



THE UNIVERSITY OF
MELBOURNE

INFO20003: Database Systems

Assignment Project Exam Help

<https://powcoder.com>

Dr Renata Borovica-Gajic
Add WeChat powcoder

Lecture 21
NoSQL Databases

Week 11



- By the end of this session, you should be able to:
 - Define what Big Data is
 - Describe why databases go beyond relational DBs
 - Understand why we need NoSQL
 - Types of NoSQL
 - CAP theorem

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

* material in this lecture is drawn from <http://martinfowler.com/books/nosql.html>, including talk at GOTO conference 2012 and Thoughtworks article at <https://www.thoughtworks.com/insights/blog/nosql-databases-overview>

Much of business data is tabular



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

EMP_RECORD...	EMP_ID	EMP_REGION	EMP_DEPT	EMP_HIRE_DATE	EMP_SALARY	EMP_NAME	EMP_SKIL
68	3715	4	153	09061987	9 6 1987 14000000	IRENE HIRSH	041085
62	39412	1	650	03119590	3 11 9590 167000000	ANN FAHEY	031099
56	1939	2	265	09281988	9 28 1988 21300000	EMILY WILM...	021077
50	3502	2	165	07041985	7 4 1985 19500000	CATHEZINE ...	011015
44	4435	2	117	05141989	5 14 1989 17000000	AGNES KING	00
68	1673	3	138	07021985	7 2 1985 16800000	MARTIN XU	041033
62	4181	3	161	02031988	2 3 1988 15900000	JOHN DURN	030045
56	1443	1	265	12028900	12 2 8900 6000000	PAT DUNN	021055
50	3607	3	127	08072000	8 7 2000 18300000	ANDREA HIN...	011014
44	1775	3	288	02051989	2 5 1989 2700000	PETER JONES	00
68	1209	2	165	05121986	5 12 1986 17300000	DIDRA WILK...	041065



- Pros of relational databases
 - simple, can capture (nearly) any business use case
 - can integrate multiple applications via shared data store
 - standard interface language SQL
 - ad-hoc queries, across and within "data aggregates"
 - fast, reliable, concurrent, consistent
- Cons of relational databases
 - Object Relational (OR) impedance mismatch
 - not good with big data
 - not good with clustered/replicated servers
- Adoption of NoSQL driven by “cons” of Relational
- but ‘polyglot persistence’ = Relational will not go away

One business object (in aggregate form) is stored across many relational tables.

Xin Cube Inc.
300 Francisco St.
San Francisco
CA 94133 US
Tel: (415) 959-1155 Fax: (415) 959-2285
Email: admin@xincube.com
Website: www.xincube.com

Invoice

Bill To: John Symax Inc.
123 AA Juanita Ave.
Glendora CA 91740 US

Ship To: John Symax Inc.
123 AA Juanita Ave.
Glendora CA 91740 US

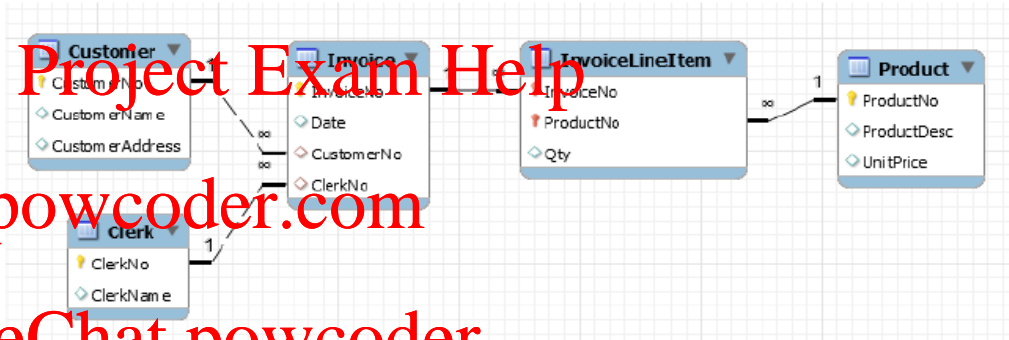
Date	14-Aug-2008	Order No		Sales Person	Charles Wooten
Shipping Date	13-Aug-2008	Shipping Terms			COD
ID	SKU / Description	Unit Price (USD)	Qty	Amount (USD)	
PGV680.065	AMD Athlon X2DO-7450 2.4GHz/1GB/160GB-SMP-DVD/VS	800.00	6.00	4,800.00	
PGV680.037	PDC ES300 2.8GHz/1GB/120GB-SMP-DVD/FDD/VS	645.00	4.00	2,580.00	
LCV800.082	LG 18.5" WLCD	230.00	10.00	2,300.00	
HPQ754.071	HP LaserJet R200	1,100.00	1.00	1,100.00	
Sub Total (USD):				10,780.00	
Discount (USD):				0.00	
Sales Tax (USD):				790.70	
Shipping (USD):				0.00	
Total (USD):				10,243.70	
Deposit (USD):				0.00	
Amount Due (USD):				10,243.76	

