# NoSQL
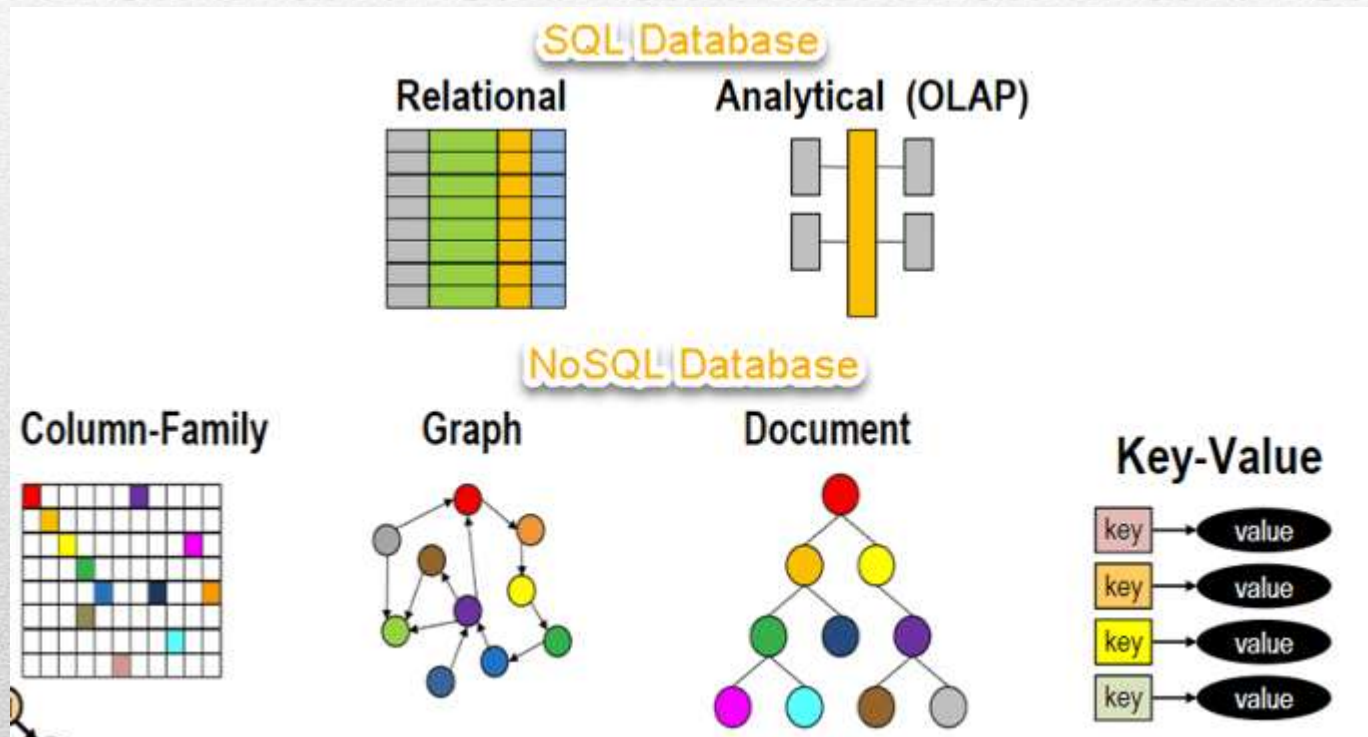
# Contents

* What is NoSQL , why NoSQL
* Understand NoSQL business Drivers
* Desirable features of NoSQL
* Need for NoSQL through case studies
* NoSQL data Architectural Pattern
* variations of NoSQL Architectural Pattern
* How NoSQL is used to manage big data
* How NoSQL system handles big data problems

- NoSQL Database is a non-relational Data Management System, that does not require a fixed schema. It avoids joins, and is easy to scale.

- NoSQL is used for Big data and real-time web apps.

- For example, companies like Twitter, Facebook and Google collect terabytes of user data every single day.

NoSQL database system encompasses a wide range of database technologies that can store structured, semi-structured, unstructured and polymorphic data.

# Brief History of NoSQL Databases

1998- Carlo Strozzi use the term NoSQL for his lightweight, open-source relational

   database

2000- Graph database Neo4j is launched

2004- Google BigTable is launched

2005- CouchDB is launched

2007- The research paper on Amazon Dynamo is released

2008- Facebooks open sources the Cassandra project

2009- The term NoSQL was reintroduced
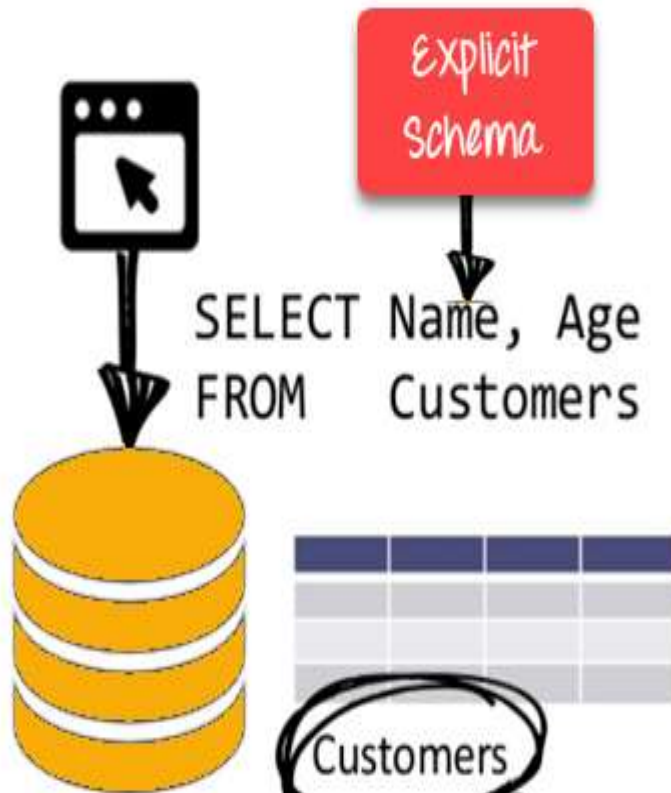
**Features of NoSQL**

**Non-relational**

- NoSQL databases never follow the relational model
- Never provide tables with flat fixed-column records
- Work with self-contained aggregates or BLOBs
- Doesn't require object-relational mapping and data normalization
- No complex features like query languages, query planners , referential integrity joins, ACID
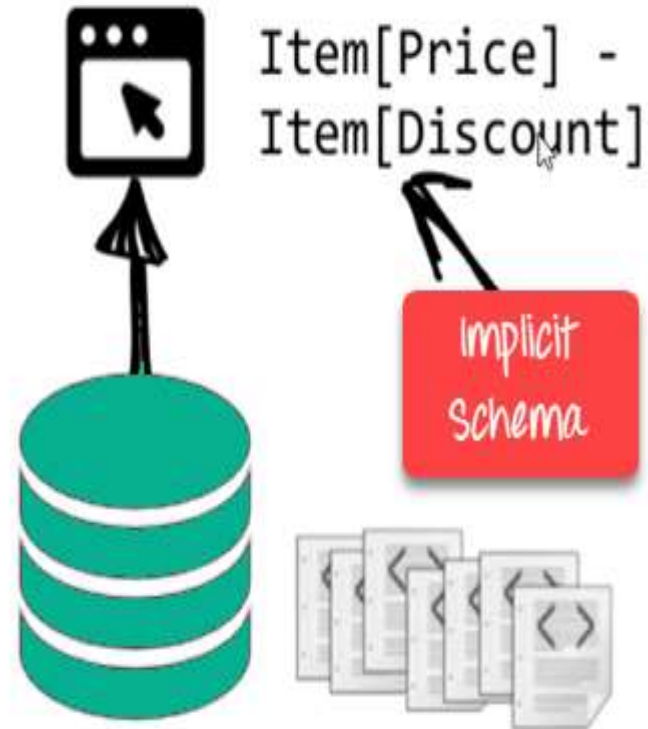
**Schema-free**

