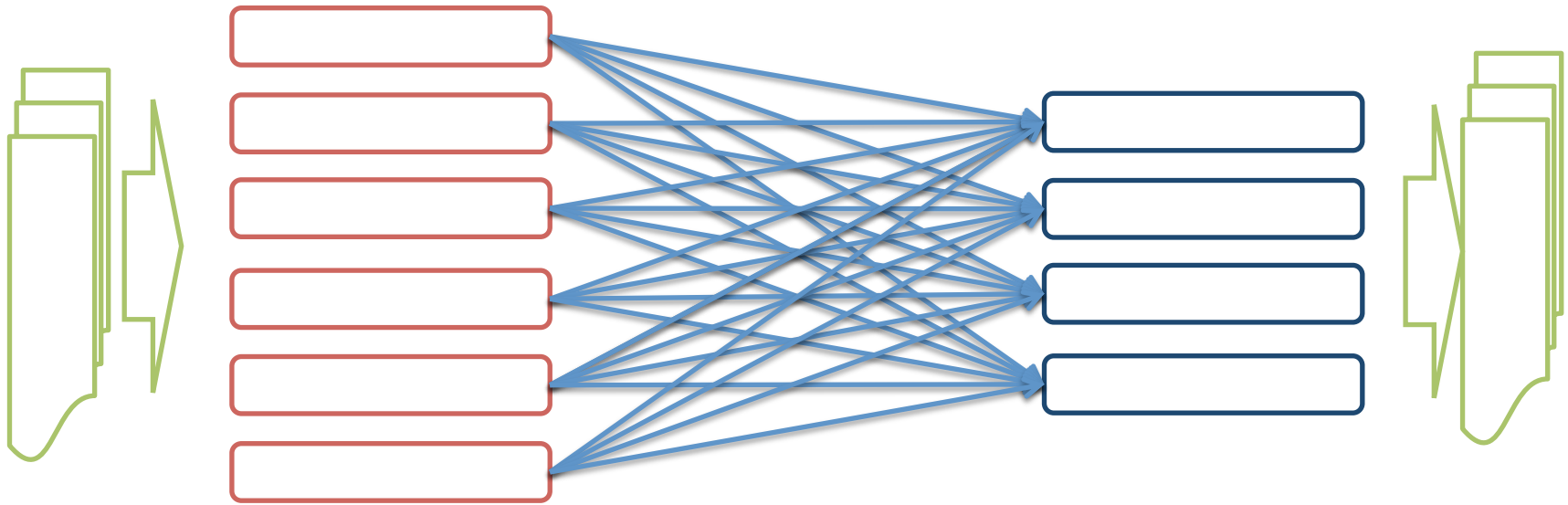
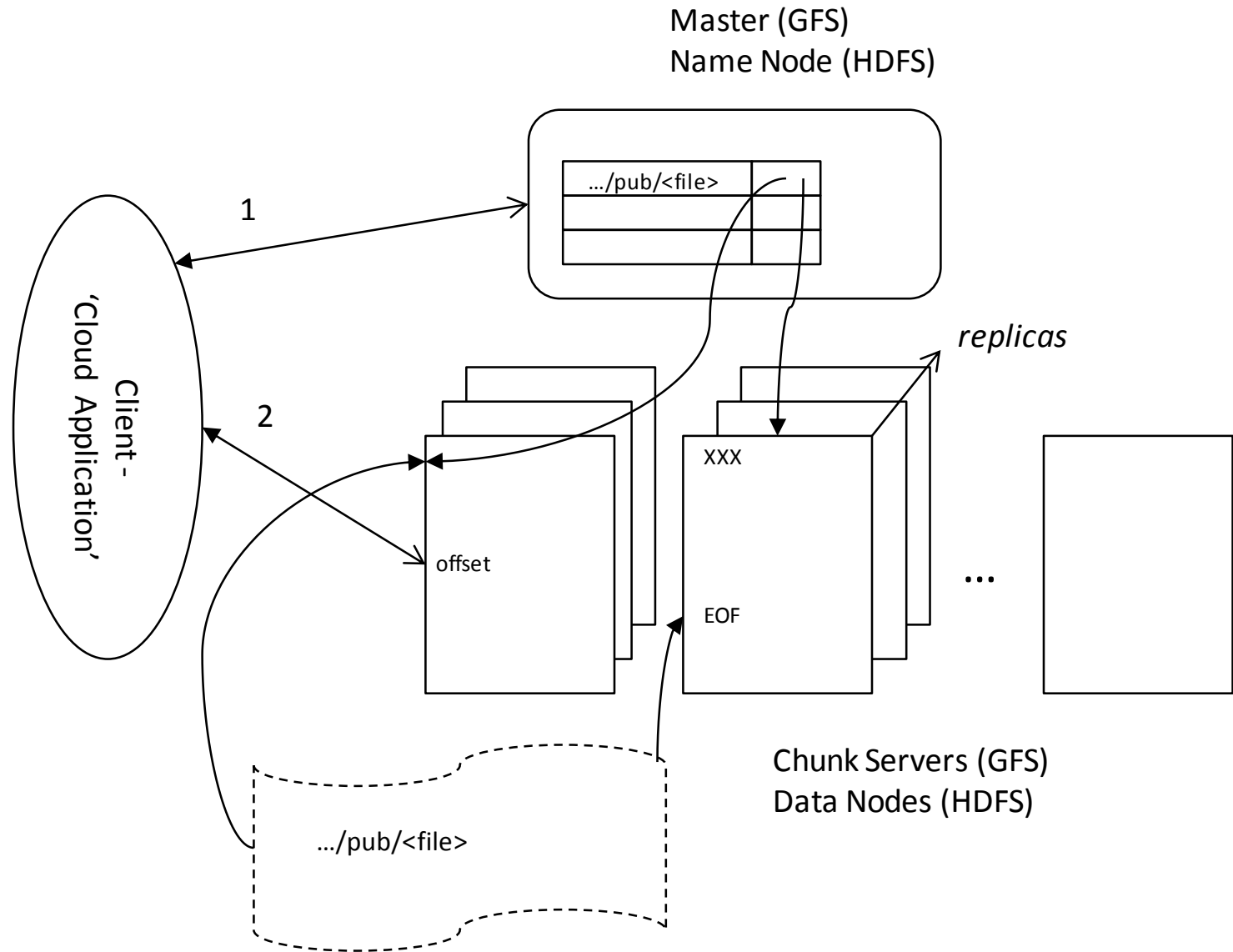