Notes:
All Payments must be made only in the form of a crossed cheque or cash payable to Xin Cube Inc.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



This enables analytical queries like:

```

select productno, sum(qty) from
InvoiceLineItem
group by productno;
    
```

But there is a lot of work to disassemble and reassemble the aggregate.



Data in Aggregate form: Examples of JSON and XML

JSON Example

```
{"products": [  
  {"number": 1, "name": "Zoom X", "Price": 10.00},  
  {"number": 2, "name": "Wheel Z", "Price": 7.50},  
  {"number": 3, "name": "Spring 10", "Price": 12.75}  
]}
```

JavaScript Object
Notation

XML Example

```
<products>  
  <product>  
    <number>1</number> <name>Zoom X</name> <price>10.00</price>  
  </product>  
  <product>  
    <number>2</number> <name>Wheel Z</name> <price>7.50</price>  
  </product>  
  <product>  
    <number>3</number> <name>Spring 10</name> <price>12.75</price>  
  </product>  
</products>
```

<https://powcoder.com>

Extensible Markup
Language

Add WeChat powcoder



- Data that exist in very large volumes and many different varieties (data types) and that need to be processed at a very high velocity (speed).
 - **Volume** – much larger quantity of data than typical for relational databases
 - **Variety** – lots of different data types and formats
 - **Velocity** – data comes at very fast rate (e.g. mobile sensors, web click stream)

- Schema on Read, rather than Schema on Write
 - Schema on Write— preexisting data model, how traditional databases are designed (relational databases)
 - Schema on Read – data model determined later, depends on how you want to use it (XML, JSON)
 - Capture and store the data, and worry about how you want to use it later
- Data Lake
 - A large integrated repository for internal and external data that does not follow a predefined schema
 - Capture everything, dive in anywhere, flexible access

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Schema on write vs. schema on read

Schema on Write

Requirements gathering and structuring

Formal data modeling process

Database schema

Database use based on the predefined schema

Traditional
database
design

Schema on Read

Collecting large amounts of data with
locally defined structures (e.g., using JSON/XML)

Storing the data in a data lake

Analyzing the stored data to identify
meaningful ways to structure it

Structuring and organizing the data
during the data analysis process

The big data
approach

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



- Features

- Doesn't use relational model or SQL language
- Runs well on distributed servers
- Most are open-source
- Built for the modern web
- Schema-less (though there may be an "implicit schema")
- Supports schema on read
- Not ACID compliant
- 'Eventually consistent'

- Goals

- to improve programmer productivity (OR mismatch)
- to handle larger data volumes and throughput (big data)

from NoSQL Databases: An Overview
by Pramod Sadalage, Thoughtworks(2014)

