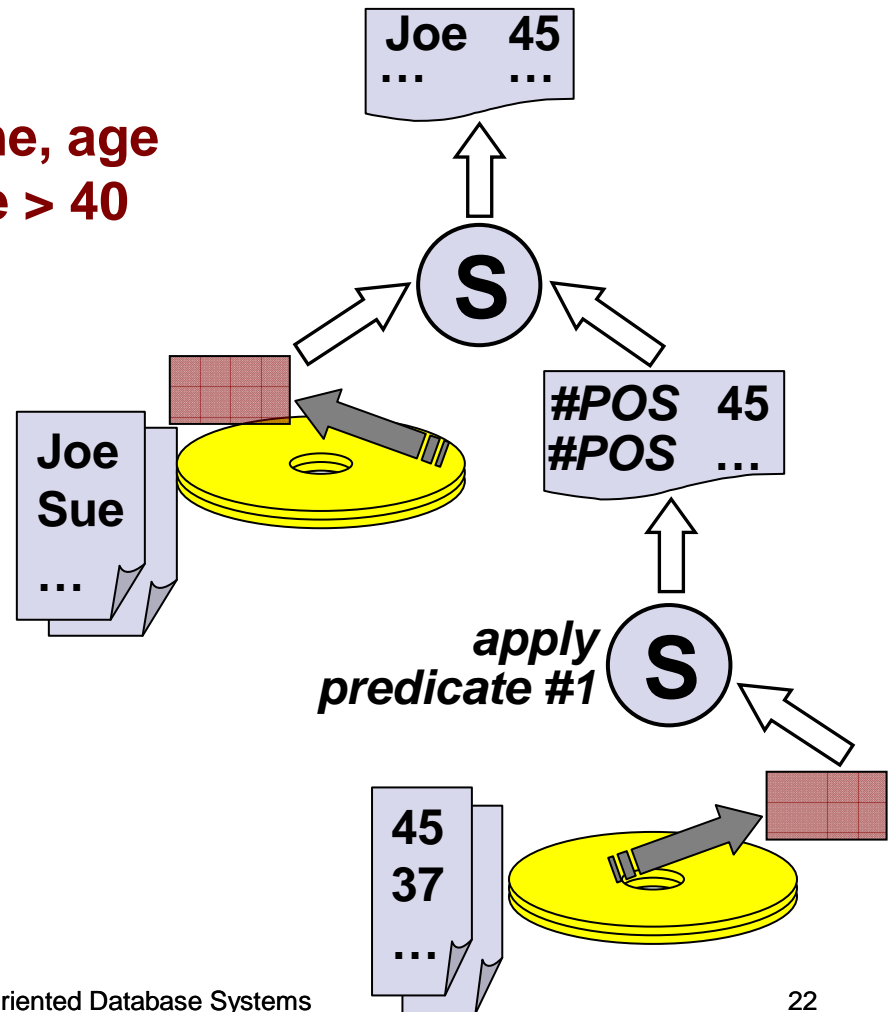


row scanner



**SELECT name, age
WHERE age > 40**

column scanner

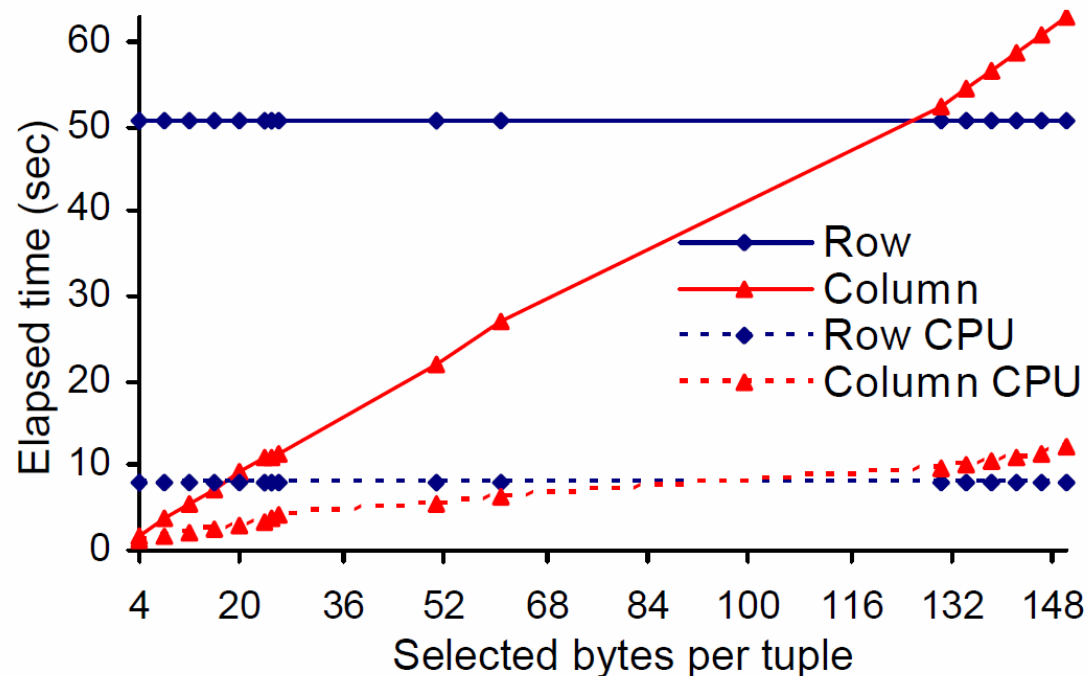




Scan performance

- 1 Large prefetch hides disk seeks in columns
- 1 Column-CPU efficiency with lower selectivity
- 1 Row-CPU suffers from memory stalls
- 1 Memory stalls disappear in narrow tuples
- 1 Compression: similar to narrow