map-reduce input and output

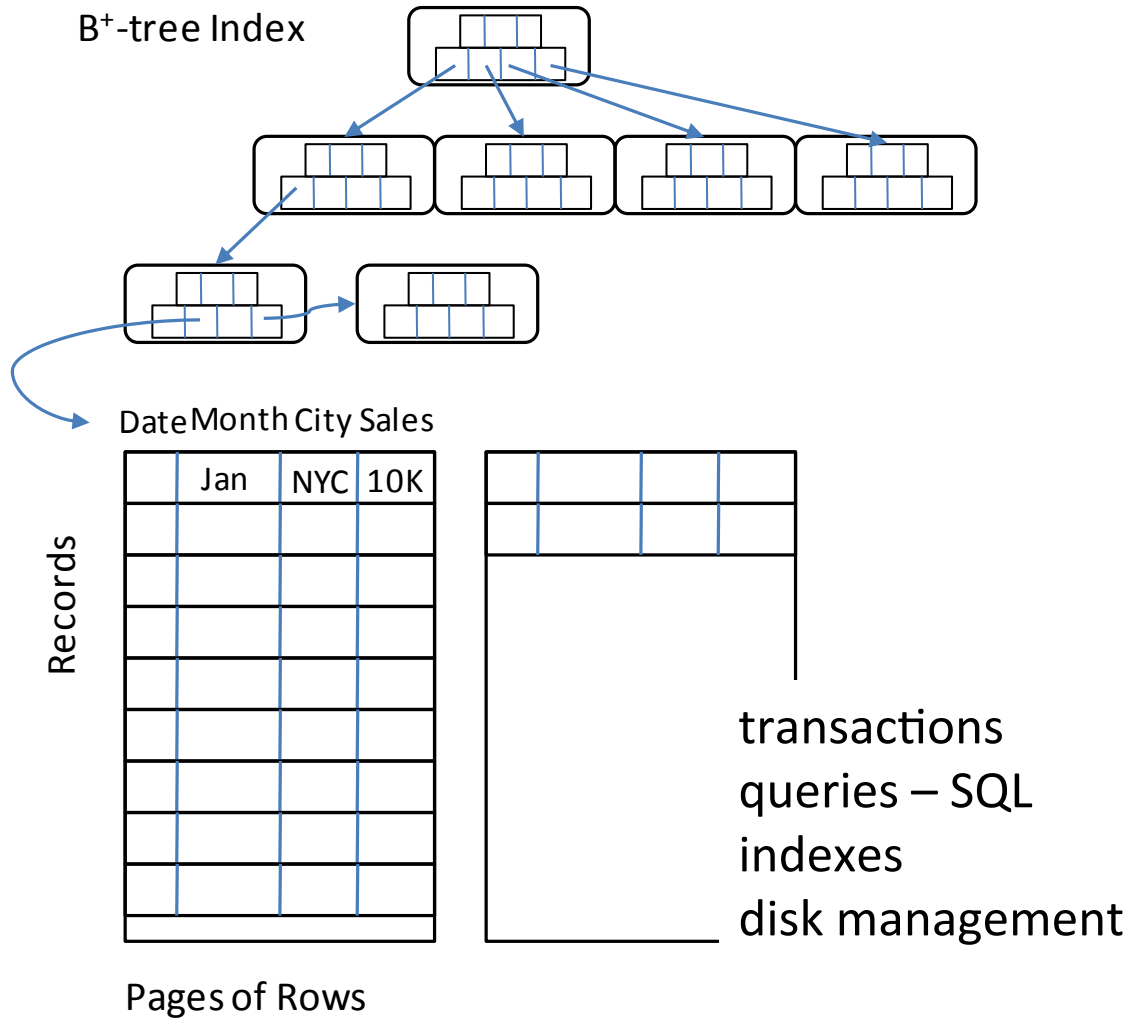


map-reduce reads data and *writes* fresh data –
from where - potential bottleneck ?
distributed data, in a distributed file system or database
parallel reading and writing
we have discussed processing-node failures; what about data?

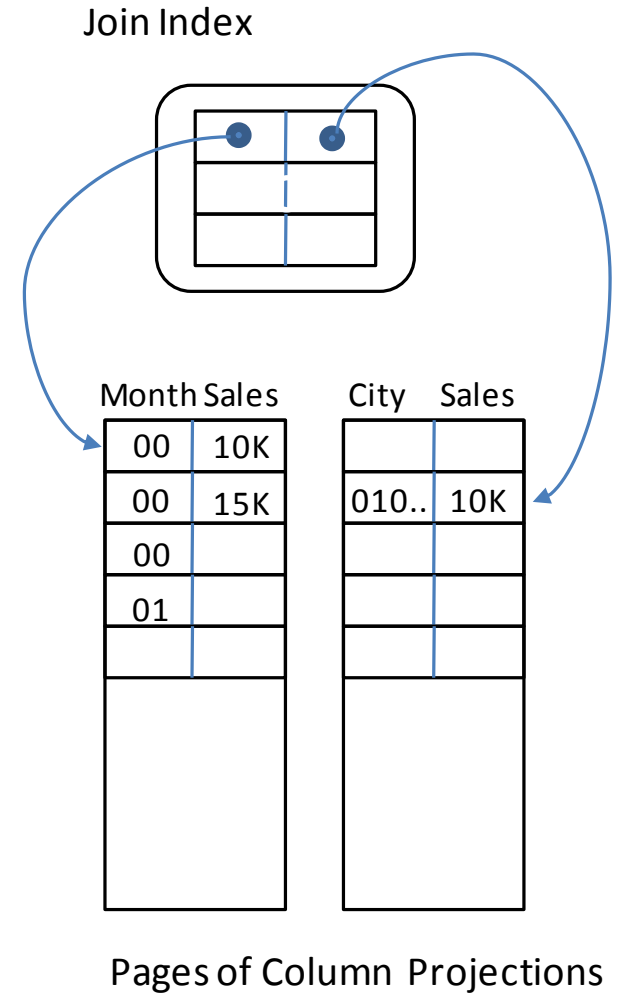
distributed file systems (GFS, HDFS)