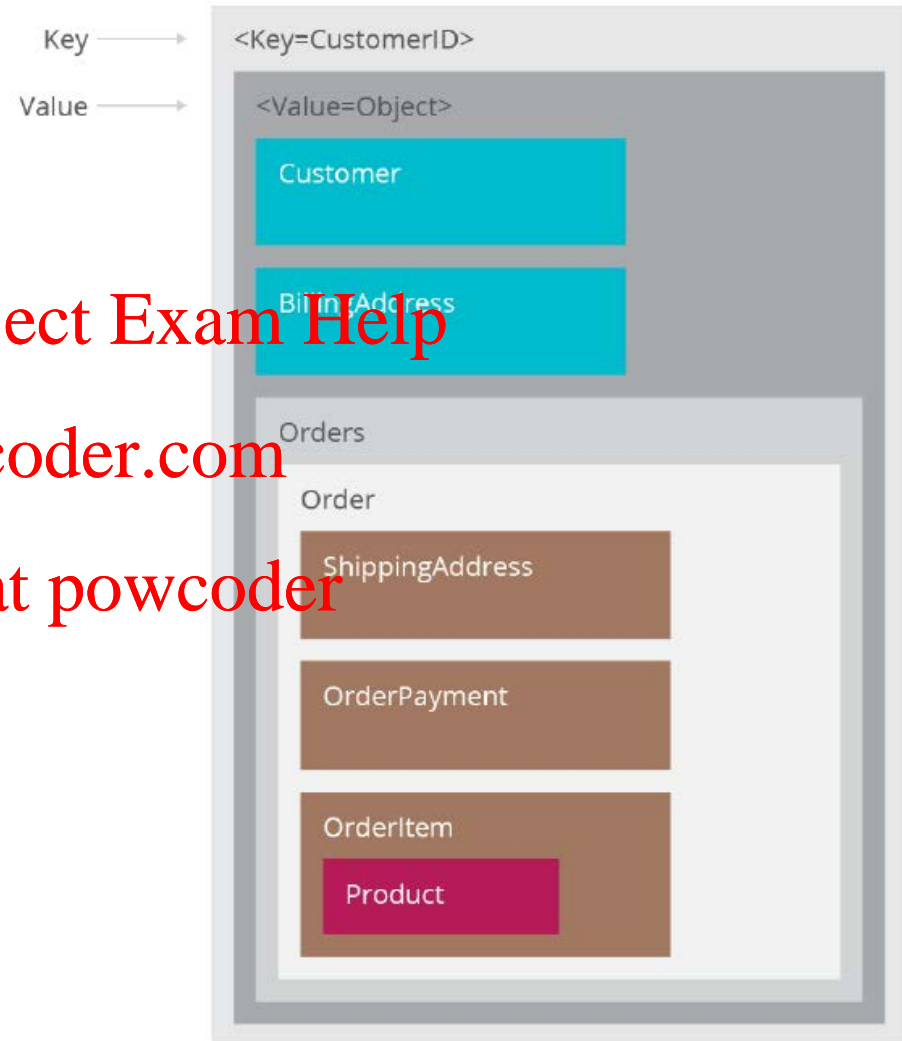
Types of NoSQL databases

WIKICOURT



(diagram from Martin Fowler)

- Key = primary key
- Value = anything (number, array, image, JSON) –*the application is in charge of interpreting what it means*
- Operations: *Put* (for storing), *Get* and *Update*
- **Examples:** Riak, Redis, Memcached, BerkeleyDB, HamsterDB, Amazon DynamoDB, Project Voldemort, Couchbase



- Similar to a key-value store except that the document is "examinable" by the databases, so its content can be queried, and parts of it updated
- Document = JSON file
- **Examples:** MongoDB, CouchDB, Terrastore, OrientDB, RavenDB

<Key=CustomerID>

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



- MongoDB documents are composed of field-and-value pairs

```
{  
  field1: value1,  
  field2: value2,  
  field3: value3,  
  ...  
  fieldN: valueN  
}
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

