Advanced Data Management (CSCI 640/490)

Scalable Databases

Dr. David Koop





Data Integration: Combine Datasets with Different Data



D. Koop, CSCI 640/490, Spring 2023

NIU















D. Koop, CSCI 640/490, Spring 2023



4











D. Koop, CSCI 640/490, Spring 2023

Northern Illinois University

NIU







"Duplicate Detection" has many Duplicates

D. Koop, CSCI 640/490, Spring 2023





Northern Illinois University



"Duplicate Detection" has many Duplicates





Record Linkage Process











Record Linkage Techniques

- Deterministic matching
 - Rule-based matching (complex to build and maintain)
- Probabilistic record linkage [Fellegi and Sunter, 1969]
 - Use available attributes for linking (often personal information, like names, addresses, dates of birth, etc.)
 - Calculate match weights for attributes
- "Computer science" approaches
 - Based on machine learning, data mining, database, or information retrieval techniques
 - Supervised classification: Requires training data (true matches) - Unsupervised: Clustering, collective, and graph based

















Data Fusion

- Problem: Given a duplicate, create a single object representation while resolving conflicting data values.
- Difficulties:
 - Null values: Subsumption and complementation
 - Contradictions in data values
 - process
 - Metadata: Preferences, recency, correctness
 - Lineage: Keep original values and their origin
 - Implementation in DBMS: SQL, extended SQL, UDFs, etc.

- Uncertainty & truth: Discover the true value and model uncertainty in this





Conflict Resolution Strategies







	SI	S2	S3	S4	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UWisc	UWisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW







	SI	S2	S3	S4	S5	
Stonebraker	MIT	Berkeley	MIT	MIT	MS	
Dewitt	MSR	MSR	UWisc	UWisc	UWisc	
Bernstein	MSR	MSR	MSR	MSR	MSR	
Carey	UCI	AT&T	BEA	BEA	BEA	
Halevy	Google	Google	UW	UW	UW	







	SI	S2	S3	S4	S5
Stonebraker	MIT	Berkeley	MIT	MIT	MS
Dewitt	MSR	MSR	UWisc	UWisc	UWisc
Bernstein	MSR	MSR	MSR	MSR	MSR
Carey	UCI	AT&T	BEA	BEA	BEA
Halevy	Google	Google	UW	UW	UW

D. Koop, CSCI 640/490, Spring 2023

2. With only a snapshot it is hard to decide which source is a copier.









I. Sharing common data does not in itself imply copying.

	SI	S2	S3	S4	S5	
Stonebraker	MIT	Berkeley	MIT	MIT	MS	
Dewitt	MSR	MSR	UWisc	UWisc	UWisc	
Bernstein	MSR	MSR	MSR	MSR	M\$R	
Carey	UCI	AT&T	BEA	BEA	BEA	
Halevy	Google	Google	UW	UW	UW	
	3. A copier can also provide or verify some data by itself, so it is inappropriate to ignore all of its data.					

D. Koop, CSCI 640/490, Spring 2023

2. With only a snapshot it is hard to decide which source is a copier.







deas

- If two sources share a lot of false values, they are more likely to be dependent.
- highly different from the accuracy of S1.

• S1 is more likely to copy from S2, if the accuracy of the common data is











Combining Accuracy and Dependence

Source-accuracy Computation

D. Koop, CSCI 640/490, Spring 2023







Northern Illinois University





Combining Accuracy and Dependence

Source-accuracy Computation

Step 3

D. Koop, CSCI 640/490, Spring 2023



Step









The Motivating Example

	SI	S2	S3	S4	S5	
Stonebraker	MIT	Berkeley	MIT	MIT	MS	
Dewitt	MSR	MSR	UWisc	UWisc	UWisc	
Bernstein	MSR	MSR	MSR	MSR	MSR	
Carey	UCI	AT&T	BEA	BEA	BEA	
Halevy	Google	Google	UW	UW	UW	
S_2 Q_2 Q_2 Q_2 Q_2 Q_3 S_3 $Rnd 2$ S_2 Q_4 Q_4 Q_5 Q_5 Q_5 S_4 Q_4 Q_5 Q_5 S_5 S_4 S_5 S_5 S_1 S_1 S_2 S_1 S_2 S_2 S_1 S_2 S_2 S_1 S_2 S_2 S_3 S_5 S_1 S_1 S_2 S_1 S_2 S_1 S_2 S_1 S_2 S_1 S_2 S_1 S_2 S_2 S_1 S_2 S_2 S_1 S_2 S_2 S_1 S_2 S_1 S_2 S_1 S_2 S_2 S_2 S_1 S_2 S_2 S_2 S_1 S_2 S_2 S_1 S_2 S_2 S_2 S_2 S_2 S_1 S_2 $S_$						
	Rnd 3	 Rnd II	S ₂ S ₄	49 S ₃ 49 44 55 55 S ₅	[X L Dong	