overview of relational databases



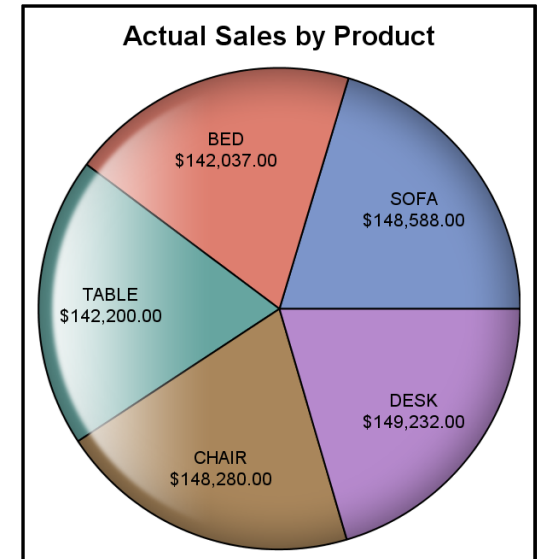
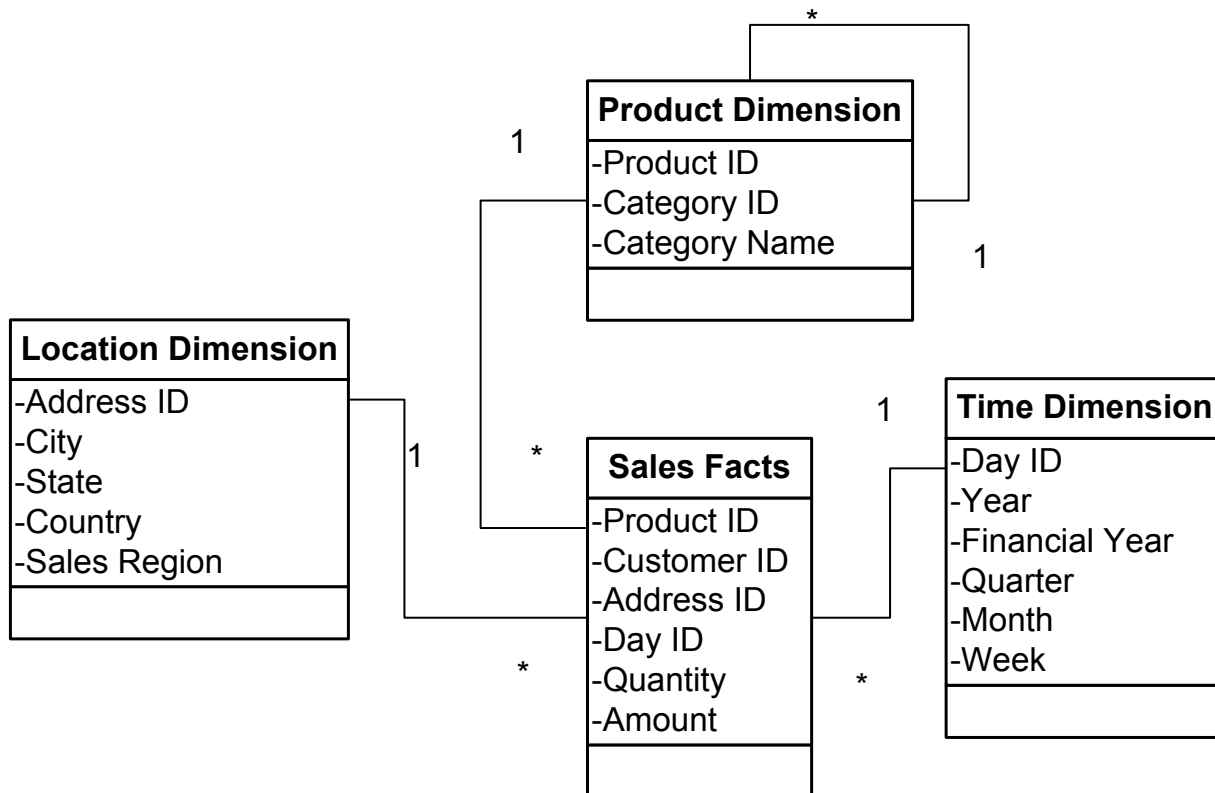
Row Oriented Database



Column Oriented Database

OLAP (“online analytical processing”)

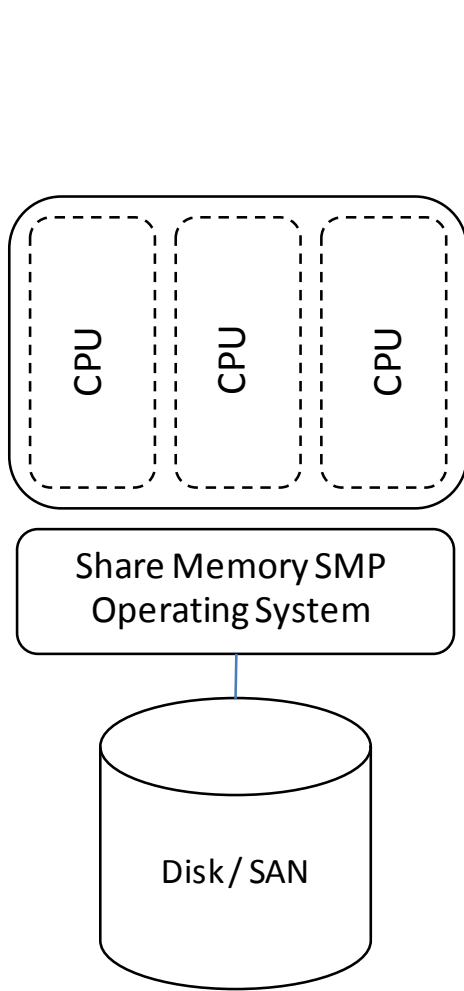
e.g.: select SUM(S.amount), S.pid, P.catname from S where S.did=T.did S.pid = P.pid and T.qrtr = 3 group by catname