- NoSQL databases are either schema-free or have relaxed schemas
- Do not require any sort of definition of the schema of the data
- Offers heterogeneous structures of data in the same domain
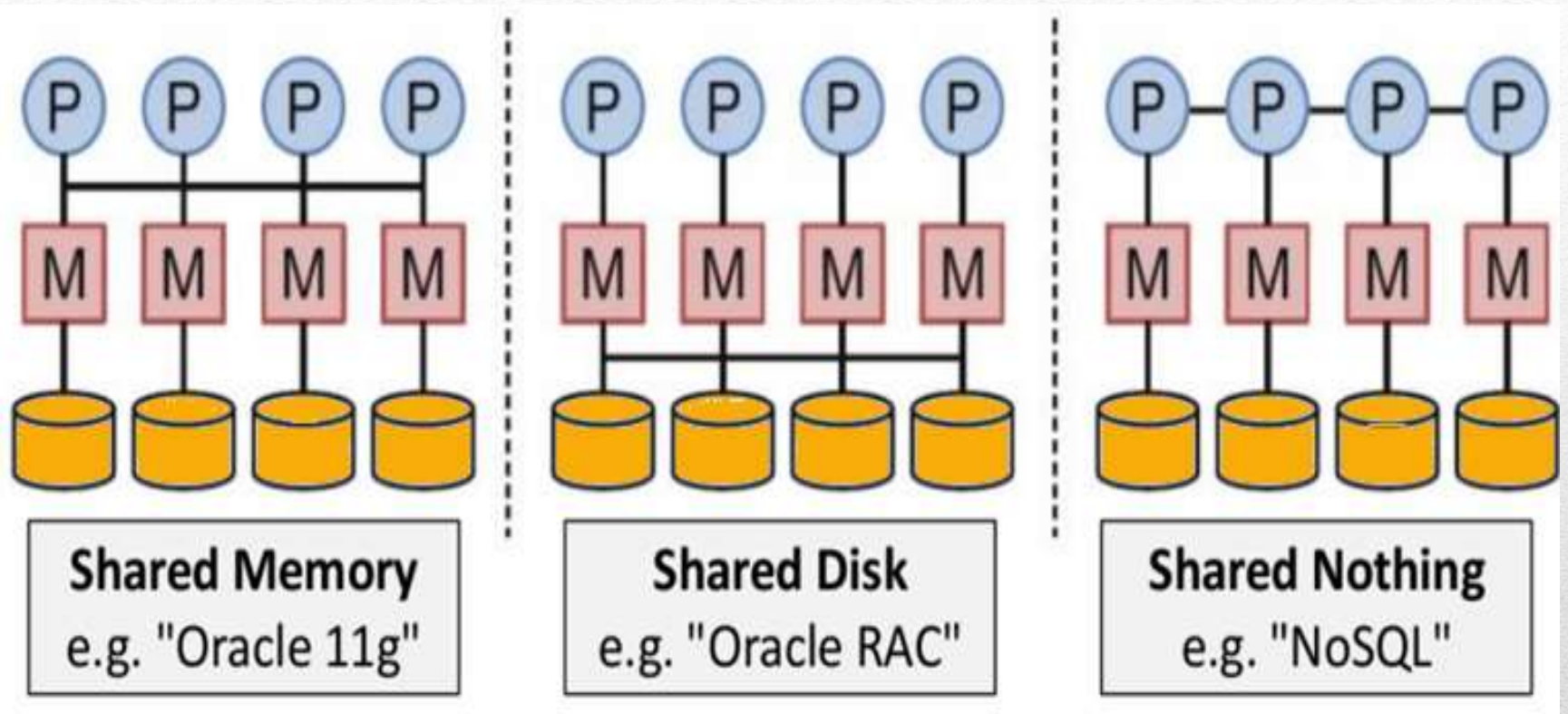
**NoSQL is Schema-Free**

## Simple API

- Offers easy to use interfaces for storage and querying data provided

- APIs allow low-level data manipulation & selection methods

- Text-based protocols mostly used with HTTP REST with JSON

- Mostly used no standard based NoSQL query language

- Web-enabled databases running as internet-facing services

## Distributed

- Multiple NoSQL databases can be executed in a distributed fashion

- Offers auto-scaling and fail-over capabilities

- Often ACID concept can be sacrificed for scalability and throughput

- Mostly no synchronous replication between distributed nodes Asynchronous Multi-Master Replication, peer-to-peer, HDFS Replication

- Only providing eventual consistency

- Shared Nothing Architecture. This enables less coordination and higher distribution.

**NoSQL is Shared Nothing.**

**Types of NoSQL Databases**

NoSQL Databases are mainly categorized into four types: Key-value pair, Column-oriented, Graph-based and Document-oriented. Every category has its unique attributes and limitations. None of the above-specified database is better to solve all the problems. Users should select the database based on their product needs.

**Types of NoSQL Databases:**

- Key-value Pair Based
- Column-oriented Graph
- Graphs based
- Document-oriented

**Key Value**

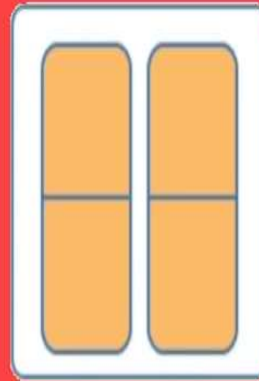Example:
Riak, Tokyo Cabinet, Redis server, Memcached, Scalaris

**Document-Based**
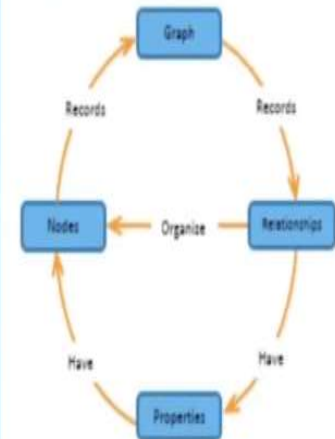
Example:
MongoDB, CouchDB, OrientDB, RavenDB

**Column-Based**

Example:
BigTable, Cassandra, Hbase, Hypertable

**Graph-Based**

Graph

Records                    Records

Nodes    Organize    Relationships

Have                          Have

Properties

Example:
Neo4J, InfoGrid, Infinite Graph, Flock DB

# Key Value Pair Based

Data is stored in key/value pairs. It is designed in such a way to handle lots of data and heavy load.

Key-value pair storage databases store data as a hash table where each key is unique, and the value can be a JSON, BLOB(Binary Large Objects), string, etc.

For example, a key-value pair may contain a key like "Website" associated with a value like "google".

| Key | Value |
|---|---|
| Name | Joe Bloggs |
| Age | 42 |
| Occupation | Stunt Double |
| Height | 175cm |
| Weight | 77kg |

It is one of the most basic NoSQL database example. This kind of NoSQL database is used as a collection, dictionaries, associative arrays, etc. Key value stores help the developer to store schema-less data. They work best for shopping cart contents.

Redis, Dynamo, Riak are some NoSQL examples of key-value store DataBases. They are all based on Amazon's Dynamo paper.

## Column-based

Column-oriented databases work on columns and are based on BigTable paper by Google. Every column is treated separately. Values of single column databases are stored contiguously.

| ColumnFamily | | | |
|---|---|---|---|
| Row Key | Column Name | | |
| | Key | Key | Key |
| | Value | Value | Value |
| | Column Name | | |
| | Key | Key | Key |
| | Value | Value | Value |

**Column based NoSQL database**

- They deliver high performance on aggregation queries like SUM, COUNT, AVG, MIN etc. as the data is readily available in a column.

- Column-based NoSQL databases are widely used to manage data warehouses, business intelligence, CRM, Library card catalogs , HBase, Cassandra, HBase, Hypertable are NoSQL query examples of column based database

## Document-Oriented:

Document-Oriented NoSQL DB stores and retrieves data as a key value pair but the value part is stored as a document. The document is stored in JSON or XML formats. The value is understood by the DB and can be queried.



**Relational Vs. Document**

In this diagram on your left you can see we have rows and columns, and in the right, we have a document database which has a similar structure to JSON. Now for the relational database, you have to know what columns you have and so on. However, for a document database, you have data store like JSON object. You do not require to define which make it flexible.
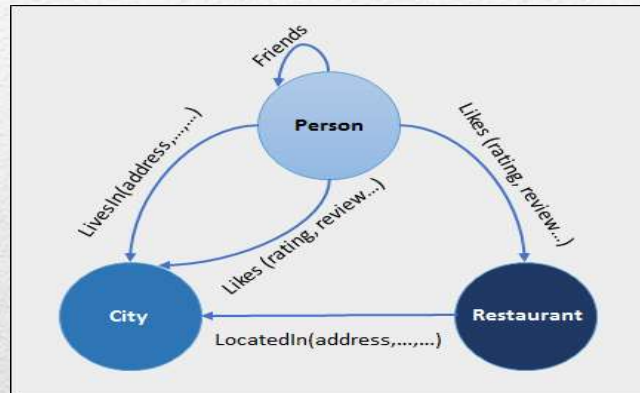
The document type is mostly used for CMS systems, blogging platforms, real-time analytics & e-commerce applications. It should not use for complex transactions which require multiple operations or queries against varying aggregate structures.