The Motivating Example

Accuracy	SI	S2	S3	S4	S5
Round I	.52	.42	.53	.53	.53
Round 2	.63	.46	.55	.55	.55
Round 3	.71	.52	.53	.53	.37
Round 4	.79	.57	.48	.48	.31
• • •					
Round 11	.97	.61	.40	.40	.21

Value		Carey	Halevy		
Confidence	UCI	AT&T	BEA	Google	UW
Round I	1.61	1.61	2.0	2.1	2.0
Round 2	1.68	1.3	2.12	2.74	2.12
Round 3	2.12	1.47	2.24	3.59	2.24
Round 4	2.51	1.68	2.14	4.01	2.14
Round 11	4.73	2.08	1.47	6.67	1.47







Assignment 4

- Data Integration & Data Fusion
- Out soon





Paper Critique

- Read <u>What's Really New with NewSQL?</u>
- Submit critique **before class** on Wednesday, March 22
- Discussion ideas:
 - What are the advantages or disadvantages of NewSQL vs NoSQL?
 - Are they really different from standard RDBMS?
 - Which category of NewSQL databases is most exciting?





Scalable Database Systems







Database Architecture













Database Architecture



D. Koop, CSCI 640/490, Spring 2023





Northern Illinois University







Database Architecture



D. Koop, CSCI 640/490, Spring 2023





Northern Illinois University







NoSQL Motivation

Scalability



D. Koop, CSCI 640/490, Spring 2023

Impedance Mismatch









Relational Database Architecture



D. Koop, CSCI 640/490, Spring 2023

[Hellerstein et al., <u>Architecture of a Database System</u>]









Relational Databases: One size fits all?

- Lots of work goes into relational database development:
 - B-trees
 - Cost-based query optimizers
 - ACID (Atomicity, Consistency, Isolation, Durability)
- Vendors largely stuck with this model from the 1980s through 2000s Having different systems leads to business problems:
- - cost problem
 - compatibility problem
 - sales problem
 - marketing problem

[Stonebraker and Cetinetmel, 2005]









ACID Transactions

- Make sure that transactions are processed reliably
- Atomicity: leave the database as is if some part of the transaction fails (e.g. don't add/remove only part of the data) using rollbacks
- Consistency: database moves from one valid state to another
- Isolation: concurrent execution matches serial execution
- Durability: endure hardware failures, make sure changes hit disk









How to Scale Relational Databases?







Shared Nothing Architecture



D. Koop, CSCI 640/490, Spring 2023

Shift towards higher distribution & less coordination:









TrafficDB: Shared-Memory Data Store

- Traffic-aware route planning
- Want up-to-date data for all
- Thousands of requests per second
 - High-Frequency Reads
 - Low-Frequency Writes
- "Data must be stored in a region of RAM that can be shared and efficiently accessed by *several* different application processes"





Northern Illinois University



Parallel DB Architecture: Shared Nothing



D. Koop, CSCI 640/490, Spring 2023



[Hellerstein et al., Architecture of a Database System]








Sharding













Stonebraker: The End of an Architectural Era

- "RDBMSs were designed for the business data processing market, which is their sweet spot"
- "They can be beaten handly in most any other market of significant enough size to warrant the investment in a specialized engine"
- Changes in markets (science), necessary features (scalability), and technology (amount of memory)
- RDBMS Overhead: Logging, Latching, and Locking
- Relational model is not necessarily the answer
- SQL is not necessarily the answer









OLTP vs. OLAP

- data entry and retrieval transactions
- OLTP Examples:
 - Add customer's shopping cart to the database of orders
 - Find me all information about John Hammond's death
- OLTP is focused on the day-to-day operations while Online Analytical Processing (OLAP) is focused on analyzing that data for trends, etc.
- OLAP Examples:

 - Find the average amount spent by each customer - Find which year had the most movies with scientists dying