not shown,
details in the paper

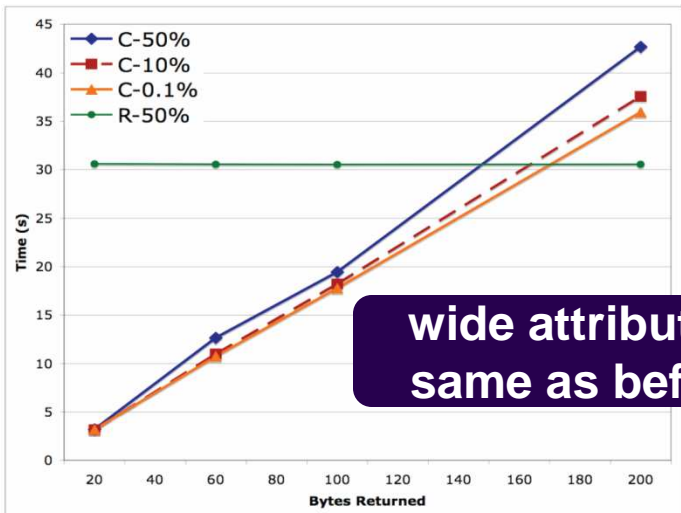




Even more results

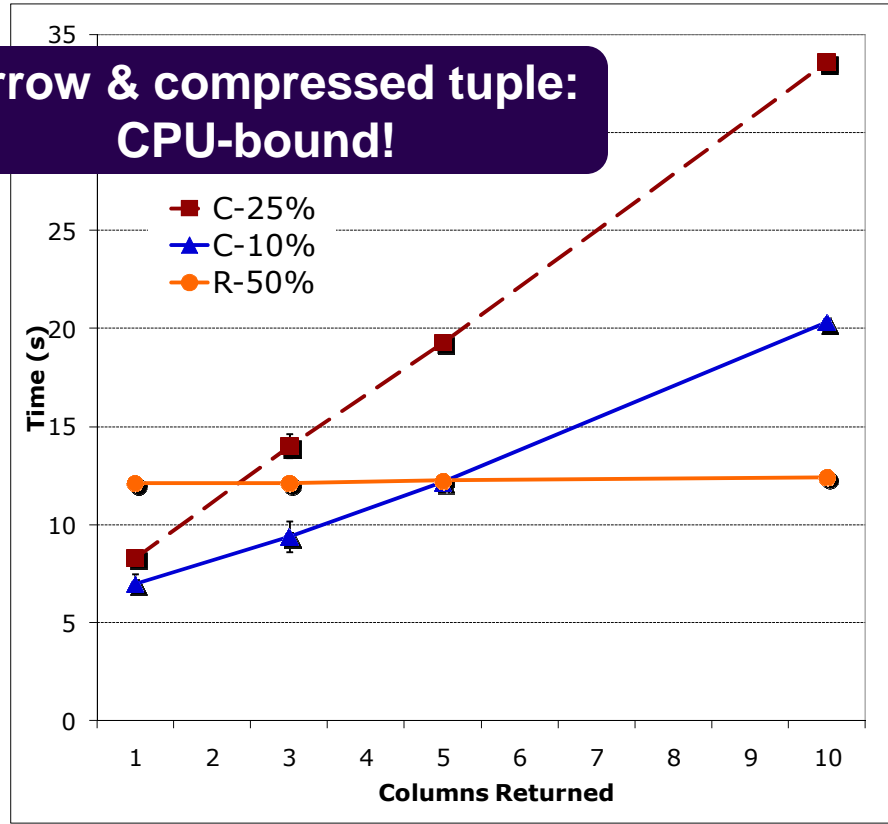
“Read-Optimized Databases, In-Depth” Holloway, DeWitt, VLDB’08

- Same engine as before
- Additional findings



**wide attributes:
same as before**

**narrow & compressed tuple:
CPU-bound!**



Non-selective queries, narrow tuples, favor well-compressed rows

Materialized views are a win

Scan times determine early materialized joins

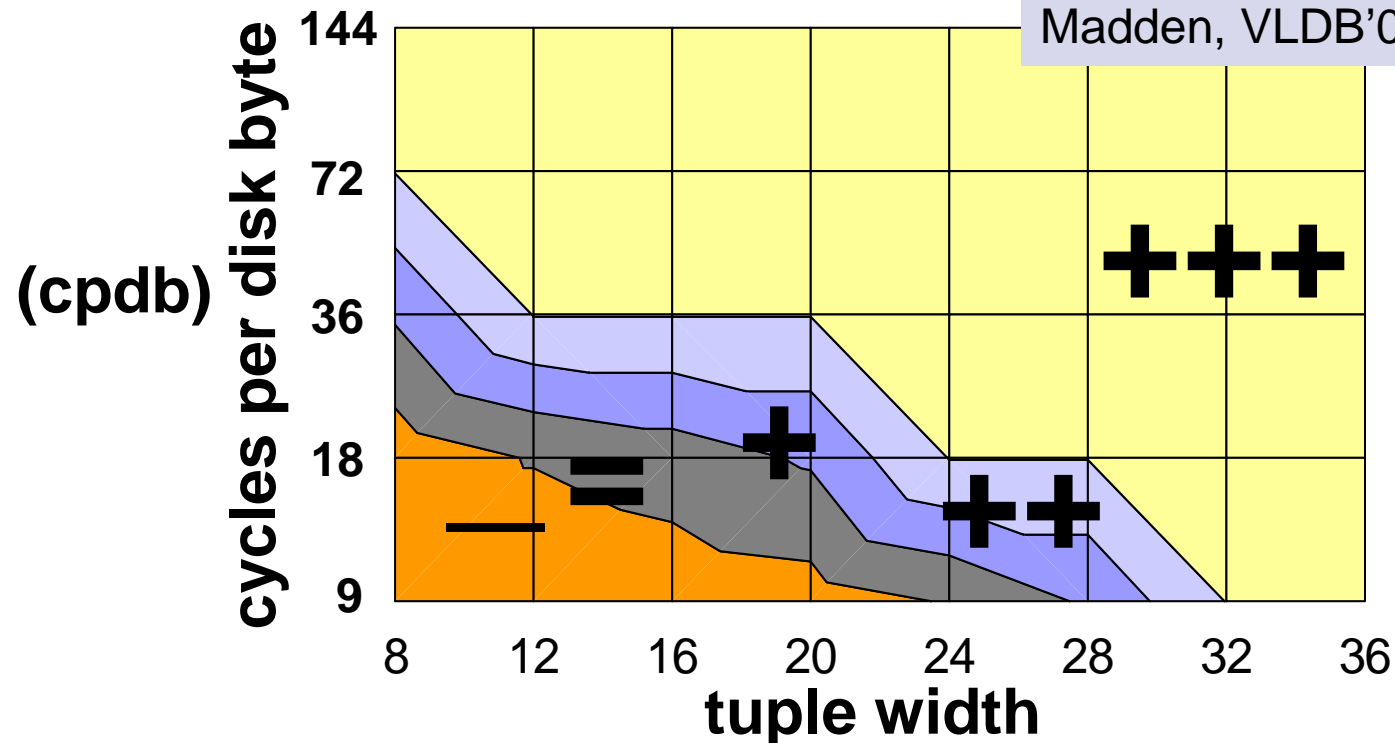
Column-joins are covered in part 2!