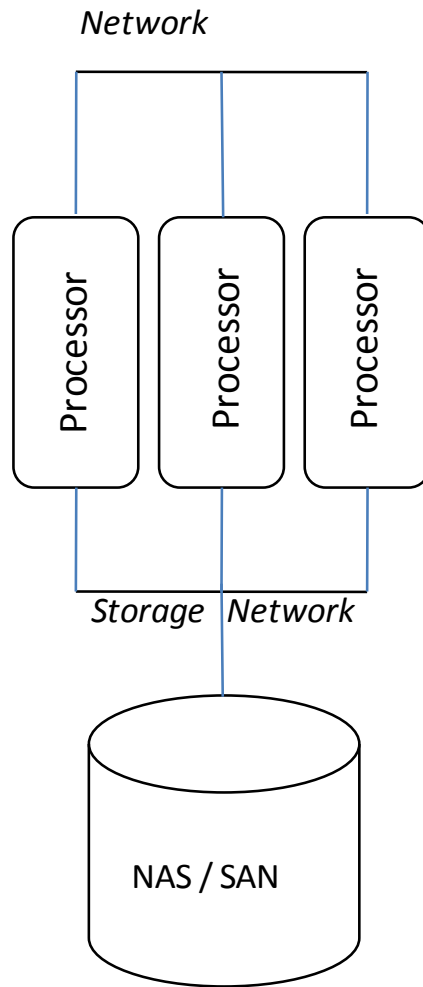
databases: why?

- transaction processing (ACID properties)
 - SQL – queries and indexing
-
- transaction processing *not* need for analytics
 - though there *may* be advantages in not having to move data out of a transaction store if avoidable
 - queries – yes, but if large volumes of data are being touched (e.g. joins, large-scale counting, building classifiers, etc.); indexes become *less* relevant
 - resilience to hardware failures, which MR provides, *is* vital.
 - *but* OLAP – can be viewed as computing a part of the joint distribution $P(f_1 \dots f_n)$ – using intuition to select

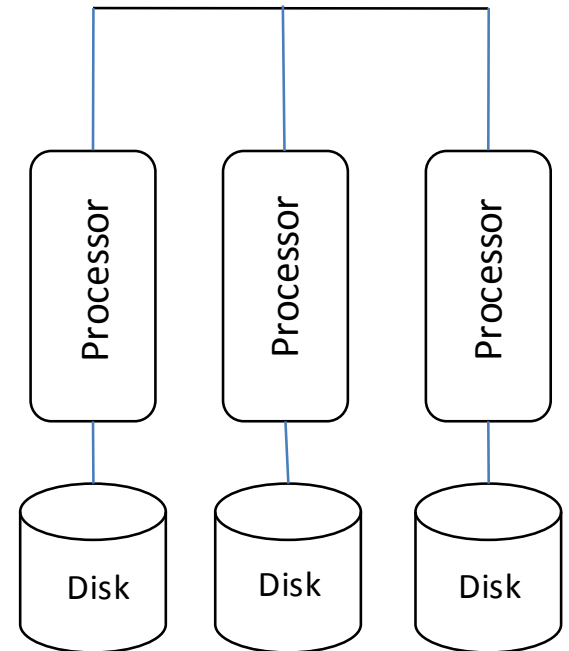
parallel databases



Shared Memory



Shared Disk



Shared Nothing

database evolution



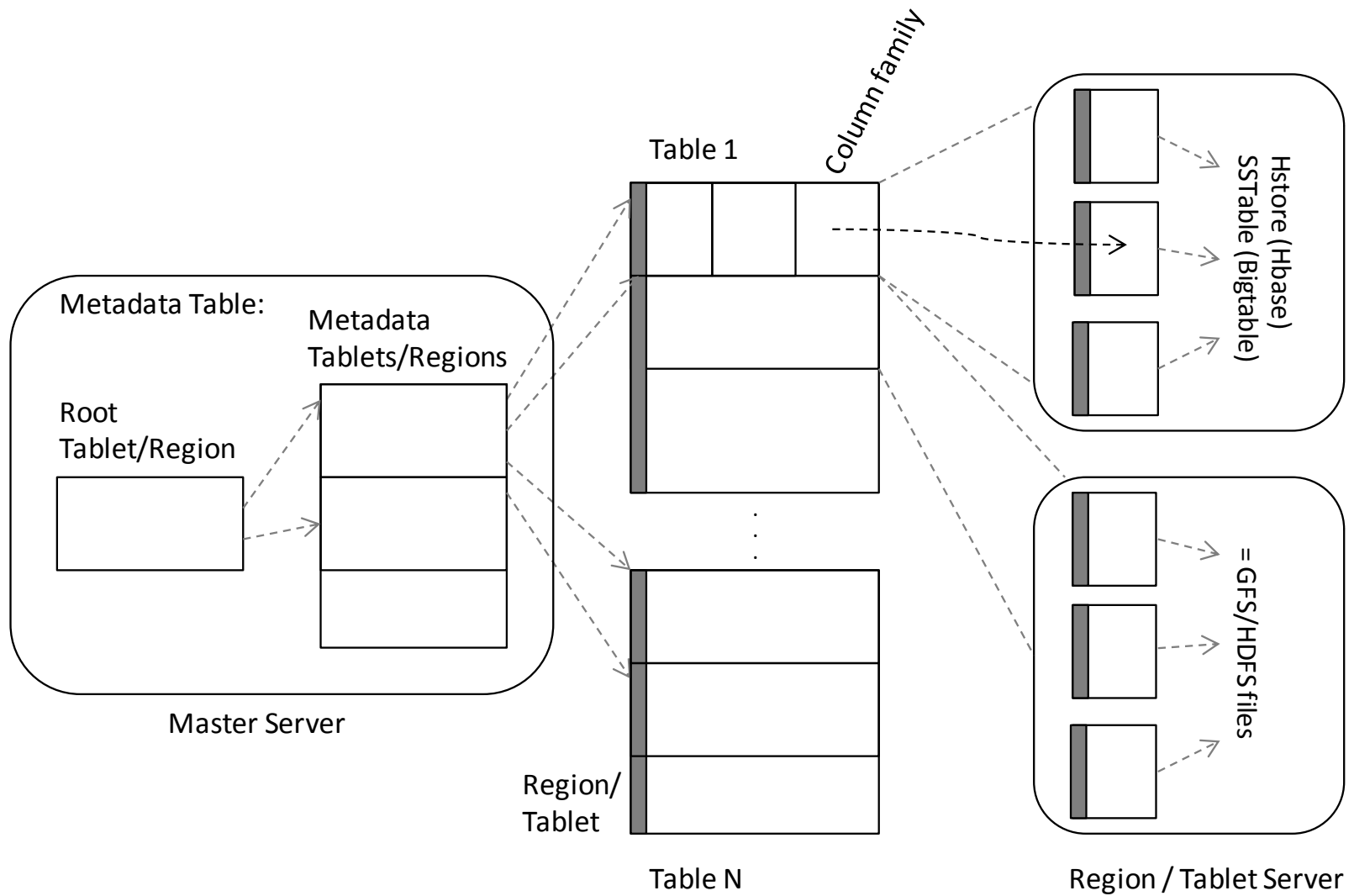
noSQL databases

- no ACID transactions
- *sharded* indexing
- restricted joins
- support columnar storage (if needed)

in-memory databases

- real-time transactions
- variety of indexes
- complex joins

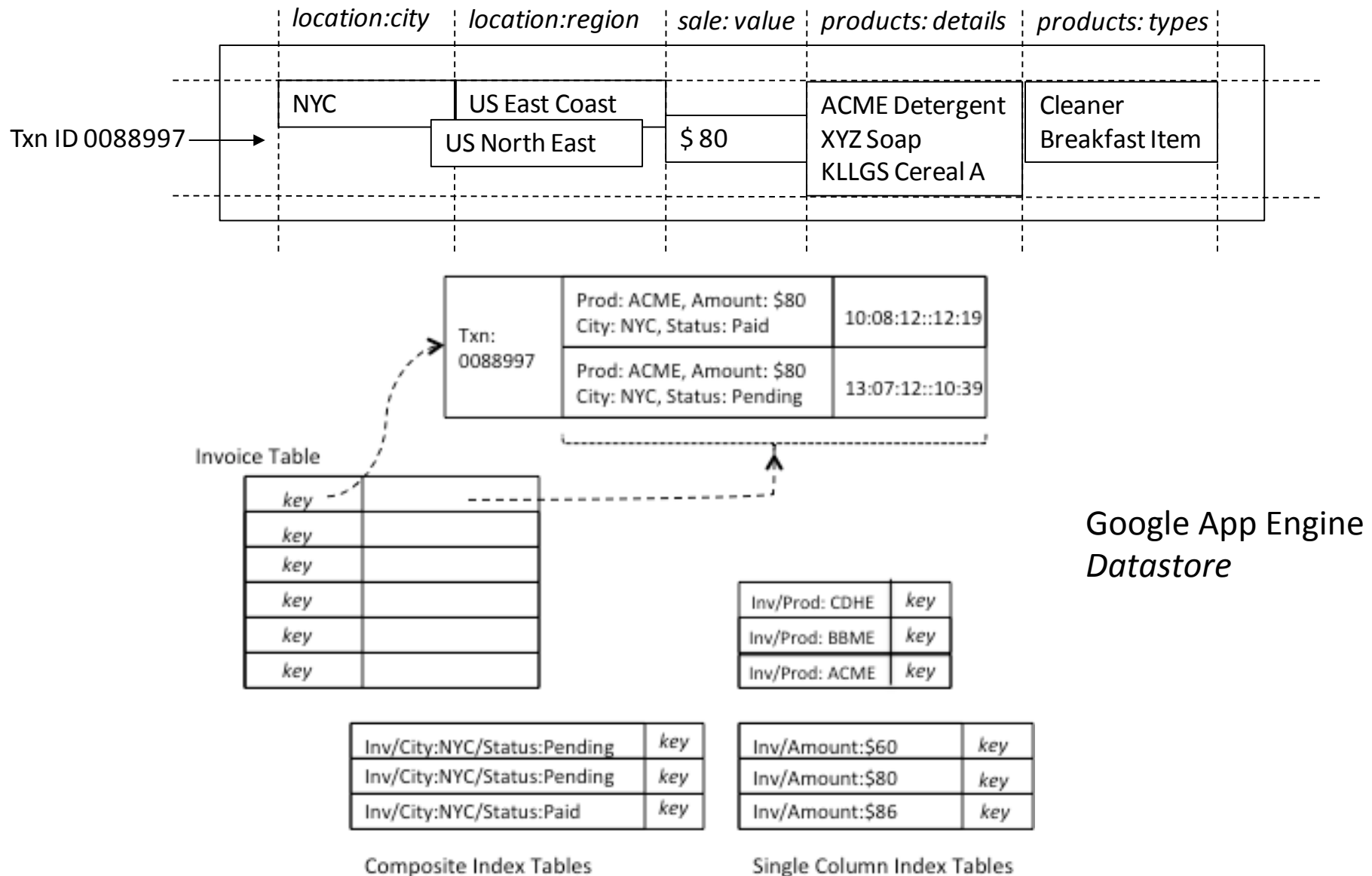
big-table (HBase)