Online Transactional Processing (OLTP) often used in business applications,







Row Stores



by	movie_name
	The Black Hole
tty	Blade Runner
Jr	Jurassic Park
	Star Trek: TNG
chine	Forbidden Planet
	Terminator 2: Judgment Day
	[J. Swanhart, Introduction to Columr









Inefficiency in Row Stores for OLAP

select sum(metric) as the_sum from fact

1. Storage engine gets a whole row from the table





D. Koop, CSCI 640/490, Spring 2023



2. SQL interface extracts only requested portion, adds it to "the_sum"

[J. Swanhart, Introduction to Column Stores]









Column Stores



Each column has a file or segment on disk

D. Koop, CSCI 640/490, Spring 2023

	Person	Genre
oubtfire	Robin Williams	Comedy
	Roy Scheider	Horror
У	Jeff Goldblum	Horror
Magnolias	Dolly Parton	Drama
rdcage	Nathan Lane	Comedy
rokovitch	Julia Roberts	Drama
K	7	

[J. Swanhart, Introduction to Column Stores]









Horizontal Partitioning vs. Vertical Partitioning

Original Table

CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR
1	TAEKO	OHNUKI	BLUE
2	O.V.	WRIGHT	GREEN
3	SELDA	BAĞCAN	PURPLE
4	JIM	PEPPER	AUBERGINE









Horizontal Partitioning vs. Vertical Partitioning

Vertical Partitions

VP1

VP2

CUSTOMER ID	FIRST NAME	LAST NAME	CUSTOMER ID	FAVORITE COLOR
1	TAEKO	OHNUKI	1	BLUE
2	O.V .	WRIGHT	2	GREEN
3	SELDA	BAĞCAN	3	PURPLE
4	JIM	PEPPER	4	AUBERGINE

Original Table				
CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR	
1	TAEKO	OHNUKI	BLUE	
2	O.V .	WRIGHT	GREEN	
3	SELDA	BAĞCAN	PURPLE	
4	JIM	PEPPER	AUBERGINE	

D. Koop, CSCI 640/490, Spring 2023

Horizontal Partitions

HP1

CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR
1	TAEKO	OHNUKI	BLUE
2	O.V .	WRIGHT	GREEN

HP2

CUSTOMER ID	FIRST NAME	LAST NAME	FAVORITE COLOR
3	SELDA	BAĞCAN	PURPLE
4	JIM	PEPPER	AUBERGINE









NoSQL Paradigm Shift



Commercial DBMS

Specialized DB hardware (Oracle Exadata, etc.)

Highly available network (Infiniband, Fabric Path, etc.)

Highly Available Storage (SAN, RAID, etc.)

D. Koop, CSCI 640/490, Spring 2023

Open-Source DBMS

Commodity hardware

Commodity network (Ethernet, etc.)

Commodity drives (standard HDDs, JBOD)







Problems with Relational Databases

	1	\cap	n	1
ID.		U	U	

Customer: Ann

Line Items:

0321293533	2	\$48	\$96
0321601912	1	\$39	\$39
0131495054	1	\$51	\$51

Payment Details:

Card: Amex **CC Number:** 12345 Expiry: 04/2001











NoSQL Classification Criteria











Key-Value Stores

Data model: (key) -> value Interface: CRUD (Create, Read, Update, Delete)



Examples: Amazon Dynamo (AP), Riak (AP), Redis (CP)







Key-Value Stores

- Always use primary-key access
- Operations:
 - Get/put value for key
 - Delete key

>	<key=customerid></key=customerid>
>	<value=object></value=object>
	Customer
	BillingAddress
	Orders
	Order
	ShippingAddress
	OrderPayment
	OrderItem Product









Wide-Column Stores

Data model: (rowkey, column, timestamp) -> value Interface: CRUD, Scan



Examples: Cassandra (AP), Google BigTable (CP), HBase (CP)







Column Stores

- Instead of having rows grouped/sharded, we group columns
- ... or families of columns
- Put similar columns together













Document Stores



Examples: CouchDB (AP), RethinkDB (CP), MongoDB









Document Stores

- Documents are the main entity
 - Self-describing
 - Hierarchical
 - Do not have to be the same
- Could be XML, JSON, etc.
- Key-value stores where values are "examinable"
- Can have query language and indices overlaid