Speedup of columns over rows

“Performance Tradeoffs in Read-Optimized Databases”
Harizopoulos, Liang, Abadi, Madden, VLDB’06

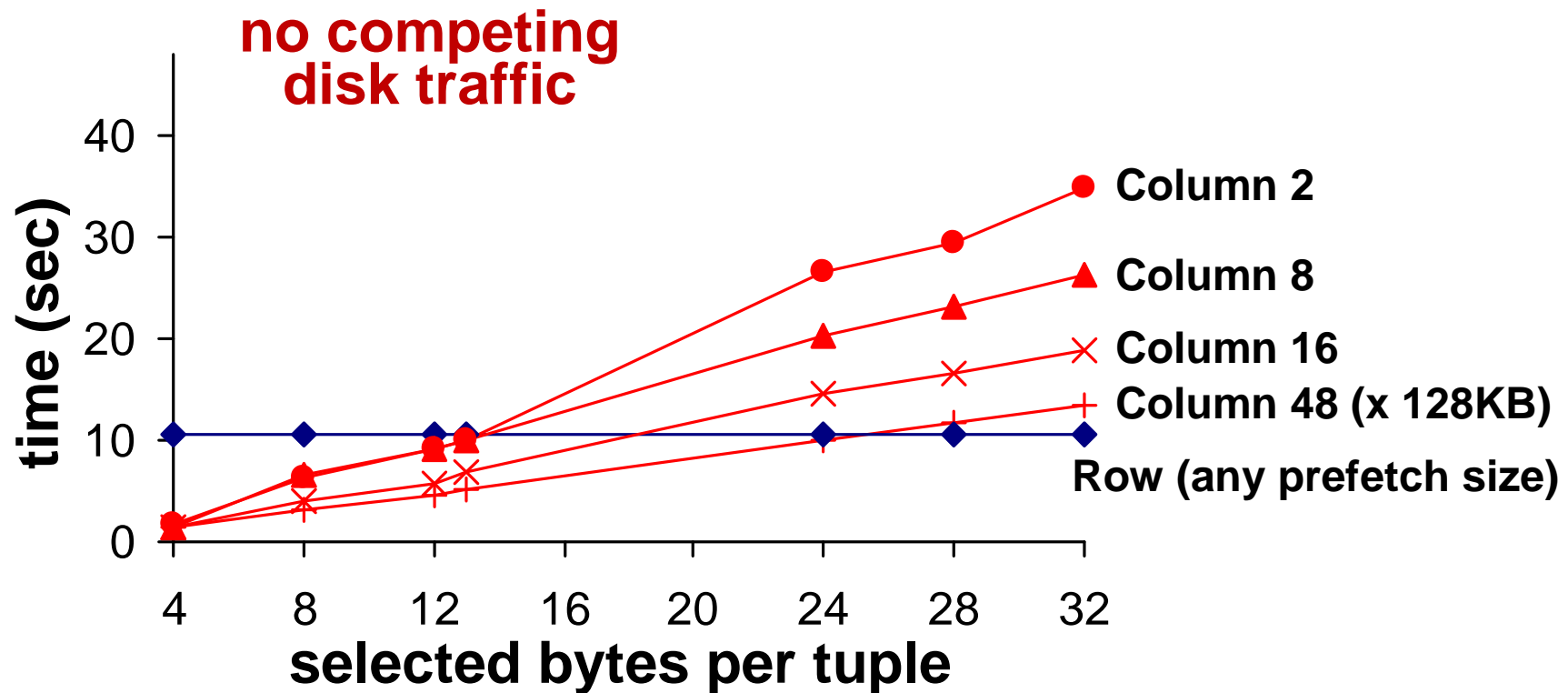


- 1 Rows favored by narrow tuples and low *cpdb*
- 1 Disk-bound workloads have higher *cpdb*