Amazon SimpleDB, CouchDB, MongoDB, Riak, Lotus Notes, MongoDB, are popular Document originated DBMS systems.

# Graph-Based

A graph type database stores entities as well the relations amongst those entities. The entity is stored as a node with the relationship as edges. An edge gives a relationship between nodes. Every node and edge has a unique identifier.



Compared to a relational database where tables are loosely connected, a Graph database is a multi-relational in nature. Traversing relationship is fast as they are already captured into the DB, and there is no need to calculate them.

Graph base database mostly used for social networks, logistics, spatial data.

Neo4J, Infinite Graph, OrientDB, FlockDB are some popular graph-based databases.

**Query Mechanism tools for NoSQL**

The most common data retrieval mechanism is the REST-based retrieval of a value based on its key/ID with GET resource

Document store Database offers more difficult queries as they understand the value in a key-value pair. For example, CouchDB allows defining views with MapReduce

**What is the CAP Theorem?**

CAP theorem is also called brewer's theorem. It states that is impossible for a distributed data store to offer more than two out of three guarantees

- Consistency
- Availability
- Partition Tolerance

**Consistency:**

The data should remain consistent even after the execution of an operation. This means once data is written, any future read request should contain that data. For example, after updating the order status, all the clients should be able to see the same data.

**Availability:**

The database should always be available and responsive. It should not have any downtime.

**Partition Tolerance:**

Partition Tolerance means that the system should continue to function even if the communication among the servers is not stable. For example, the servers can be partitioned into multiple groups which may not communicate with each other. Here, if part of the database is unavailable, other parts are always unaffected.
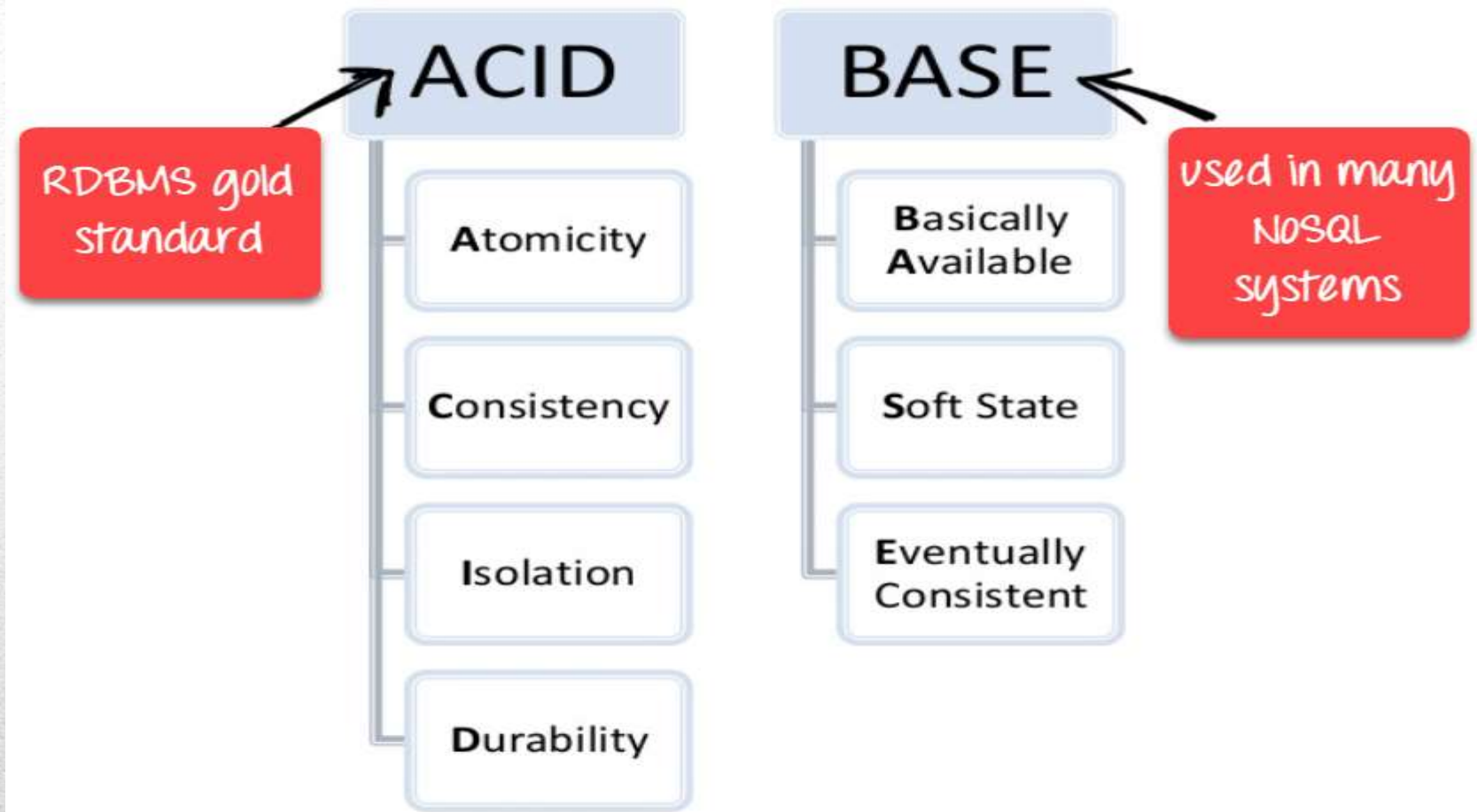
**Eventual Consistency**

The term "eventual consistency" means to have copies of data on multiple machines to get high availability and scalability. Thus, changes made to any data item on one machine has to be propagated to other replicas.

Data replication may not be instantaneous as some copies will be updated immediately while others in due course of time. These copies may be mutually, but in due course of time, they become consistent. Hence, the name eventual consistency.

BASE: **B**asically **A**vailable, **S**oft state, **E**ventual consistency

- Basically, available means DB is available all the time as per CAP theorem
- Soft state means even without an input; the system state may change
- Eventual consistency means that the system will become consistent over time

**ACID**

RDBMS gold standard

- **A**tomicity
- **C**onsistency
- **I**solation
- **D**urability

**BASE**

used in many NOSQL systems

- **B**asically **A**vailable
- **S**oft State
- **E**ventually Consistent

**Advantages of NoSQL**

- Can be used as Primary or Analytic Data Source

- Big Data Capability

- No Single Point of Failure

- Easy Replication

- No Need for Separate Caching Layer

- It provides fast performance and horizontal scalability.

- Can handle structured, semi-structured, and unstructured data with equal effect

- Object-oriented programming which is easy to use and flexible

- NoSQL databases don't need a dedicated high-performance server

- Support Key Developer Languages and Platforms

- Simple to implement than using RDBMS

- It can serve as the primary data source for online applications.

- Handles big data which manages data velocity, variety, volume, and complexity

- Excels at distributed database and multi-data center operations

- Eliminates the need for a specific caching layer to store data

- Offers a flexible schema design which can easily be altered without downtime or service

## Disadvantages of NoSQL

- No standardization rules

- Limited query capabilities

- RDBMS databases and tools are comparatively mature

- It does not offer any traditional database capabilities, like consistency when multiple transactions are performed simultaneously.

- When the volume of data increases it is difficult to maintain unique values as keys become difficult

- Doesn't work as well with relational data

- The learning curve is stiff for new developers

- Open source options so not so popular for enterprises.

# CAP Theorem

- Conjectured by Prof. Eric Brewer at PODC (Principle of Distributed Computing) 2000 keynote talk

- Described the *trade-offs involved in distributed system*

- It is impossible for a web service to provide following *three guarantees at the same time*:
  - **Consistency**
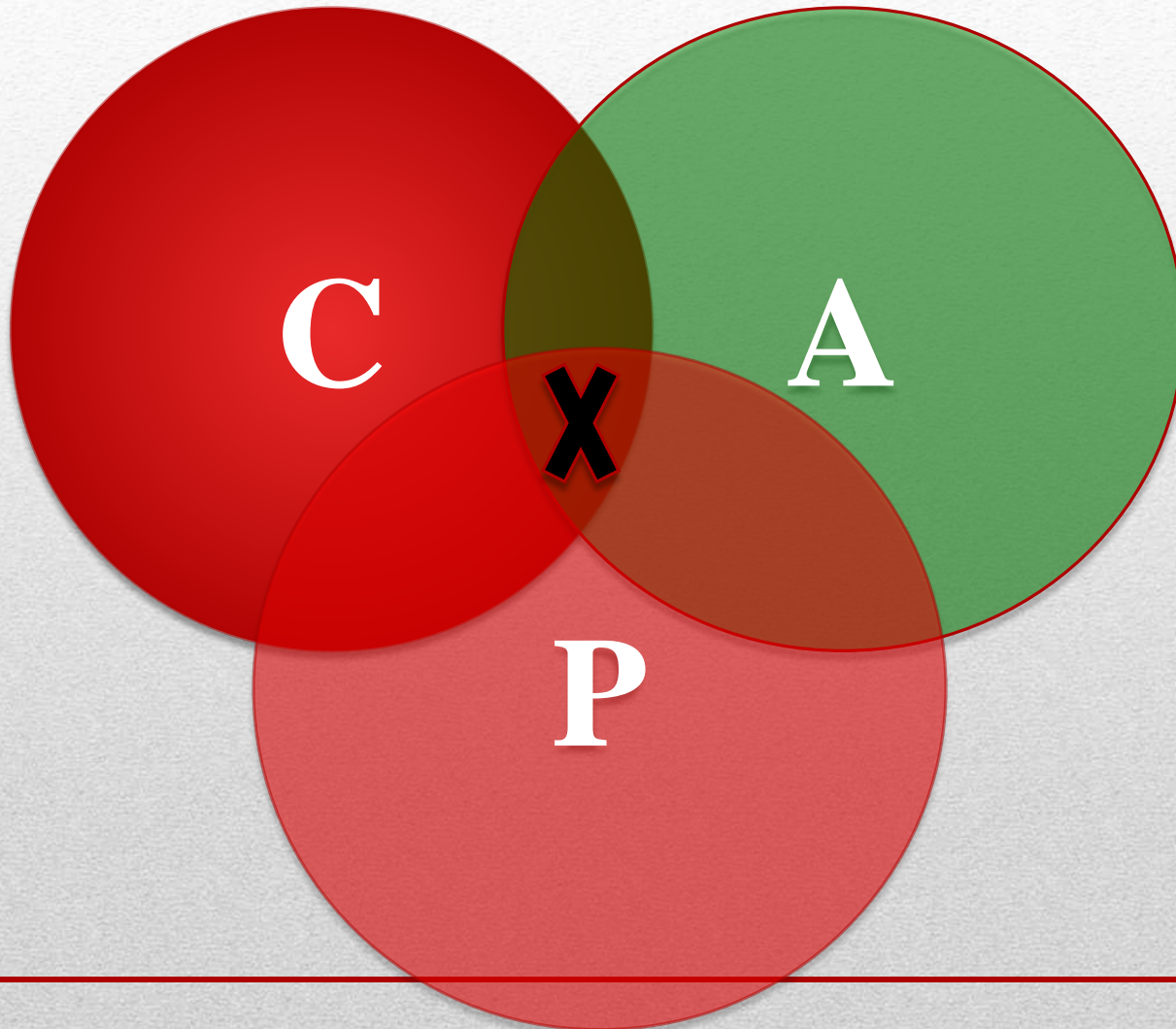  - **Availability**
  - **Partition-tolerance**

- **C**onsistency:
  - All nodes should see the same data at the same time
- **A**vailability:
  - Node failures do not prevent survivors from continuing to operate
- **P**artition-tolerance:
  - The system continues to operate despite network partitions
- A distributed system can satisfy any two of these guarantees at the same time **but not all three**
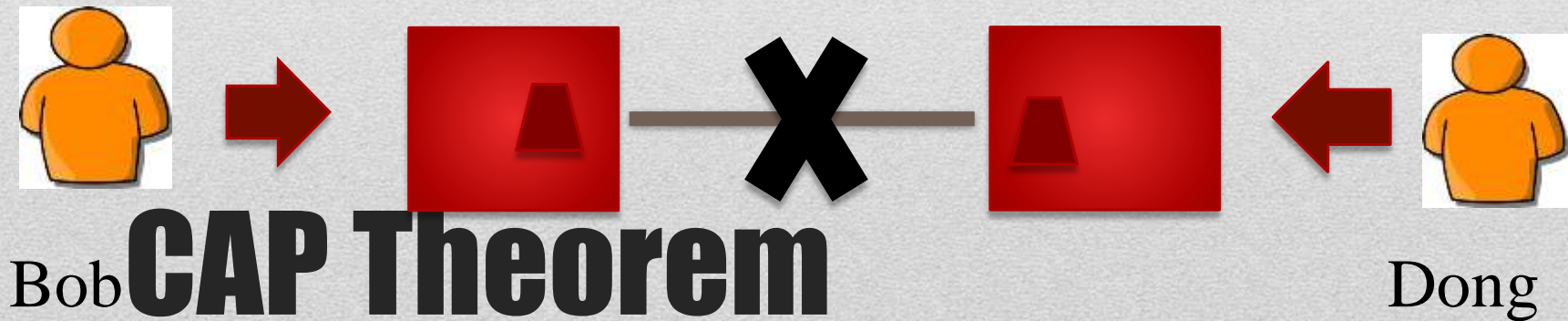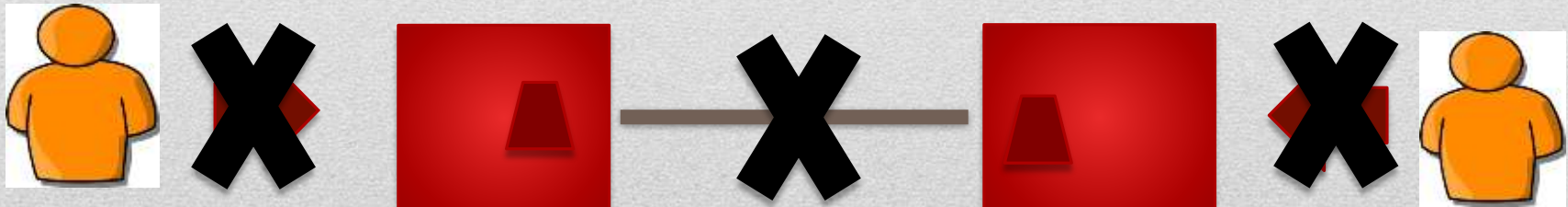
# CAP Theorem

# CAP Theorem

**Hotel Booking**: are we double-booking the same room?

- A simple example:



Bob **CAP Theorem**                                    Dong

# **Hotel Booking**: are we double-booking the same room?

- A simple example:



Bob **CAP Theorem**                              Dong

**Hotel Booking**: are we double-booking the same room?

- A ~~APPROVED~~ ample:

Bob **CAP Theorem** Dong

- 2002: Proven by research conducted by Nancy Lynch and Seth Gilbert at MIT

Gilbert, Seth, and Nancy Lynch. "Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services." ACM SIGACT News 33.2 (2002): 51-59.
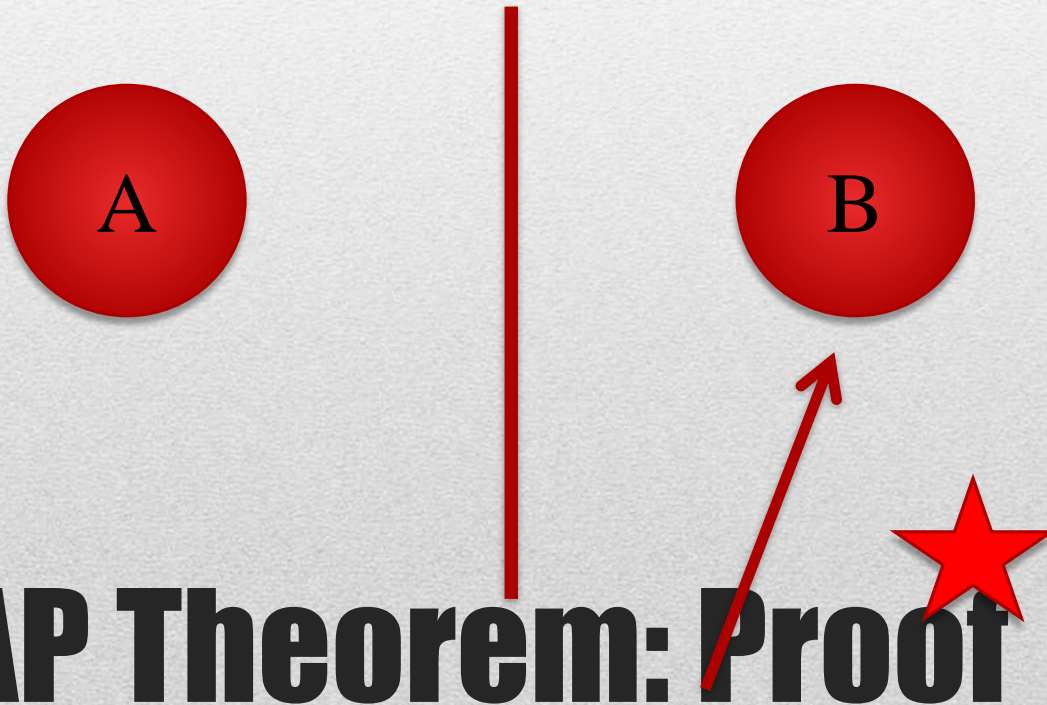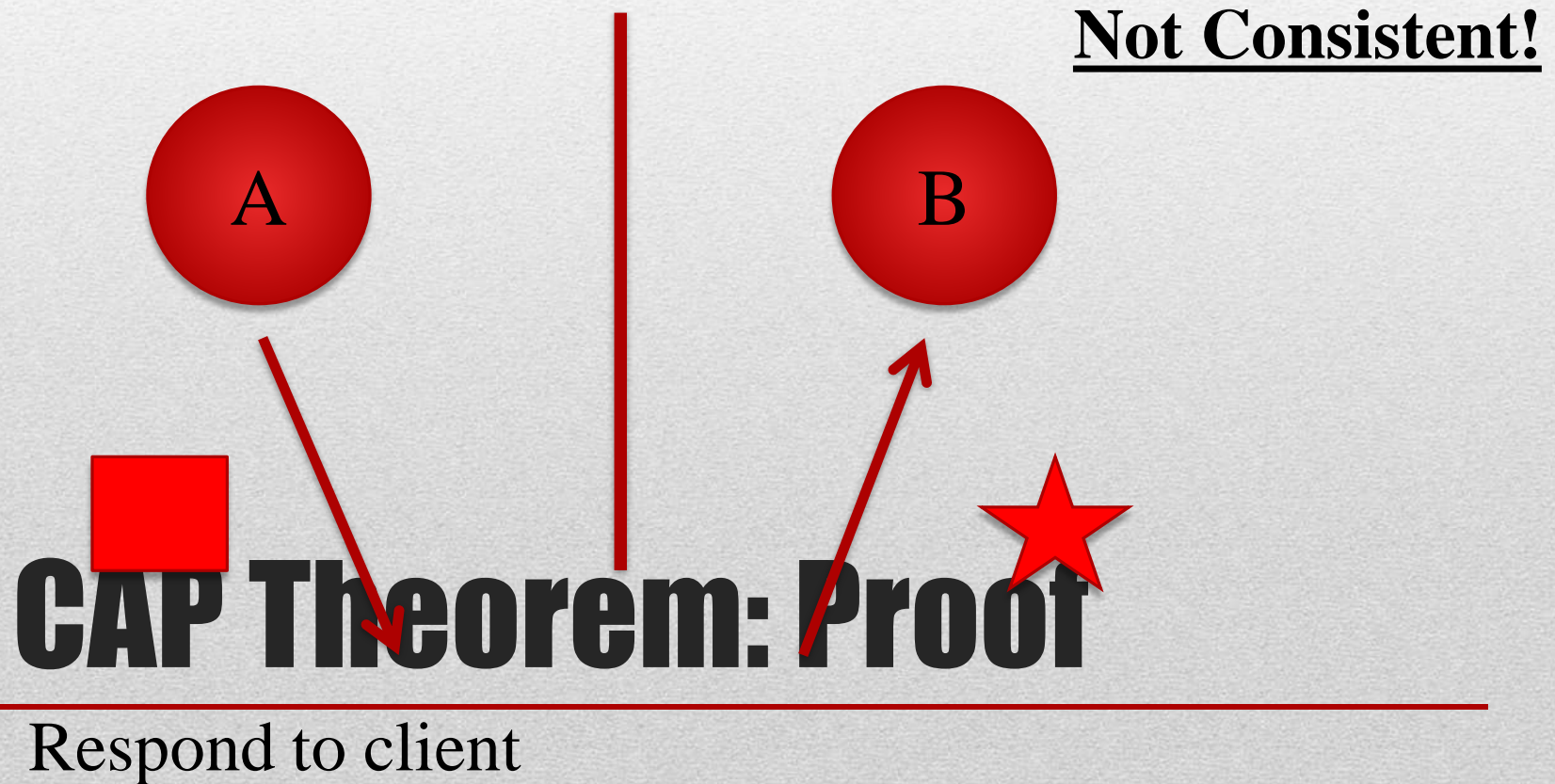
# CAP Theorem: Proo

- A simple proof using two nodes:



CAP Theorem: Proof

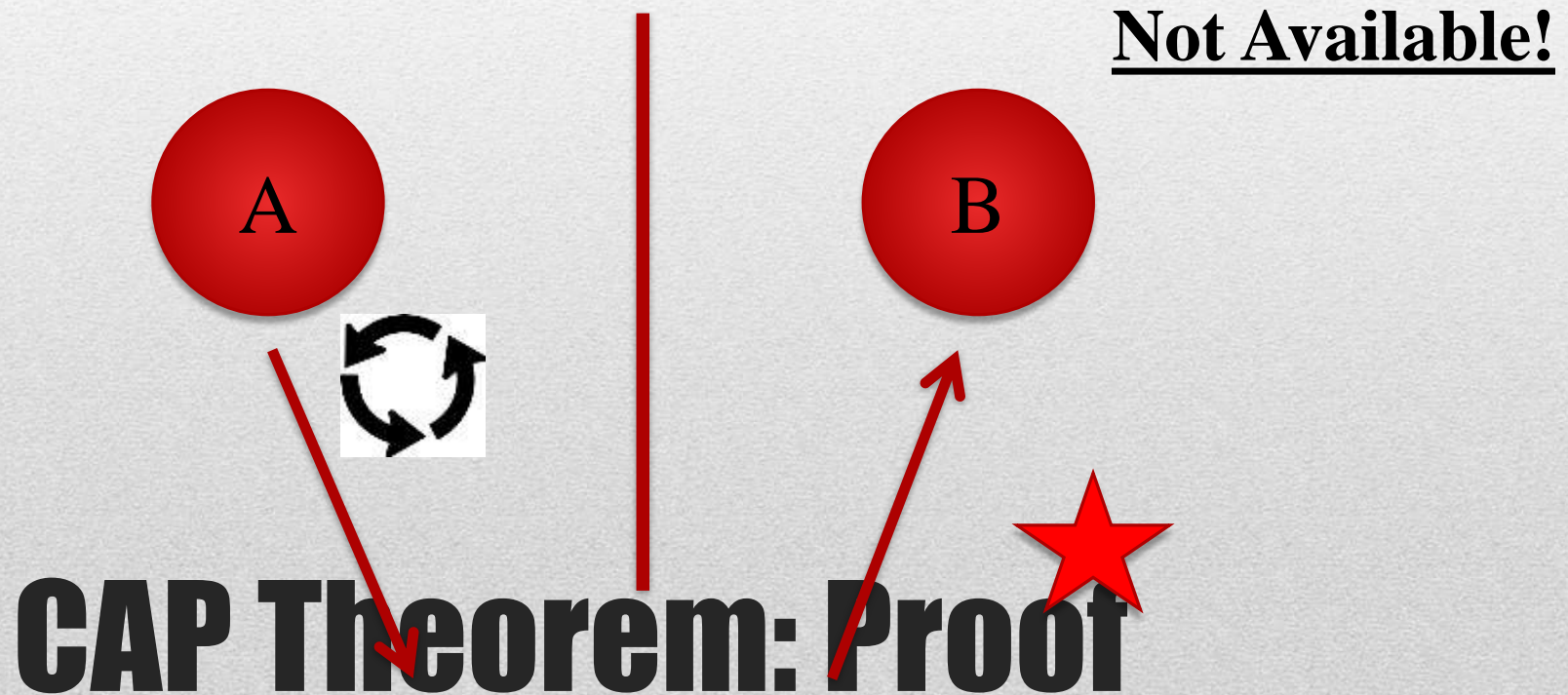- A simple proof using two nodes:

**Not Consistent!**

A

B

# CAP Theorem: Proof

Respond to client

- A simple proof using two nodes:

**Not Available!**

A

B

# CAP Theorem: Proof

Wait to be updated

- A simple proof using two nodes:



**Not Partition Tolerant!**

A

B

# CAP Theorem: Proof

A gets updated from B

- The future of databases is **distributed** (Big Data Trend, etc.)
- CAP theorem describes the **trade-offs** involved in distributed systems
- A proper understanding of CAP theorem is essential to **making decisions** about the future of distributed database **design**
- Misunderstanding can lead to **erroneous or inappropriate** design choices
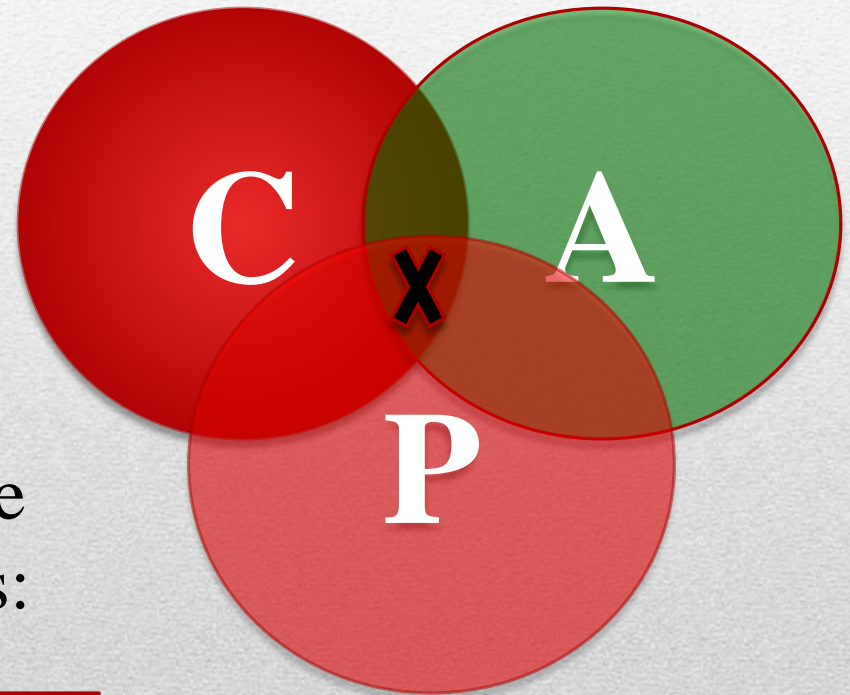
# Why this is important?

# Problem for Relational Database to Scale

- The Relational Database is built on the principle of **ACID** (Atomicity, Consistency, Isolation, Durability)

- It implies that a truly distributed relational database should have **availability, consistency and partition tolerance**.

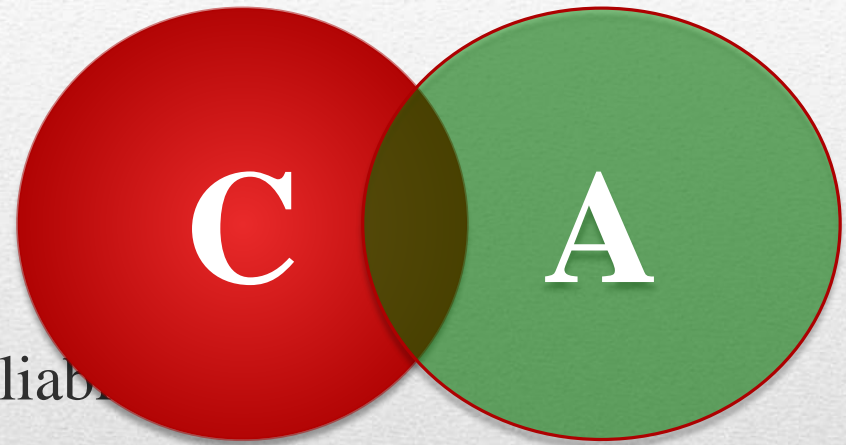- Which unfortunately is **impossible** …

# Revisit CAP Theorem

- Of the following three guarantees potentially offered a by distributed systems:
    - Consistency
    - Availability
    - Partition tolerance

- Pick two

- This suggests there are three kinds of distributed systems:
    - CP
    - AP
    - CA

**C**

**A**

**X**

**P**

*Any problems?*

- How about CA?
- Can a distributed system (with unreliable network) really be not tolerant of partitions?

# A popular misconception: 2 out 3

- Coda Hale, Yammer software engineer:
  - "Of the CAP theorem's Consistency, Availability, and Partition Tolerance, **Partition Tolerance is mandatory in distributed systems**. You cannot not choose it."

# A few witnesses

- Werner Vogels, Amazon CTO
  - "An important observation is that in larger distributed-scale systems, network partitions are a given; therefore, **consistency and availability cannot be achieved at the same time**."

# A few witnesses

- Daneil Abadi, Co-founder of Hadapt
  - So in reality, there are only two types of systems ... I.e., if there is a partition, **does the system give up availability or consistency?**

# A few witnesses

- Prof. Eric Brewer: father of CAP theorem
  - "The "2 of 3" formulation was always **misleading** because it tended to oversimplify the tensions among properties. ...
  - **CAP prohibits only a tiny part of the design space**: *perfect availability and consistency in the presence of partitions*, which are rare."

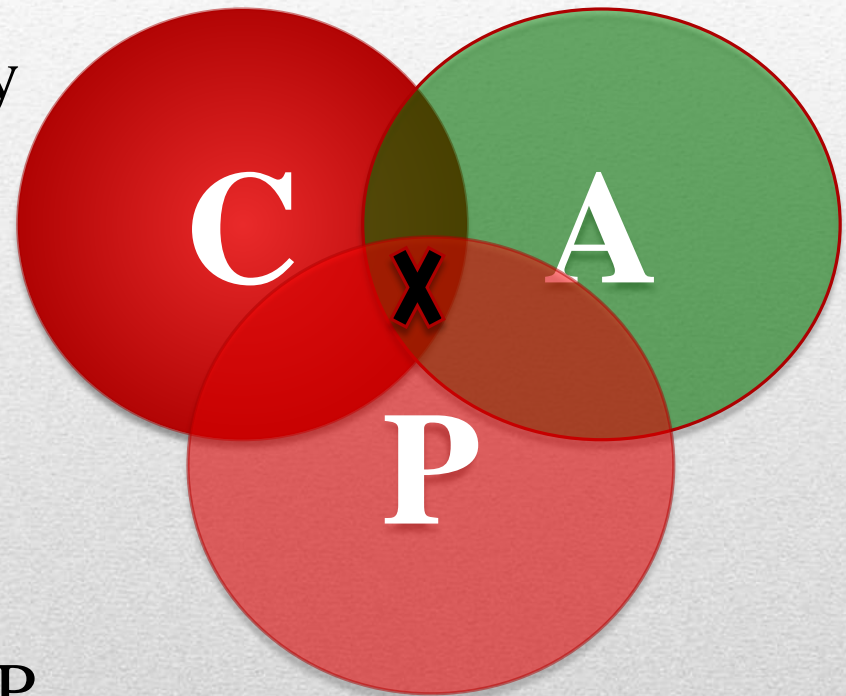# CAP Theorem 12 year later

# Consistency or Availability

- Consistency and Availability is not "binary" decision

- AP systems relax consistency in favor of availability – but are not inconsistent

- CP systems sacrifice availability for consistency- but are not unavailable

- This suggests both AP and CP systems can offer a degree of consistency, and availability, as well as partition tolerance

- Example:
  - Web Caching
  - DNS
- Trait:
  - Optimistic
  - Expiration/Time-to-live
  - Conflict resolution

# AP: Best Effort Consistency

- Example:
  - Majority protocols
  - Distributed Locking (Google Chubby Lock service)
- Trait:
  - Pessimistic locking
  - Make minority partition unavailable

# CP: Best Effort Availability

- Strong Consistency
  - After the update completes, **any subsequent access** will return the **same** updated value.
- Weak Consistency
  - It is **not guaranteed** that subsequent accesses will return the updated value.
- **Eventual Consistency**
  - Specific form of weak consistency
  - It is guaranteed that if **no new updates** are made to object, eventually all accesses will return the last updated value (e.g., *propagate updates to replicas in a lazy fashion)*

# Types of Consistency

# Eventual Consistency Variations

- Causal consistency
  - Processes that have causal relationship will see consistent data
- Read-your-write consistency
  - A process always accesses the data item after it's update operation and never sees an older value
- Session consistency
  - As long as session exists, system guarantees read-your-write consistency
  - Guarantees do not overlap sessions

- Monotonic read consistency
  - If a process has seen a particular value of data item, any subsequent processes will never return any previous values
- Monotonic write consistency
  - The system guarantees to serialize the writes by the *same* process
- In practice

# Eventual Consistency Variations

  - A number of these properties can be combined
  - Monotonic reads and read-your-writes are most desirable

- Bob finds an interesting story and shares with Alice by posting on her Facebook wall

- Bob asks Alice to check it out

- Alice logs in her account, checks her Facebook wall but finds:

    - **Nothing is there!**

# Eventual Consistency

- Bob tells Alice to wait a bit and check out later
- Alice waits for a minute or so and checks back:

   - **She finds the story Bob shared with her!**

# Eventual Consisten

- Reason: it is possible because Facebook uses an **eventual consistent model**
- Why Facebook chooses eventual consistent model over the strong consistent one?
  - Facebook has more than 1 billion active users
  - It is non-trivial to efficiently and reliably store the huge amount of data generated at any given time
  - Eventual consistent model offers the option to **reduce the load and improve availability**

# Eventual Consistency - A Facebook Example

- Dropbox enabled immediate consistency via synchronization in many cases.
- However, what happens in case of a network partition?