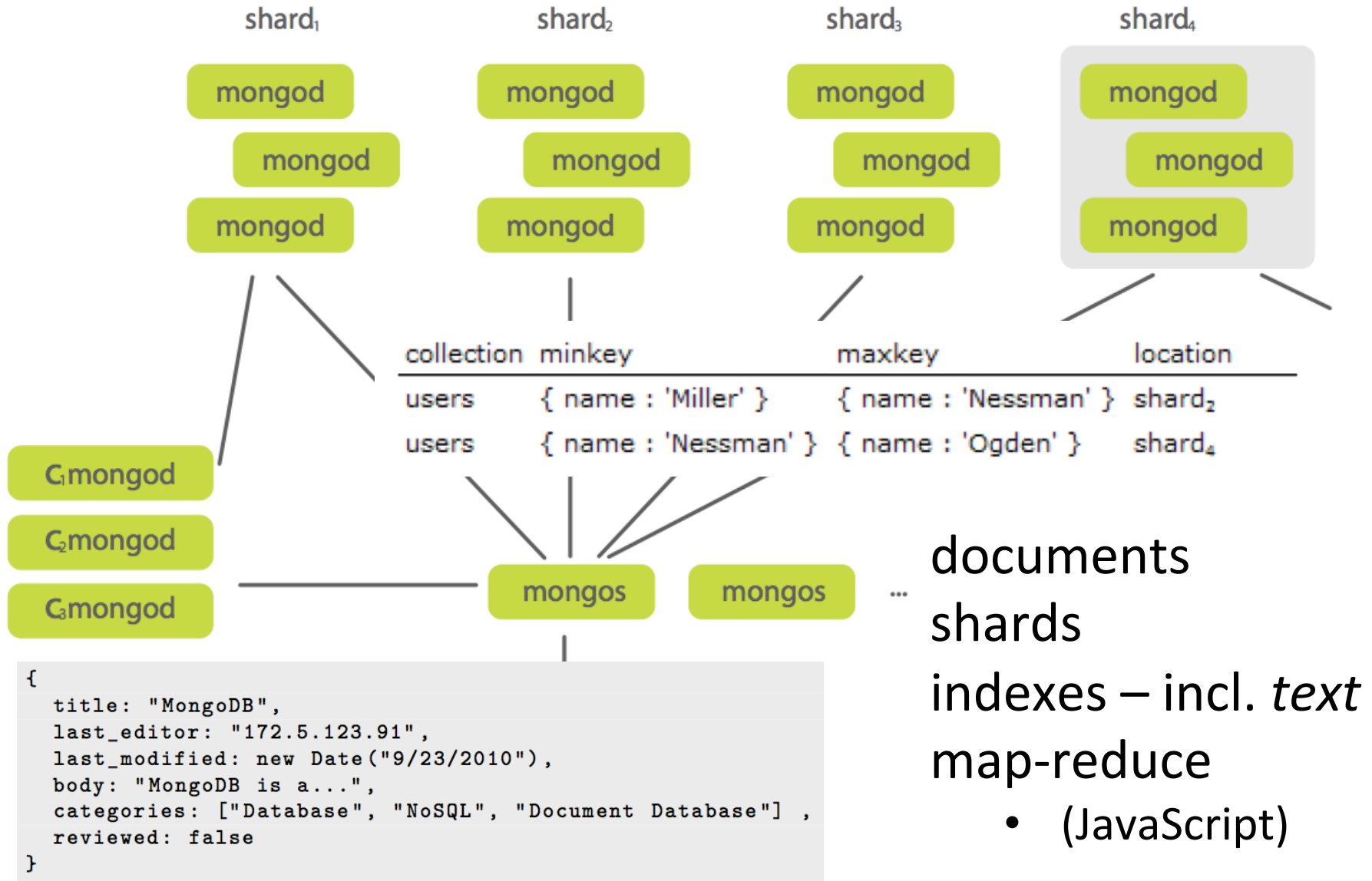
e.g. data in big-table

	<i>location:city</i>	<i>location:region</i>	<i>sale:value</i>	<i>products:details</i>	<i>products:types</i>
Txn ID 0088997 →	NYC	US East Coast	\$ 80	ACME Detergent	Cleaner
		US North East		XYZ Soap KLLGS Cereal A	Breakfast Item

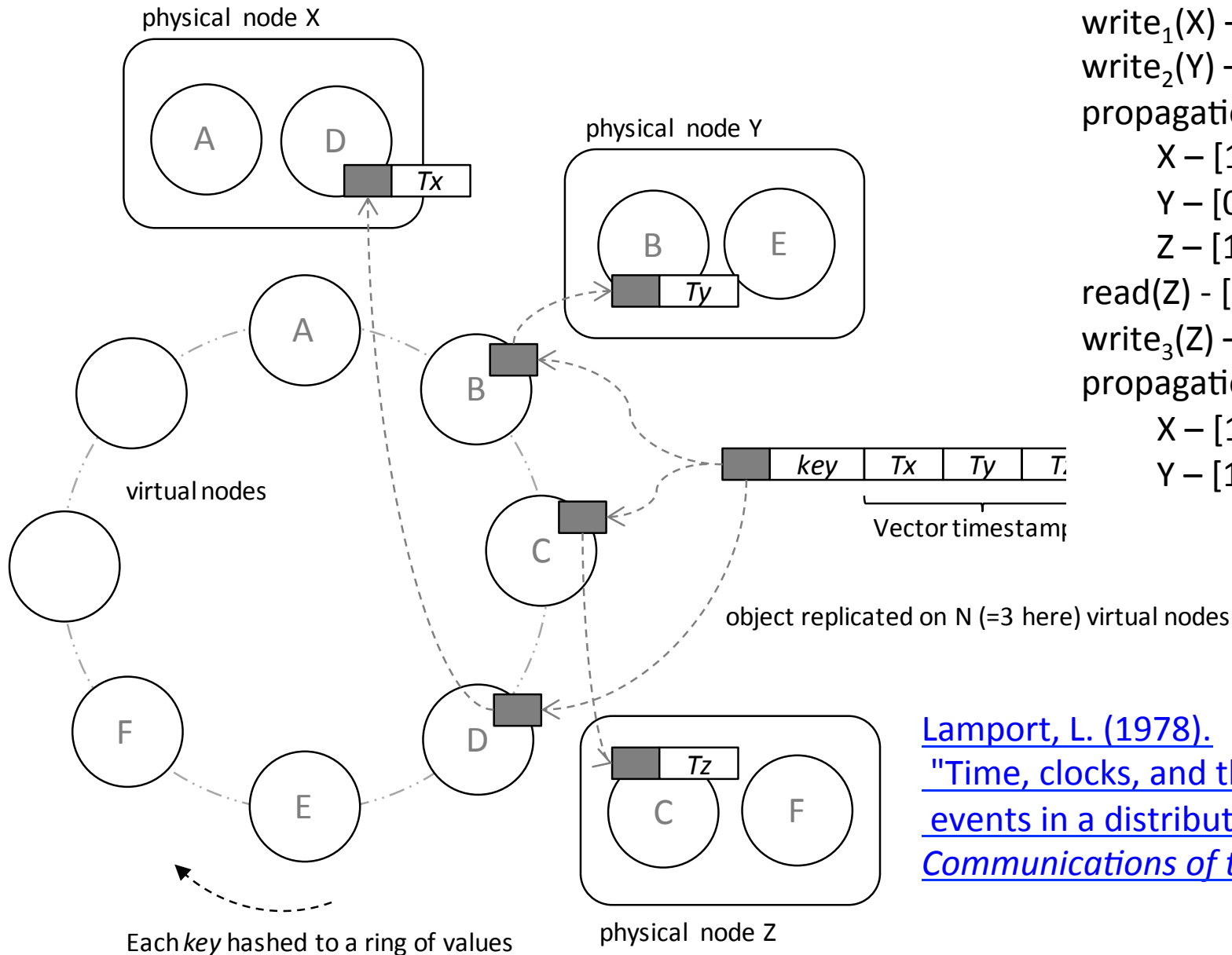
e.g. indexing *sharded* data (big-table)