Varying prefetch size

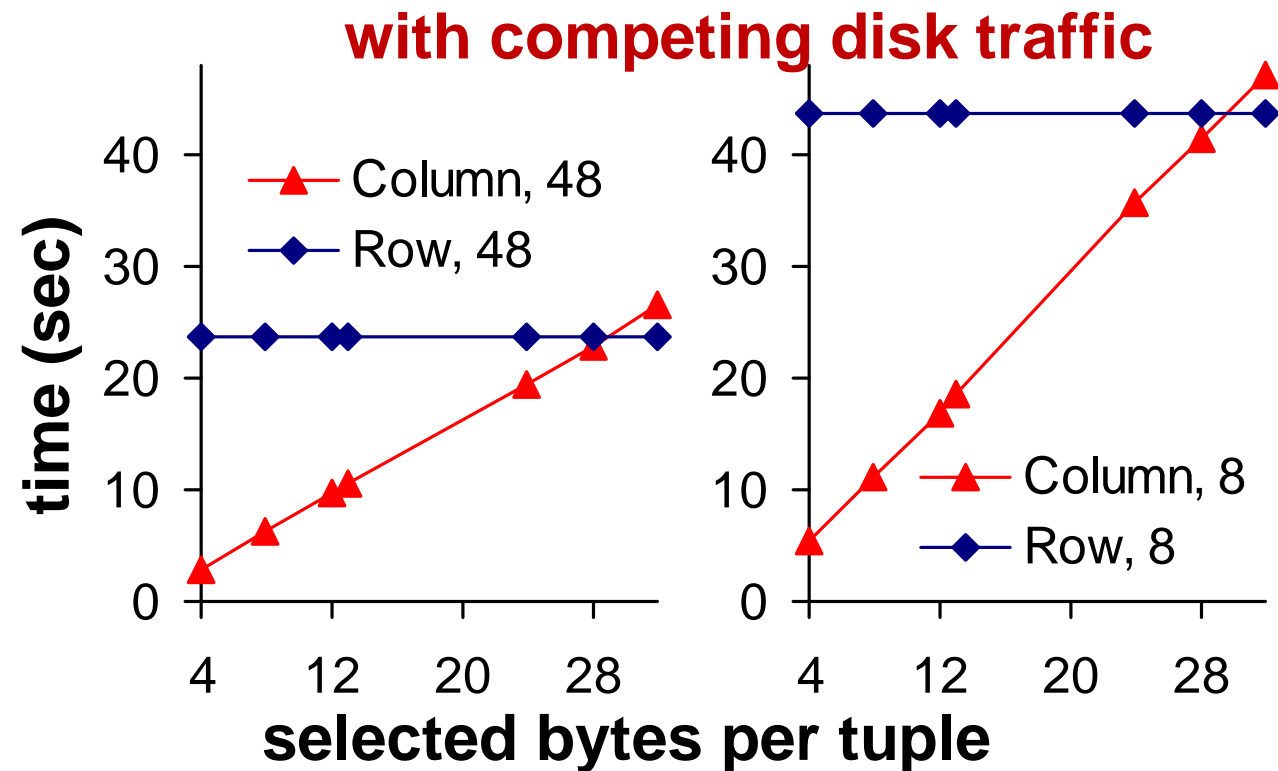


- 1 No prefetching hurts columns in single scans





Varying prefetch size



- 1 No prefetching hurts columns in single scans
- 1 Under competing traffic, columns outperform rows for any prefetch size





CPU Performance

“DSM vs. NSM: CPU performance trade offs in block-oriented query processing”
Boncz, Zukowski, Nes, DaMoN’08

- 1 Benefit in on-the-fly conversion between NSM and DSM
- 1 DSM: sequential access (block fits in L2), random in L1
- 1 NSM: random access, SIMD for grouped Aggregation

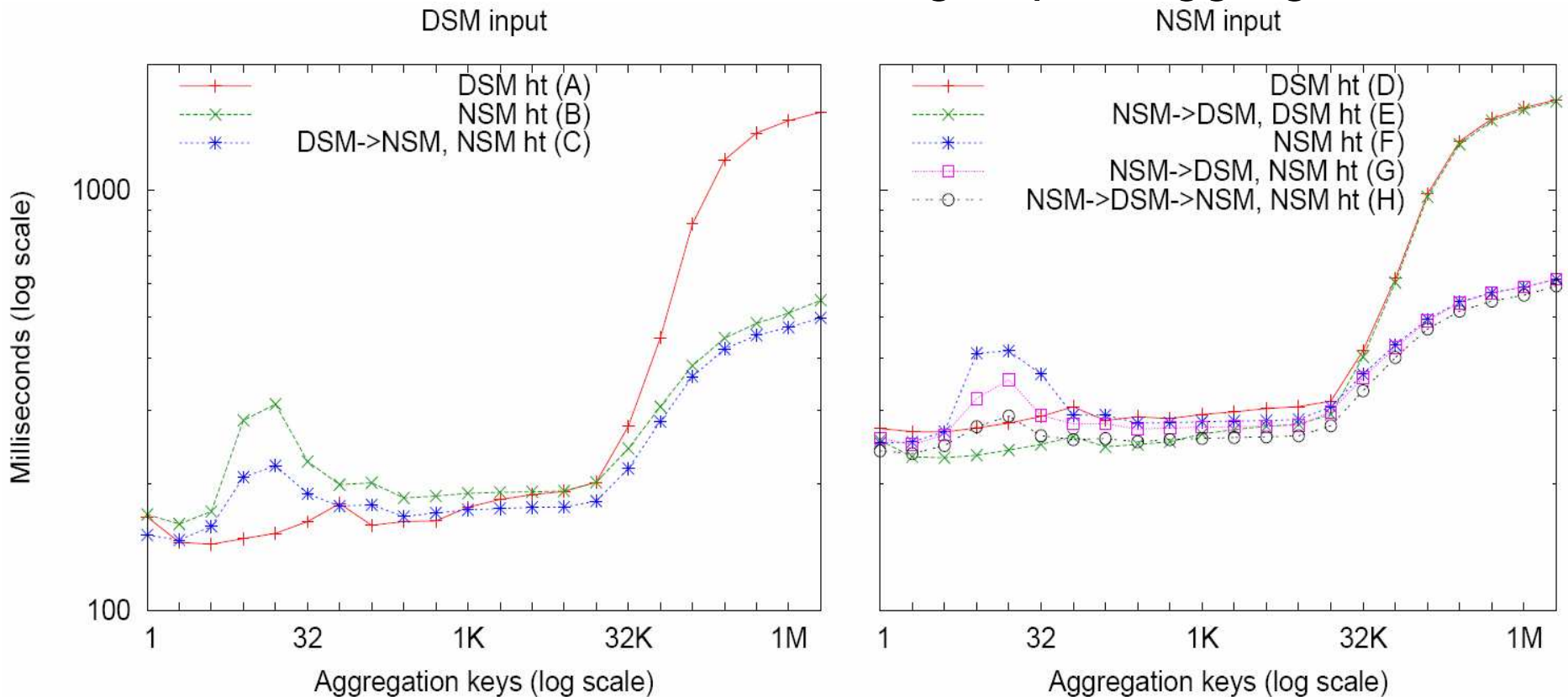


Figure 5: TPC-H Q1, with a varying number of keys and different data organizations (ht – hash table)



New storage technology: Flash SSDs

- 1 Performance characteristics
 - 1 very fast random reads, slow random writes
 - 1 fast sequential reads and writes
 - 1 Price per bit (capacity follows)
 - 1 cheaper than RAM, order of magnitude more expensive than Disk
 - 1 Flash Translation Layer introduces unpredictability
 - 1 avoid random writes!
 - 1 Form factors not ideal yet
 - 1 SSD (⌘ small reads still suffer from SATA overhead/OS limitations)
 - 1 PCI card (⌘ high price, limited expandability)
-
- 1 Boost Sequential I/O in a simple package
 - 1 Flash RAID: very tight bandwidth/cm³ packing (4GB/sec inside the box)
 - 1 Column Store Updates
 - 1 useful for delta structures and logs
 - 1 Random I/O on flash fixes unclustered index access
 - 1 still suboptimal if I/O block size > record size
 - 1 therefore column stores profit much less than horizontal stores
 - 1 Random I/O useful to exploit secondary, tertiary table orderings
 - 1 the larger the data, the deeper clustering one can exploit





Even faster column scans on flash SSDs

1 New-generation SSDs

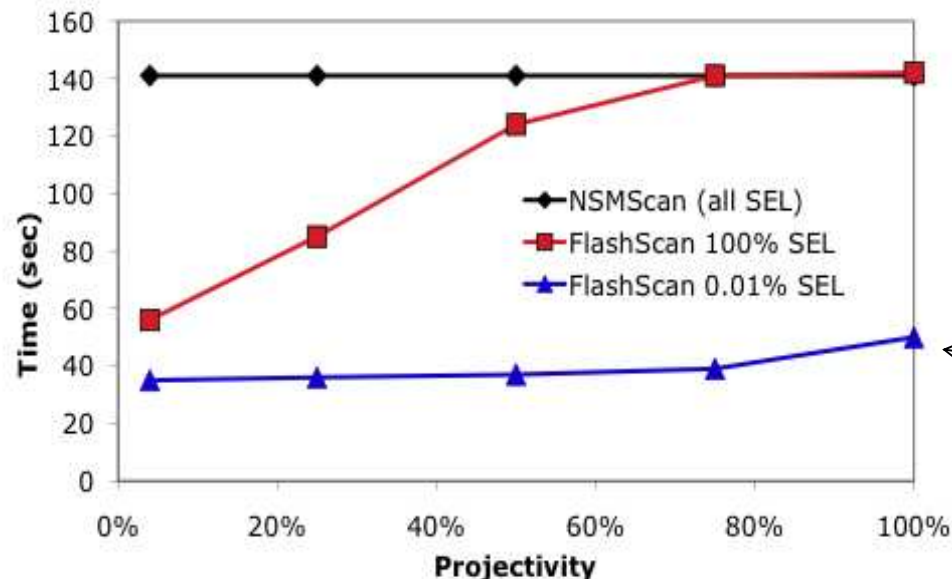
30K Read IOps, 3K Write Iops
250MB/s Read BW, 200MB/s Write

1 Very fast random reads, slower random writes

1 Fast sequential RW, comparable to HDD arrays

1 No expensive seeks across columns

1 FlashScan and Flashjoin: PAX on SSDs, inside Postgres



“Query Processing Techniques for Solid State Drives” Tsirogiannis, Harizopoulos, Shah, Wiener, Graefe, SIGMOD’09

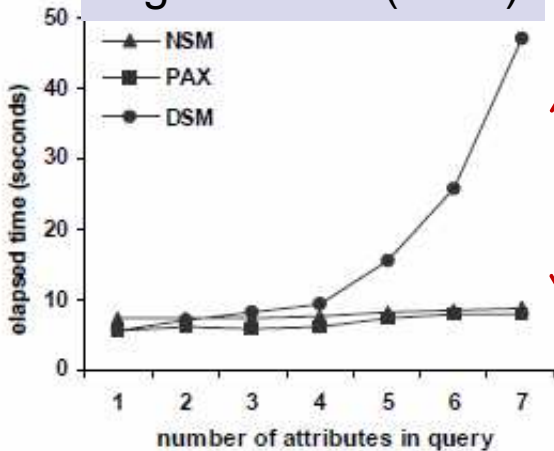
mini-pages with no qualified attributes are not accessed





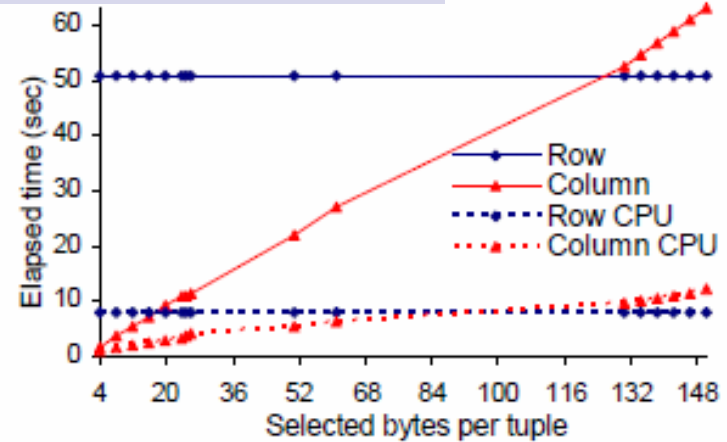
Column-scan performance over time

regular DSM (2001)



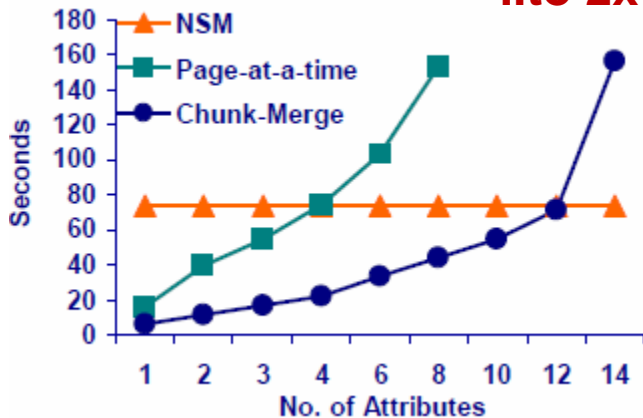
from 7x slower

column-store (2006)