D. Koop, CSCI 640/490, Spring 2023

<Key=CustomerID>

```
"customerid": "fc986e48ca6"
"customer":
"firstname": "Pramod",
"lastname": "Sadalage",
"company": "ThoughtWorks",
"likes": [ "Biking","Photography" ]
"billingaddress":
{ "state": "AK",
  "city": "DILLINGHAM",
   "type": "R"
```







Graph Databases

Data model: G = (V, E): Graph-Property Modell Interface: Traversal algorithms, querys, transactions



Examples: Neo4j (CA), InfiniteGraph (CA), OrientDB









Graph Databases

- Focus on entities and relationships
- Edges may have properties
- Relational databases required a set traversal
- Traversals in Graph DBs are faster











NoSQL Classification Criteria













CAP Theorem











CAP Theorem

- Consistency: every read would get you the most recent write Availability: every node (if not failed) always executes queries Partition tolerance: system continues to work even if nodes are down • Theorem (Brewer): It is impossible for a distributed data store to simultaneously provide more than two of Consistency, Availability, and

- Partition Tolerance









CAP Theorem "Proof"

possible



Network partition

D. Koop, CSCI 640/490, Spring 2023

• If there is a network partition, one of consistency or availability will not be









NoSQL Techniques













Distributing Data

- Aggregate-oriented databases
- Sharding (horizontal partitioning): Sharding distributes different data across multiple servers, so each server acts as the single source for a subset of data
- Replication: Replication copies data across multiple servers, so each bit of data can be found in multiple places. Replication comes in two forms,
 - Source-replica replication makes one node the authoritative copy that handles writes, replica synchronizes with the source and may handle reads. - Peer-to-peer replication allows writes to any node; the nodes coordinate to synchronize their copies of the data.











Sharding













Sharding Approaches

- Hash-based Sharding
 - Hash of data values (e.g. key) determines partition (shard)
 - Pro: Even distribution, Con: No data locality



D. Koop, CSCI 640/490, Spring 2023

[D. DeWitt & J. Gray, 1992, via F. Gessert, Image: MongoDB]









Sharding Approaches

- Range-based Sharding
 - Assigns ranges defined over fields (shard keys) to partitions - Pro: Enables Range Scans & Sorting, Con: Repartitioning/balancing reg'd









Sharding Approaches

- Entity-Group Sharding
 - Explicit data co-location for single-node-transactions



D. Koop, CSCI 640/490, Spring 2023

- Pro: Enables ACID Transactions, Con: Partitioning not easily changable

[D. DeWitt & J. Gray, 1992, via F. Gessert, Image: J. Kim]







Replication

- Store N copies of each data item
- Consistency model: synchronous vs. asynchronous
- Coordination: Multiple Primary, Primary/Replica















Replication: When

- Asynchronous (lazy)
 - Writes are acknowledged immdediately
 - Performed through log shipping or update propagation
 - Pro: Fast writes, no coordination needed
 - Con: Replica data potentially stale (inconsistent)
- Synchronous (eager)
 - The node accepting writes synchronously propagates updates/transactions before acknowledging
 - Pro: Consistent
 - Con: needs a commit protocol (more roundtrips), unavailable under certain network partitions











Replication: Where

- Primary-Replica (Primary Copy)
 - Only a dedicated primary is allowed to accept writes, replicas are read-replicas
 - Pro: reads from the primary are consistent - Con: primary is a bottleneck and SPOF
- Multi-Primary (Update anywhere)
 - The server node accepting the writes synchronously propagates the update or transaction before acknowledging
 - Pro: fast and highly-available
 - Con: either needs coordination protocols (e.g. Paxos) or is inconsistent









Consistency Levels









Slides: Introduction to Cassandra

Robert Stupp





What is Cassandra?

- Fast Distributed (Column Family NoSQL) Database
 - High availability
 - Linear Scalability
 - High Performance
- Fault tolerant on Commodity Hardware
- Multi-Data Center Support
- Easy to operate
- Proven: CERN, Netflix, eBay, GitHub, Instagram, Reddit









Cassandra and CAP








Cassandra: Ring for High Availability



D. Koop, CSCI 640/490, Spring 2023









Next Class's Paper Critique

- Read What's Really New with NewSQL?
- Submit critique **before class** on Wednesday, March 22
- Discussion ideas:
 - What are the advantages or disadvantages of NewSQL vs NoSQL?
 - Are they really different from standard RDBMS?
 - Which category of NewSQL databases is most exciting?

D. Koop, CSCI 640/490, Spring 2023