- Let's do a simple experiment here:
  - Open a file in your drop box
  - Disable your network connection (e.g., WiFi, 4G)
  - Try to edit the file in the drop box: can you do that?
  - Re-enable your network connection: what happens to your dropbox folder?

# Eventual Consistency - A Dropbox Example

- Dropbox embraces eventual consistency:
  - Immediate consistency is impossible in case of a network partition
  - Users will feel bad if their word documents freeze each time they hit Ctrl+S , simply due to the large latency to update all devices across WAN
  - Dropbox is oriented to **personal syncing**, not on collaboration, so it is not a real limitation.

# Eventual Consistency - A Dropbox Example

# Eventual Consistency - An ATM Example

- In design of automated teller machine (ATM):
  - Strong consistency appear to be a nature choice
  - However, in practice, **A beats C**
  - Higher availability means **higher revenue**
  - ATM will allow you to withdraw money *even if the machine is partitioned from the network*
  - However, it puts **a limit** on the amount of withdraw (e.g., $200)
  - The bank might also charge you a fee when a overdraft happens

- An airline reservation system:
  - When most of seats are available: it is ok to rely on somewhat out-of-date data, availability is more critical
  - When the plane is close to be filled: it needs more accurate data to ensure the plane is not overbooked, consistency is more critical
- Neither strong consistency nor guaranteed availability, but it may significantly increase the tolerance of network disruption

# Dynamic Tradeoff between C and A

- No single uniform requirement
  - Some aspects require strong consistency
  - Others require high availability
- Segment the system into different components
  - Each provides different types of guarantees
- Overall guarantees neither consistency nor availability
  - Each part of the service gets exactly what it needs
- Can be partitioned along different dimensions

# Heterogeneity: Segmenting C and A

- In an e-commercial system (e.g., Amazon, e-Bay, etc), what are the trade-offs between consistency and availability you can think of? What is your strategy?
- Hint -> Things you might want to consider:
  - Different types of data (e.g., shopping cart, billing, product, etc.)
  - Different types of operations (e.g., query, purchase, etc.)
  - Different types of services (e.g., distributed lock, DNS, etc.)
  - Different groups of users (e.g., users in different geographic areas, etc.)

# Discussion

- Data Partitioning
- Operational Partitioning
- Functional Partitioning
- User Partitioning
- Hierarchical Partitioning

# Partitioning Examples

Data Partitioning

- Different data may require different consistency and availability
- Example:
  - Shopping cart: high availability, responsive, can sometimes suffer anomalies
  - Product information need to be available, slight variation in inventory is sufferable
  - Checkout, billing, shipping records must be consistent

# Partitioning Examples

Operational Partitioning

- Each operation may require different balance between consistency and availability

- Example:
  - Reads: high availability; e.g.., "query"
  - Writes: high consistency, lock when writing; e.g., "purchase"

# Partitioning Examples

Functional Partitioning

- System consists of sub-services
- Different sub-services provide different balances
- Example: A comprehensive distributed system
  - Distributed lock service (e.g., Chubby) :
    - Strong consistency
  - DNS service:
    - High availability

# Partitioning Examples

User Partitioning

- Try to keep related data close together to assure better performance
- Example: Craglist
  - Might want to divide its service into several data centers, e.g., east coast and west coast
    - Users get high performance (e.g., high availability and good consistency) if they query servers closet to them
    - Poorer performance if a New York user query Craglist in San Francisco

# Partitioning Examples

Hierarchical Partitioning

- Large global service with local "extensions"
- Different location in hierarchy may use different consistency
- Example:
  - Local servers (better connected) guarantee more consistency and availability
  - Global servers has more partition and relax one of the requirement

# Partitioning Examples

- Tradeoff between **Consistency** and **Latency**:
- Caused by the **possibility of failure** in distributed systems
  - High availability -> replicate data -> consistency problem
- Basic idea:
  - Availability and latency are arguably **the same thing**: unavailable -> extreme high latency

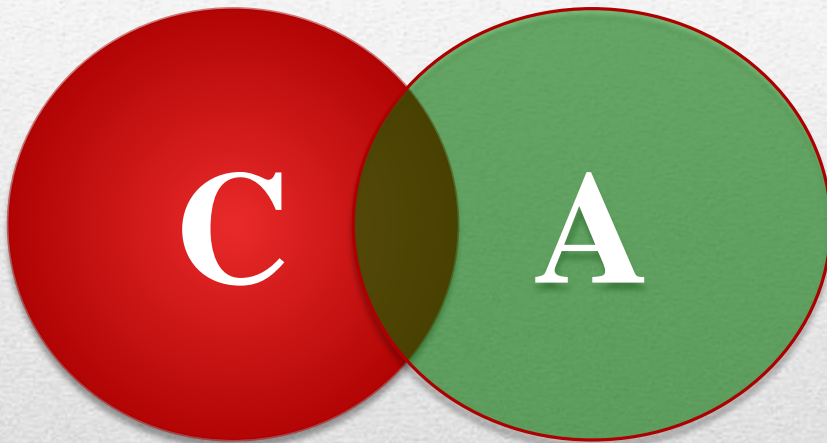# What if there are no partitions?

  - Achieving different levels of consistency/availability takes different amount of time

- A more complete description of the space of potential tradeoffs for distributed system:
  - If there is a **partition (P)**, how does the system trade off **availability and consistency (A and C)**; **else (E)**, when the system is running normally in the absence of partitions, how does the system trade off **latency (L) and consistency (C)**?

# CAP -> PACELC

Abadi, Daniel J. "Consistency tradeoffs in modern distributed database system design." Computer-IEEE Computer Magazine 45.2 (2012): 37.
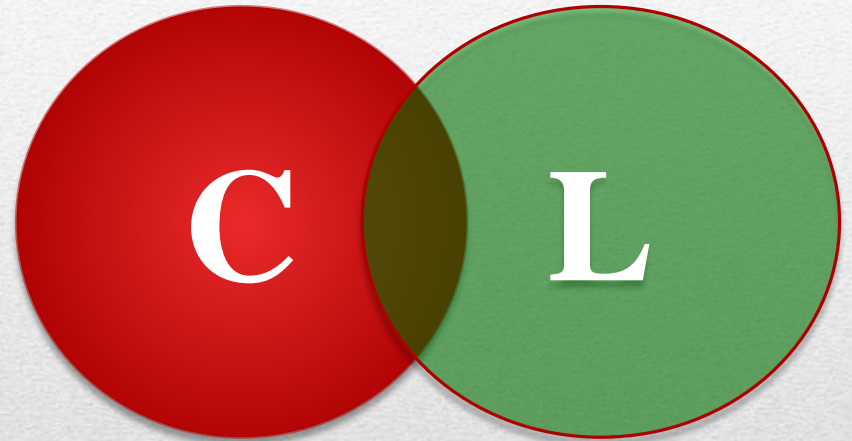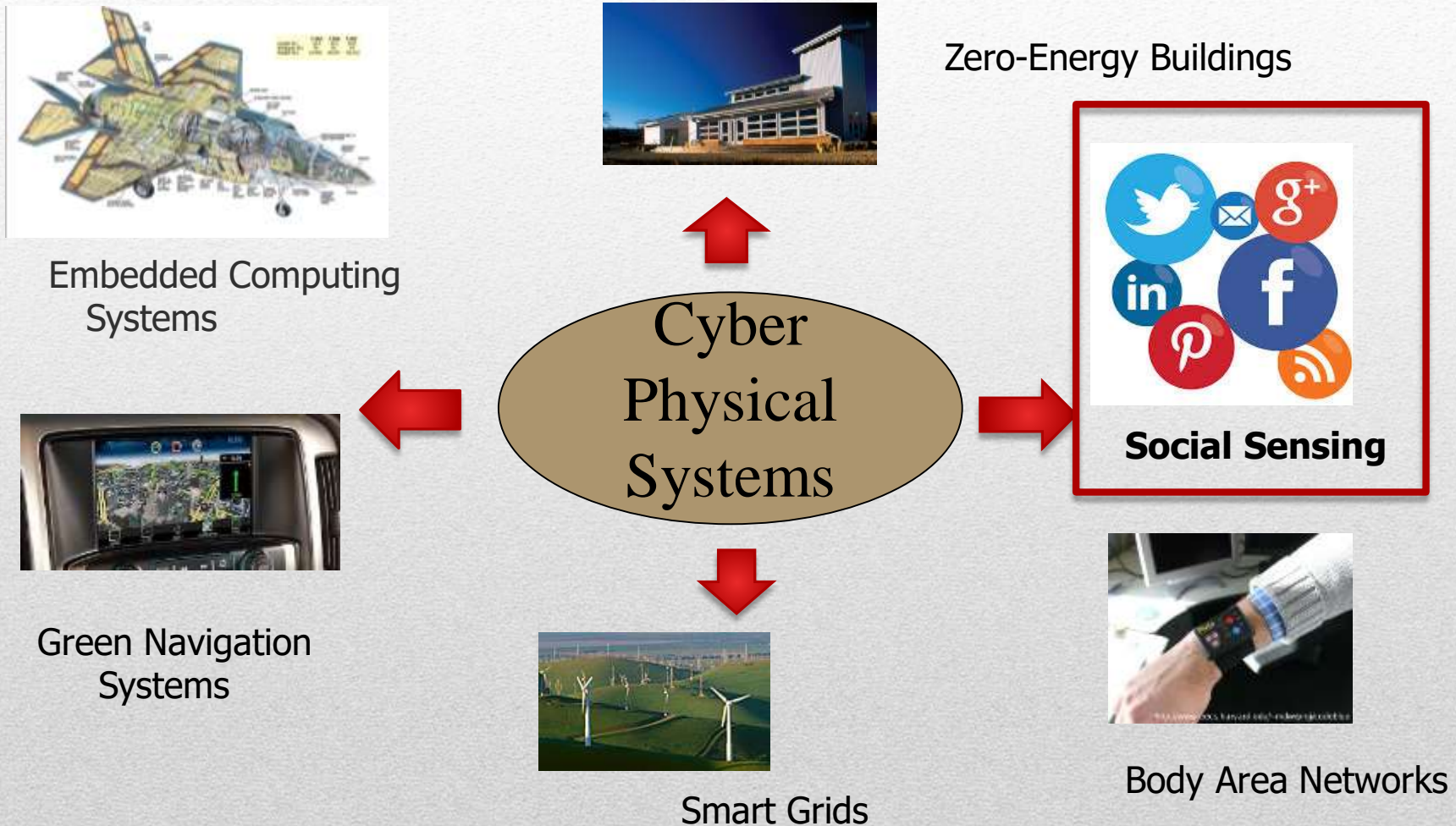
# PACELC