..to 1.2x slower

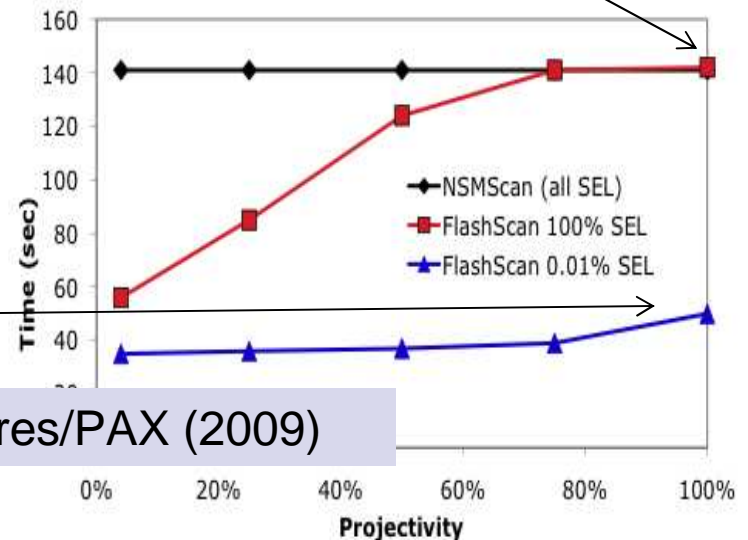
Lineitem (TPCH) 1GB



..to 2x slower

and 3x faster!

..to same



SSD Postgres/PAX (2009)

optimized DSM (2002)





Outline

- 1 Part 1: Basic concepts — *Stavros*
 - 1 Introduction to key features
 - 1 From DSM to column-stores and performance tradeoffs
 - 1 Column-store architecture overview
 - 1 Will rows and columns ever converge?
- 1 Part 2: Column-oriented execution — *Daniel*
- 1 Part 3: MonetDB/X100 and CPU efficiency — *Peter*





Architecture of a column-store

storage layout

- 1 read-optimized: dense-packed, compressed
- 1 organize in extends, batch updates
- 1 multiple sort orders
- 1 sparse indexes

engine

- 1 block-tuple operators
- 1 new access methods
- 1 optimized relational operators

system-level

- 1 system-wide column support
- 1 loading / updates
- 1 scaling through multiple nodes
- 1 transactions / redundancy



“C-Store: A Column-Oriented DBMS.” Stonebraker et al. VLDB 2005.



C-Store

- 1 Compress columns
- 1 No alignment
- 1 Big disk blocks
- 1 Only materialized views (perhaps many)
- 1 Focus on Sorting not indexing
- 1 Data ordered on anything, not just time
- 1 Automatic physical DBMS design
- 1 Optimize for grid computing
- 1 Innovative redundancy
- 1 Xacts – but no need for Mohan
- 1 Column optimizer and executor





C-Store: only materialized views (MVs)

- 1 **Projection** (MV) is some number of columns from a fact table
- 1 Plus columns in a dimension table – with a 1-n join between Fact and Dimension table
- 1 Stored in order of a storage key(s)
- 1 Several may be stored!
- 1 With a **permutation**, if necessary, to map between them
- 1 Table (as the user specified it and sees it) is not stored!
- 1 No secondary indexes (they are a one column sorted MV plus a permutation, if you really want one)

User view:

EMP (name, age, salary, dept)
Dept (dname, floor)

Possible set of MVs:

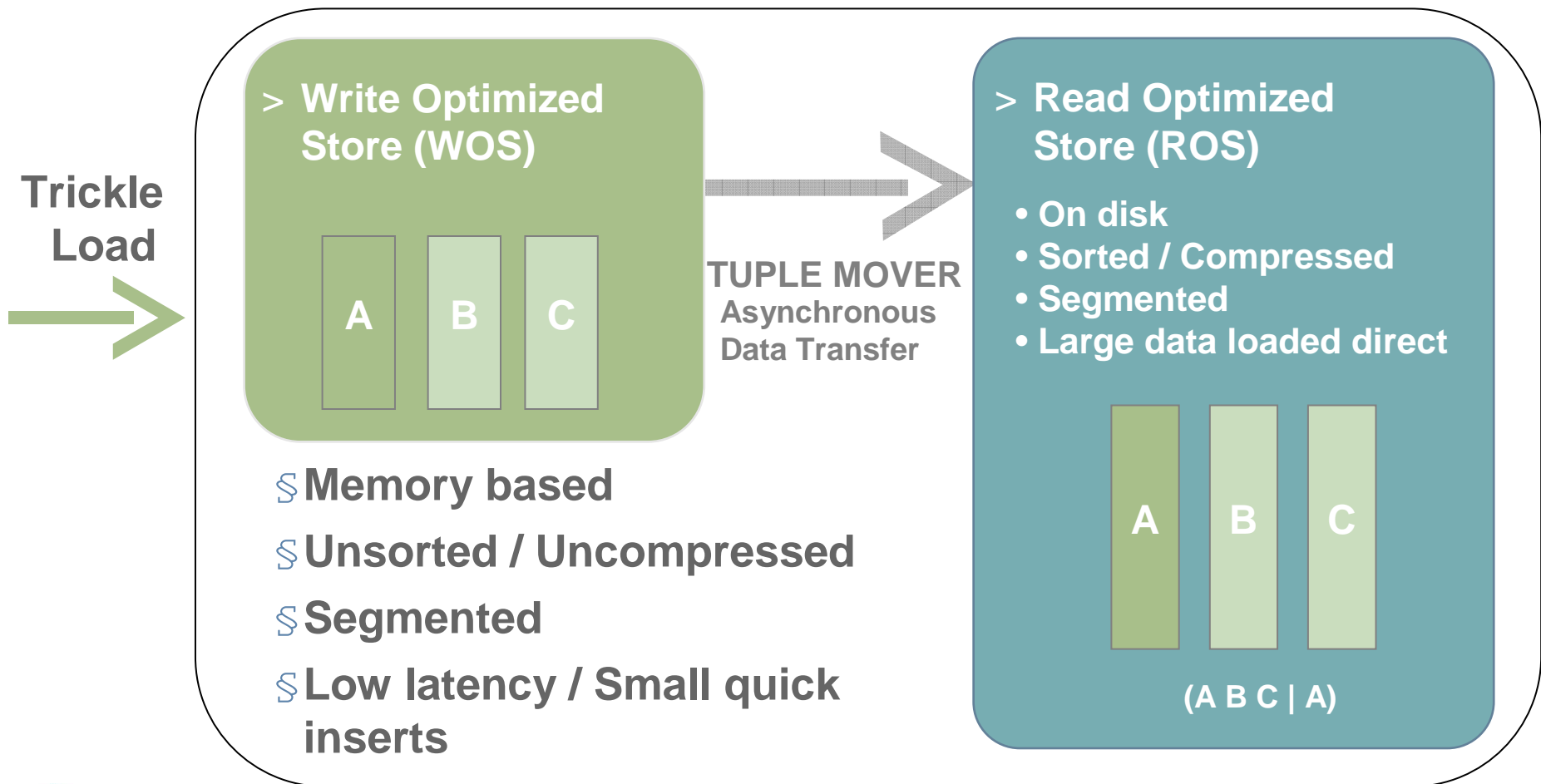
MV-1 (name, dept, floor) in floor order
MV-2 (salary, age) in age order
MV-3 (dname, salary, name) in salary order