mongo DB



eventual consistency



$\text{write}_1(X) - [1 \ 0 \ 0]$

$\text{write}_2(Y) - [0 \ 1 \ 0]$

propagation

$X - [1 \ 0 \ 0], [0 \ 1 \ 0]$

$Y - [0 \ 1 \ 0], [1 \ 0 \ 0]$

$Z - [1 \ 0 \ 0], [0 \ 1 \ 0]$

$\text{read}(Z) - [1 \ 0 \ 0], [0 \ 1 \ 0]$

$\text{write}_3(Z) - [1 \ 1 \ 1]$

propagation

$X - [1 \ 1 \ 1]$

$Y - [1 \ 1 \ 1]$

[Lamport, L. \(1978\).](#)

["Time, clocks, and the ordering of events in a distributed system".](#)

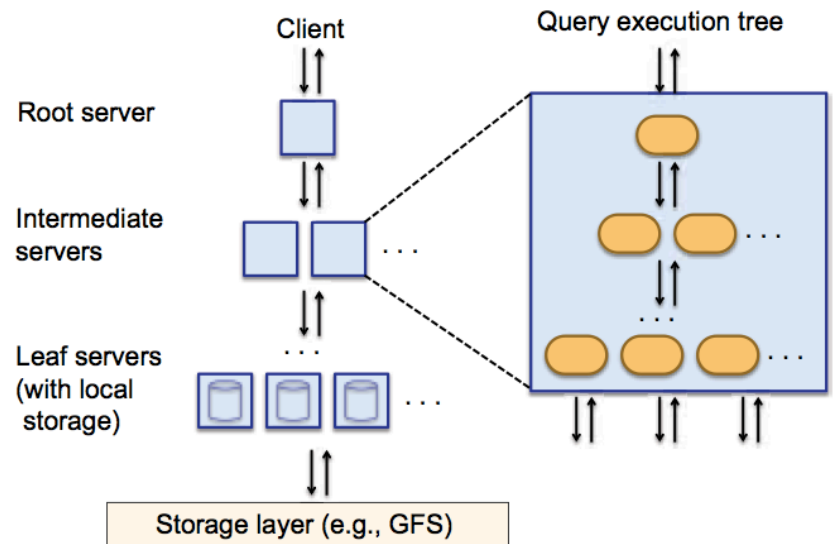
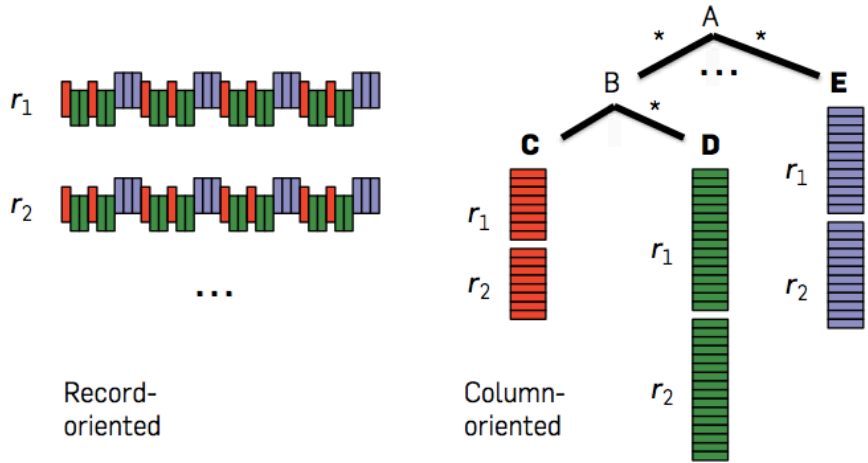
[Communications of the ACM](#)

Dremel – new ‘kid’ on the block?

powers Google’s “BigQuery”

two important innovations:

- columnar storage for *nested*, possibly non-unique fields – *leaf servers*
 - tree of *query servers* pass intermediate results from root to leaves and back
- orders of magnitude better than MR on petabytes of data – speed and storage