# Examples

- **PA/EL Systems:** Give up both Cs for availability and lower latency
  - Dynamo, Cassandra, Riak
- **PC/EC Systems:** Refuse to give up consistency and pay the cost of availability and latency
  - BigTable, Hbase, VoltDB/H-Store
- **PA/EC Systems:** Give up consistency when a partition happens and keep consistency in normal operations
  - MongoDB
- **PC/EL System:** Keep consistency if a partition occurs but gives up consistency for latency in normal operations
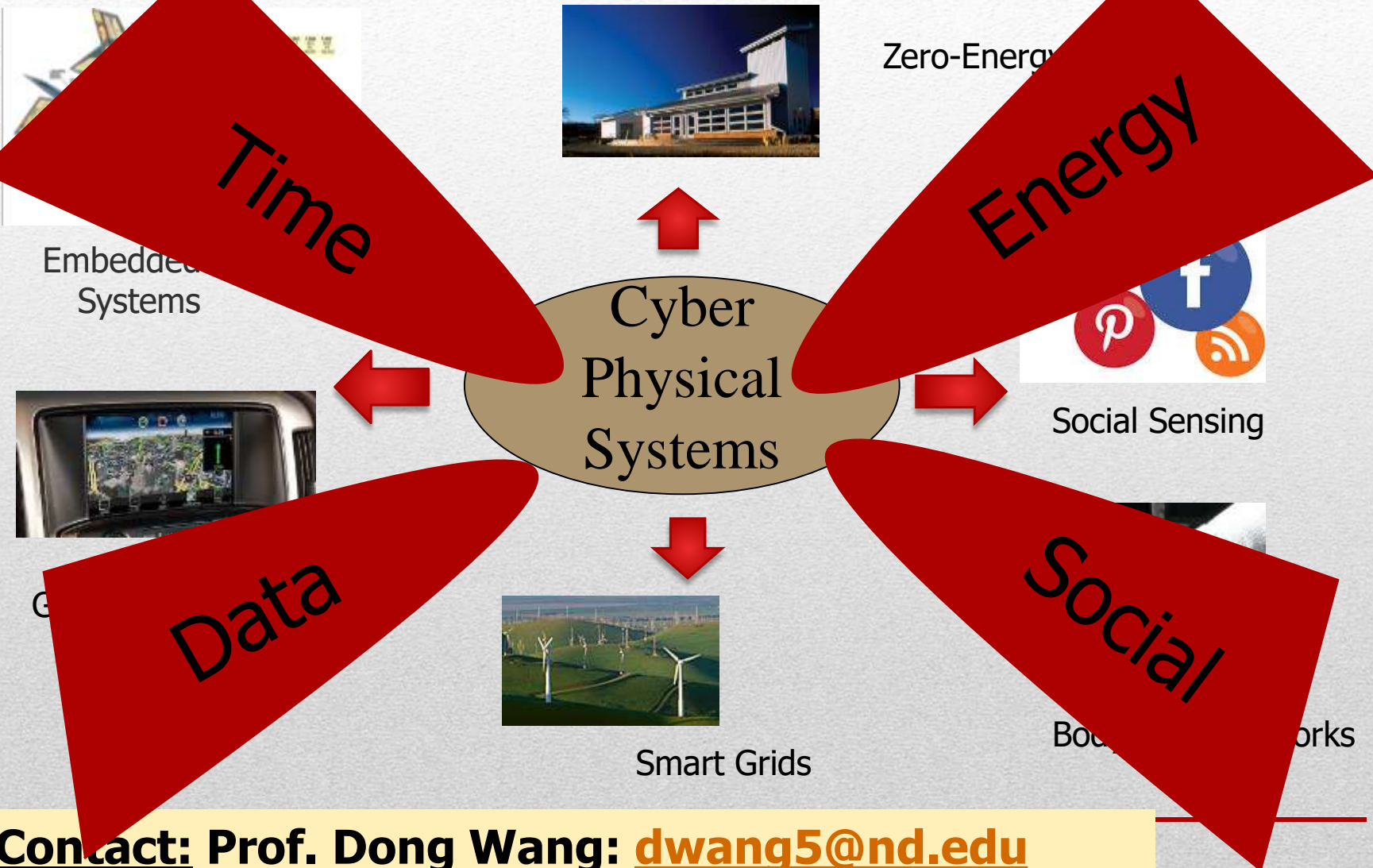  - Yahoo! PNUTS

# CSE 40437/60437, Spring 2015:
## Social Sensing & Cyber-Physical Systems



Embedded Computing Systems



Zero-Energy Buildings



**Social Sensing**

Cyber Physical Systems



Green Navigation Systems



Smart Grids



Body Area Networks

# CSE 40437/60437, Spring 2015:
## Social Sensing & Cyber-Physical Systems



Zero-Energy

Embedded Systems

Time

Energy

Cyber Physical Systems

Social Sensing

Data

Social

Smart Grids

Bod... ...orks

**Contact: Prof. Dong Wang: dwang5@nd.edu**
http://www3.nd.edu/~dwang5/teach/spring15/spring15_wang_flyer.pdf

# CSE 40437/60437, Spring 2015:
## Social Sensing & Cyber-Physical Systems

**Events**

Boston Bombing

Hurricane Sandy

Egypt unrest

News and Public Sources

twitter

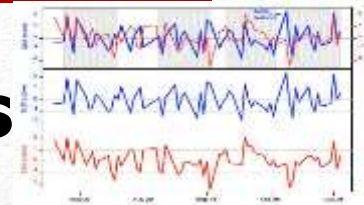facebook  YouTube

People

Sensors

**Analytics**

Data

**Decision Support**

**Applications**

**Stock Prediction (Money)**

**Traffic Monitoring (Time)**

**Disaster Response (Lives)**

**Geo Tagging (Smart City)**

**Contact: Prof. Dong Wang: dwang5@nd.edu**
http://www3.nd.edu/~dwang5/teach/spring15/spring15_wang_flyer.pdf

- Blogs website

- Multiplayer online games

- Stock trading platforms

- Video streaming sites

- Ticket booking system

- Video chat application

- Bank

1. Consistency

2. Availability

3. Partition Tolerance

# Thank you!