Continuous Load and Query (Vertica)

Hybrid Storage Architecture





Loading Data (Vertica)

> INSERT, UPDATE, DELETE

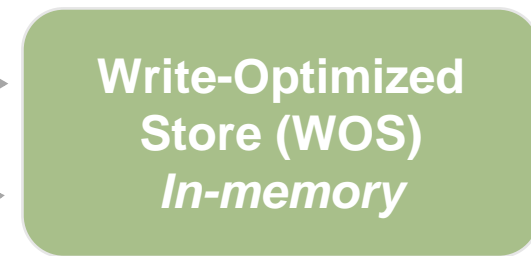
> Bulk and Trickle Loads

§ COPY

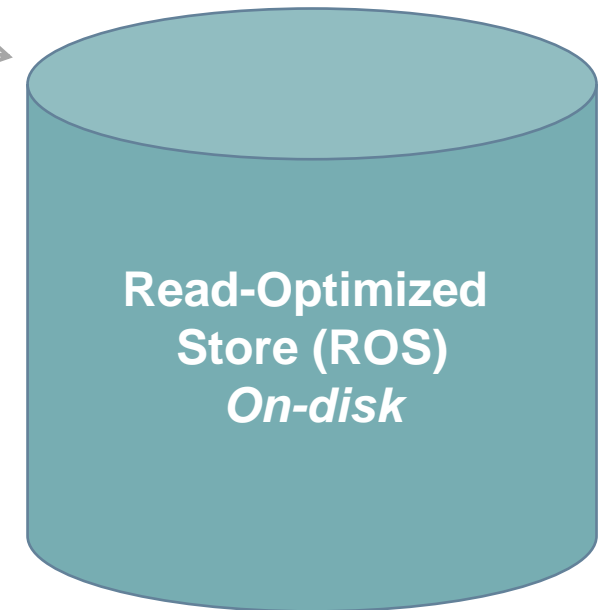
§ COPY DIRECT

> User loads data into logical Tables

> Vertica loads atomically into storage



Automatic
Tuple Mover





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

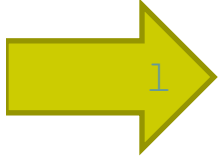
- 1 Cantor (history)
- 1 Sybase IQ
- 1 SenSage (former Addamark Technologies)
- 1 Kdb
- 1 1010data
- 1 MonetDB
- 1 C-Store/Vertica
- 1 X100/VectorWise
- 1 KickFire
- 1 SAP Business Accelerator
- 1 Infobright
- 1 ParAccel
- 1 Exasol





Outline

- 1 Part 1: Basic concepts — *Stavros*
 - 1 Introduction to key features
 - 1 From DSM to column-stores and performance tradeoffs
 - 1 Column-store architecture overview
 - 1 Will rows and columns ever converge?
- 1 Part 2: Column-oriented execution — *Daniel*
- 1 Part 3: MonetDB/X100 and CPU efficiency — *Peter*





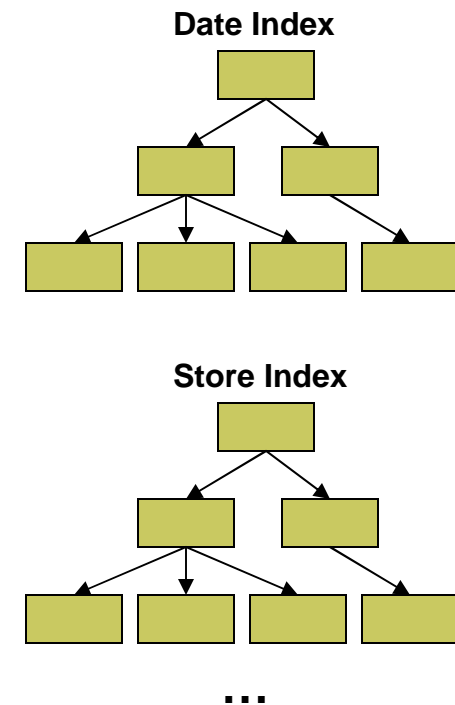
Simulate a Column-Store inside a Row-Store