database evolution - summary

relational row-store: “one size fits all”; gigabytes

column-store data warehouse; terabytes

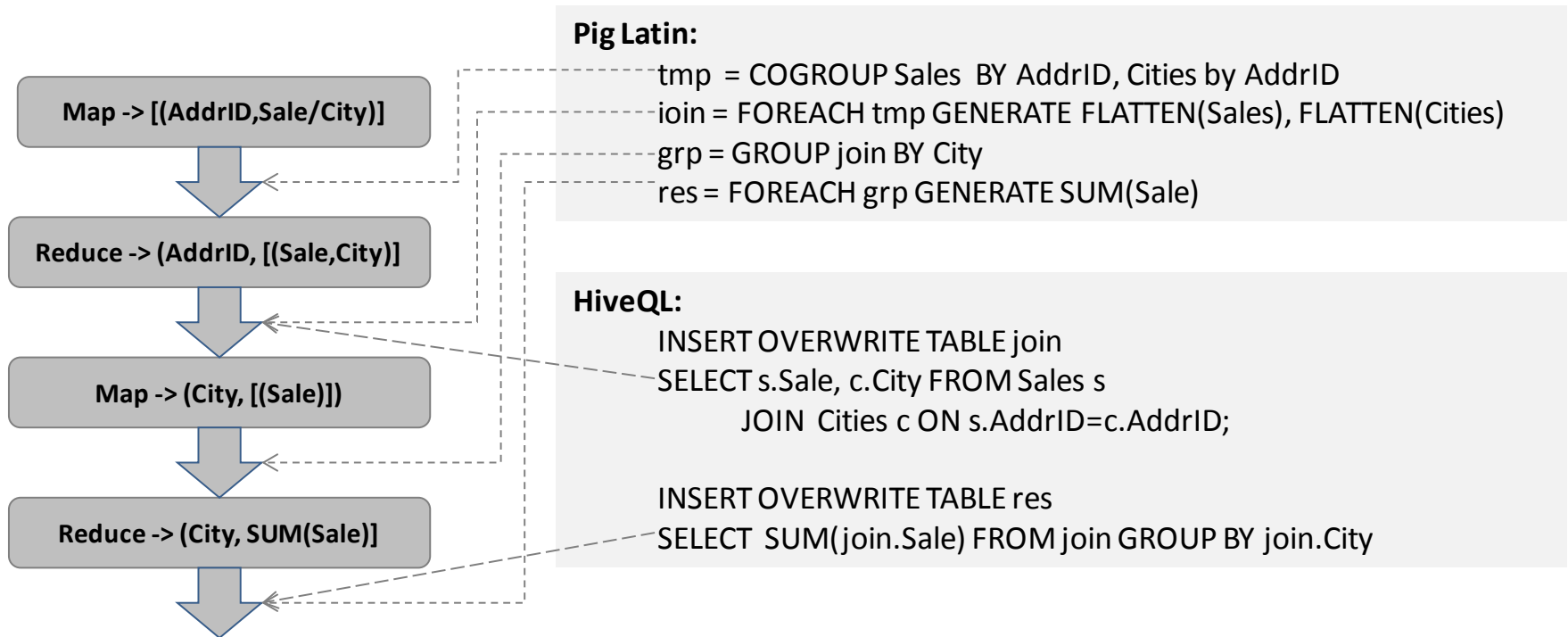
distributed noSQL row/column stores

+ map-reduce – for bulk analysis; 10s of terabytes

in-memory “one size fits all” databases; gigabytes

dremel – “one size fits all” for petabytes

SQL evolution: SQL-like MR coding



SQL evolution: in-DB statistics, in parallel

```
sql> SELECT * FROM training;
```

id	class	attributes
1	1	{1,2,3}
2	1	{1,2,1}
3	1	{1,4,3}
4	2	{1,2,2}
5	2	{0,2,2}
6	2	{0,1,3}

(6 rows)

```
sql> select * from toclassify;
```

id	attributes
1	{0,2,1}
2	{1,2,3}

(2 rows)

```
sql> SELECT madlib.create_nb_prepared_data_tables(  
'training', 'class', 'attributes', 3, 'nb_feature_probs', 'nb_class_priors');
```

```
sql> SELECT madlib.create_nb_probs_view (  
'nb_feature_probs', 'nb_class_priors', 'toclassify', 'id', 'attributes', 3, 'nb_probs_view_fast');
```

```
sql> SELECT * FROM nb_probs_view_fast;
```

key	class	nb_prob
1	1	0.4
1	2	0.6
2	1	0.75
2	2	0.25

(4 rows)

map-reduce evolution: iteration

many applications require repeated MR:

e.g. page-rank, continuous machine-learning ...

1. iterate MR

but make it more efficient: avoid data copy (HaLoop, Twister)

2. generalized data-flow graph of map->reduce tasks

tasks are 'blocking' for fault-tolerance (Dryad/LINQ, Hyracks ...)

3. direct implementation of recursion in MR

how to recover from non-blocking tasks failing?

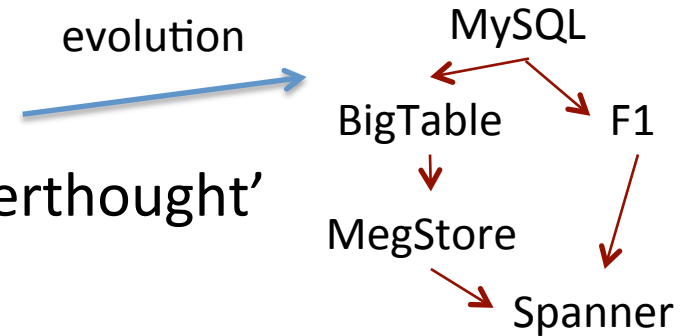
graph model: (Pregel, Giraph)

stream model: (S4)

Enterprise Data Management and Big Data

Databases

designed for transaction processing
reporting and analytics came later: 'afterthought'
(big data – exactly the reverse!)



Big Data Technologies

designed for analytics (counting, not queries)
for retrieval – inverted-index (unstructured) search
data captured in the 'raw' (logs), no transactions

- price-performance advantage
even for 'standard' ETL and OLAP (in the bulk)

[Analytical] Database evolution & trends



noSQL databases

- no ACID transactions**
- *sharded* indexing
- restricted joins
- support columnar storage

in-memory data / databases

- real-time transactions
 - variety of indexes, complex joins
- challenge:** *access patterns still matter*

relational row-store: “one size fits all”; gigabytes

column-store data warehouse; terabytes

distributed noSQL row/column stores

+ map-reduce – for bulk analysis; 10s of terabytes

in-memory “one size fits all” databases; gigabytes

Dremel (Impala*) – “one size fits all” for petabytes

for those who still want fast aggregation queries!

*Cloudera 2012

summary

- distributed files – 2nd basic element of big-data
- what databases are good for
 - and why traditional DBs were a happy compromise
- evolution of databases
- evolution of SQL
- evolution of map-reduce
- ‘big picture’ for database evolution