```
var mydoc = {  
  _id: ObjectId("5099803df3f4948bd2f98391"),  
  name: { first: "Alan", last: "Turing" },  
  birth: new Date('Jun 23, 1912'),  
  death: new Date('Jun 07, 1954'),  
  contribs: [ "Turing machine", "Turing test", "Turingery" ],  
  views : NumberLong(1250000)  
}
```


start the mongodb server, then start the mongo shell with "mongo"
show dbs// show a list of all databases
use test// use the database called 'test'
show **collections**// show all collections in the database 'test'
db.students.**insert**({name: "Jack", born: 1992})// add a doc to collection
db.students.insert({name: "Jill", born: 1990})// add a doc to collection
db.students.**find**()// list all docs in students
db.students.find({name: "Jill"})// list all docs where name field = 'Jill'
db.students.**update**({name: "Jack"}, { \$set: {born: 1990}})// change Jack's year
db.students.**remove**({born: 1990}) // delete docs where year = 1990

// now insert complex documents from file - note repeating group, no schema
db.students.find().forEach(printjson)// print all docs in neat JSON format
db.students.find({born:1990}, {name: true})// print names for all born in 1990
db.students.update({id:222222}, { \$addToSet:
{subjects: {subject: "English", result: "H1"}}})
db.students.find({id:222222}, { _id:false, subjects:true}).forEach(printjson)

db.students.insert({name: "John", color: "blue"})
// add a new student – different schema but still works



Types of NoSQL: column families

- Columns rather than rows are stored together on disk.
- Makes analysis faster, as less data is fetched.
- This is like automatic vertical partitioning.
- Related columns grouped together into 'families'.
- **Examples:** Cassandra, BigTable, HBase

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



<https://www.youtube.com/watch?v=8KGVFB3kVHQ>

- *Key-value, document store and column-family* are “aggregate-oriented- store business object in its entirety” databases (in Fowler’s terminology)

Assignment Project Exam Help

- Pros:
 - *entire* aggregate of data is stored together (no need for transactions)
 - efficient storage on clusters / distributed databases
- Cons:
 - hard to analyse across subfields of aggregates
 - e.g. sum over products instead of orders



- A 'graph' is a node-and-arc network
- Social graphs (e.g. friendship graphs) are common examples
- Graphs are difficult to program in relational DB
- A *graph DB* stores entities and their relationships
- Graph queries deduce knowledge from the graph

Assignment Project Exam Help

<https://powcoder.com>

• Examples:

Neo4J

Infinite Graph

OrientDBv

FlockDB

TAO

Add WeChat powcoder