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

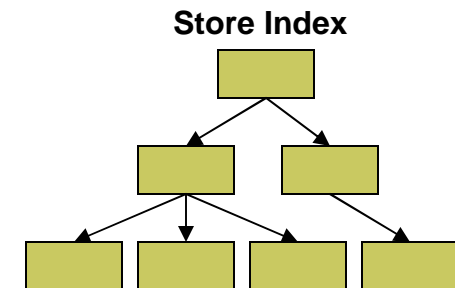
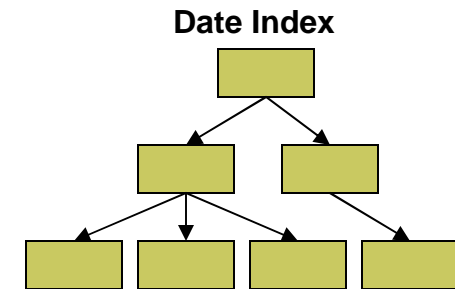
**Option A:
Vertical Partitioning**

Date		
Value	StartPos	Length
01/01	1	3

Can explicitly run-length encode date

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

**Option B:
Index Every Column**



...

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





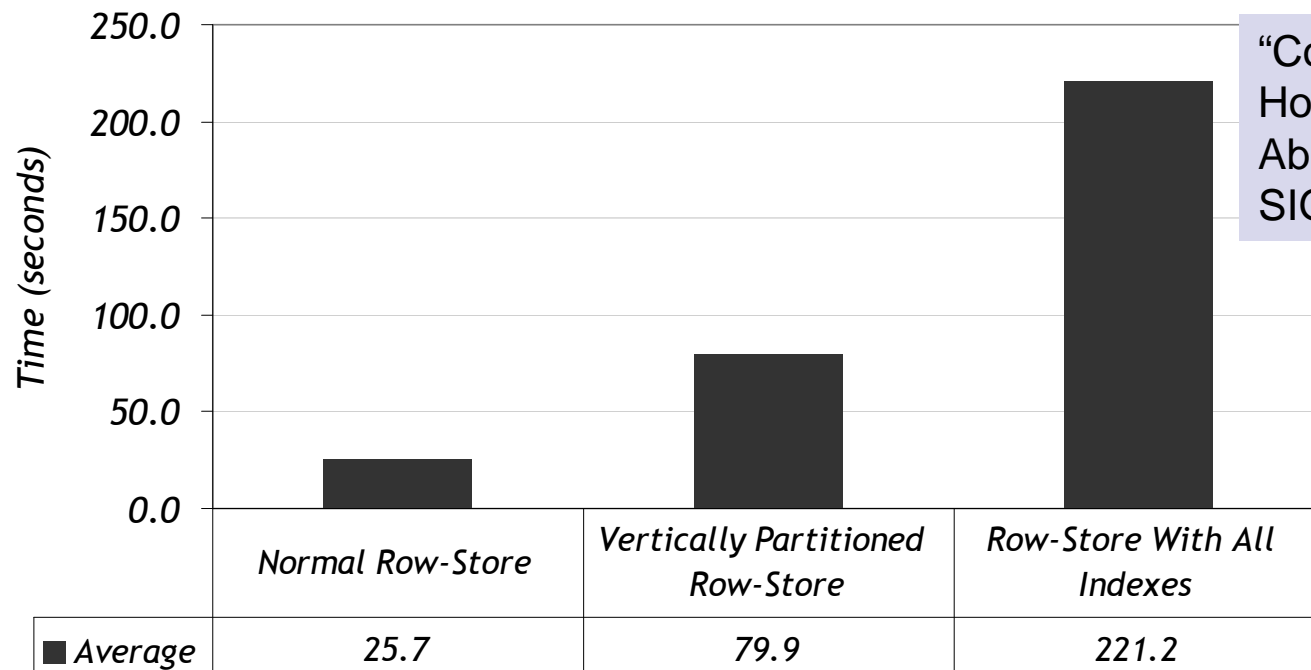
Experiments

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.





What's Going On? Vertical Partitions

- 1 Vertical partitions in row-stores:
 - 1 Work well when workload is known
 - 1 ..and queries access disjoint sets of columns
 - 1 See automated physical design

- 1 Do not work well as full-columns
 - 1 TupleID overhead significant
 - 1 Excessive joins

Tuple Header	TID	Column Data
	1	
	2	
	3	

Queries touch 3-4 foreign keys in fact table,
1-2 numeric columns

Complete fact table takes up ~4 GB
(compressed)

Vertically partitioned tables take up 0.7-1.1
GB (compressed)

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





What's Going On? All Indexes Case

1 Tuple construction

1 Common type of query:

```
SELECT store_name, SUM(revenue)
FROM Facts, Stores
WHERE fact.store_id = stores.store_id
      AND stores.country = "Canada"
GROUP BY store_name
```

1 Result of lower part of query plan is a set of TIDs that passed all predicates

1 Need to extract SELECT attributes at these TIDs

1 BUT: index maps value to TID

1 You really want to map TID to value (i.e., a vertical partition)

Tuple construction is SLOW





So....

- 1 All indexes approach is a poor way to simulate a column-store
- 1 Problems with vertical partitioning are NOT fundamental
 - 1 Store tuple header in a separate partition
 - 1 Allow virtual TIDs
 - 1 Combine clustered indexes, vertical partitioning
- 1 So can row-stores simulate column-stores?
 - 1 Might be possible, BUT:
 - 1 Need better support for vertical partitioning at the storage layer
 - 1 Need support for column-specific optimizations at the executor level
 - 1 Full integration: buffer pool, transaction manager, ..
 - 1 When will this happen?
 - 1 Most promising features = soon
 - 1 ..unless new technology / new objectives change the game (SSDs, Massively Parallel Platforms, Energy-efficiency)

See Part 2, Part 3
for most promising
features





End of Part 1

- 1 Basic concepts — *Stavros*
 - 1 Introduction to key features
 - 1 From DSM to column-stores and performance tradeoffs
 - 1 Column-store architecture overview
 - 1 Will rows and columns ever converge?

~~1 Part 2: Column-oriented execution — *Daniel*~~

~~1 Part 3: MonetDB/X100 and CPU efficiency — *Peter*~~