Table 2-1. Finding extended friends in a relational database versus efficient finding in Neo4j

Depth	RDBMS execution time(s)	Neo4j execution time(s)	Records returned
2	0.016	0.01	~2500
3	30.267	0.168	~110,000
4	1543.505	1.359	~600,000
5	Unfinished	2.132	~800,000



- Key-value stores
 - A simple pair of a key and an associated collection of values. Key is usually a string. The database has no knowledge of the structure or meaning of the values.
- Document Stores
 - Like a key-value store, but “document” goes further than “value”. The documents are structured, so specific elements can be manipulated separately.
- Column-family stores
 - Data is grouped in “column groups/families” for efficiency reasons.
- Graph-oriented databases
 - Maintain information regarding the relationships between data items. Nodes with properties.

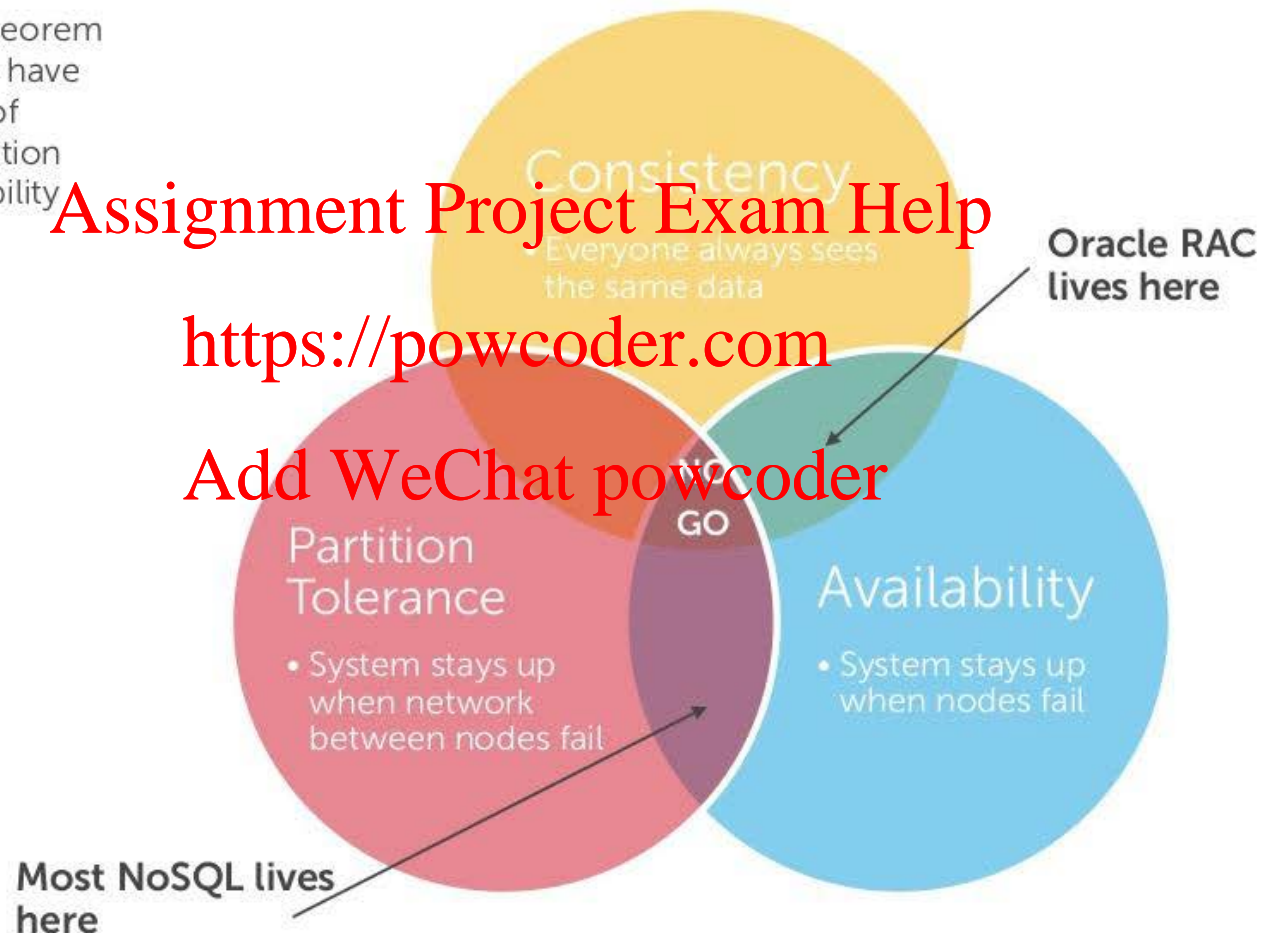
CAP Theorem says something has to give

- CAP (Brewer's) Theorem says you can only have two out of three of Consistency, Partition Tolerance, Availability

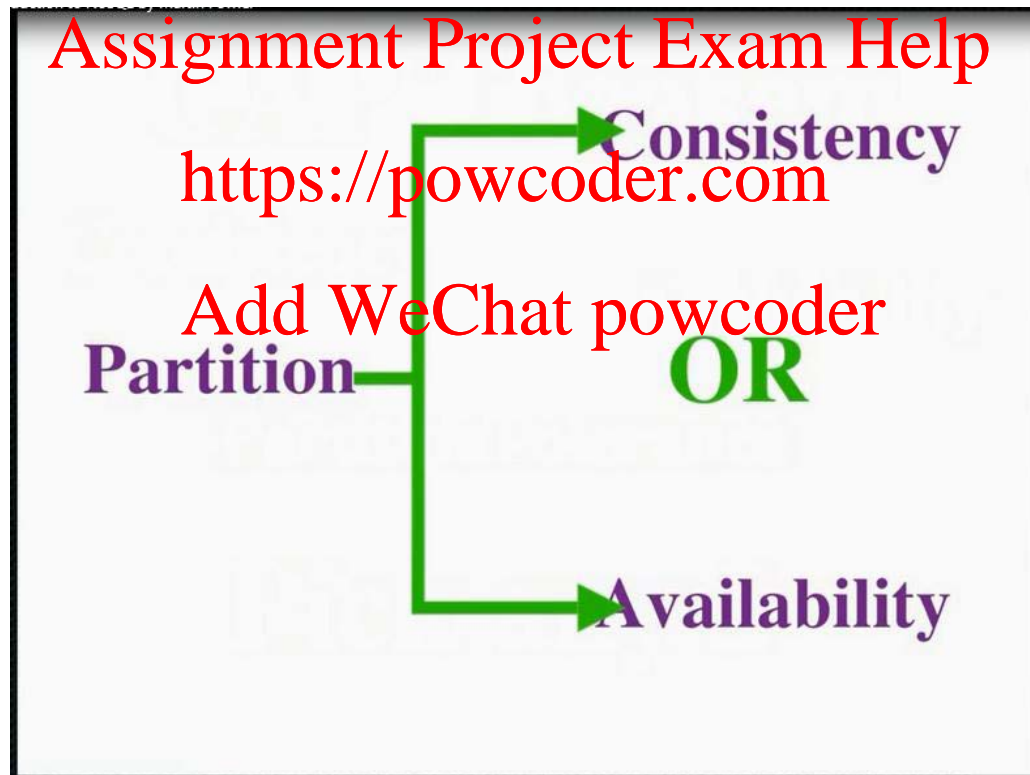
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



- Fowler's version of CAP theorem: *If you have a distributed database, when a partition occurs, you must then choose consistency OR availability.*



ACID (Atomic, Consistent, Isolated, Durable)

vs

Base (Basically Available, Soft State, Eventual Consistency)

- **Basically Available:** This constraint states that the system does guarantee the availability of the data; there will be a response to any request. But data may be in an *inconsistent* or *changing* state.

<https://powcoder.com>

- **Soft state:** The state of the system could change over time -even during times without input there may be changes going on due to 'eventual consistency'.
- **Eventual consistency:** The system will eventually become consistent once it stops receiving input. The data will propagate to everywhere it needs to, sooner or later, but the system will continue to receive input and is not checking the consistency of every transaction before it moves onto the next one.

More technical details (I won't ask you these things):

Assignment Project Exam Help

<https://powcoder.com>

https://www.youtube.com/watch?v=YUWUH_7aWHs&index=11&list=PLdQddgMBv5zHcEN9RrhADc3VBCOhY2h

Add WeChat powcoder



What is examinable?

- What is big data/NoSQL?
- What are the characteristics of NoSQL databases
- Types of NoSQL databases
- CAP theorem/BASE

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



- Databases of the future (non-examinable research avenues)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder