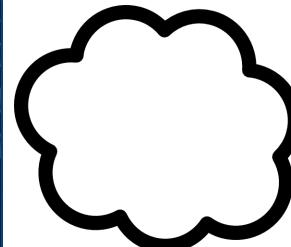
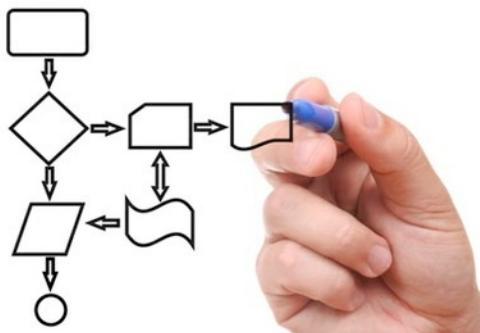


# San Jose State University

CMPE 281  
Cloud Technologies

## Introduction to NoSQL & Map Reduce



# The Papers

## MapReduce: Simplified Data Processing on Large Clusters

### MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat  
jeff@google.com, sanjay@google.com  
Google, Inc.

#### Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all values associated with the same intermediate key. Many map/reduce tasks are expressible in this model, as shown in the paper.

Programs written in this functional style are automatically parallelized and executed on a large cluster of commodity machines. The run-time system takes care of the details of partitioning the input data, scheduling the programs, reading and writing data to and from disk, managing failures, and managing the required inter-machine communication. This allows programmers without any experience with parallel and distributed systems to easily utilize large clusters.

The implementation of MapReduce runs on a large cluster of commodity machines and is highly scalable: a typical MapReduce computation processes many terabytes of data on thousands of machines. Programmers find the system easy to use; hundreds of MapReduce programs have been implemented and upwards of one thousand MapReduce jobs are executed on Google's clusters every day.

#### 1 Introduction

Over the past five years, the authors and many others at Google have implemented hundreds of special-purpose computation engines to process large amounts of data, such as crawled documents, search request logs, etc., to compute various kinds of derived data, such as inverted indices, various representations of the graph structure of web documents, summaries of the number of pages crawled per host, the set of most frequent queries in a

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across many hosts of machines in order to complete them in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures correctly to obscure the original simple computation and large amounts of complex code to deal with these issues.

As a reaction to this complexity, we designed a new abstraction that allows users to express simple computations as if they were running on a single machine, hiding the details of parallelization, fault-tolerance, data distribution and load balancing in a library. Our abstraction is inspired by the map and reduce primitives found in Lisp and mapreduce [1]. Bigtable [2] is a system that succeeds at a flexible, high-performance solution for all of these Google products. In this paper we describe the simple data model provided by Bigtable, which gives clients dynamic control over the reliability and scalability of the data out of memory or disk.

Section 2 describes the data model in more detail, and Section 3 provides an overview of the client API. Section 4 describes the distributed storage system infrastructure on which Bigtable depends. Section 5 describes the fundamentals of the Bigtable implementation, and Section 6 describes some of the refinements that we made to the original Bigtable prototype. Section 7 provides measurements of Bigtable's performance. We describe several examples of how Bigtable is used at Google and other Google products in the projects, including Google Analytics, Google Finance, Orkut, Personalized Search, Writely, and Google Earth. These products illustrate a wide variety of data models, which range from throughput-oriented batch-processing jobs to latency-sensitive serving of data to our users. The Bigtable clusters used by these products span a wide range of sizes, from a handful to thousands of servers, and store up to several hundred terabytes of data.

In many ways, Bigtable resembles a database: it stores many implementation strategies with databases. Parallel databases [14] and main-memory databases [13] have

## Bigtable: A Distributed Storage System for Structured Data

### Bigtable: A Distributed Storage System for Structured Data

Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach  
Mike Burrows, Tushar Chandra, Andrew Fikes, Robert E. Gruber  
{fay,jeff,sanjay,wilson,hc,deborah,mike,tushar,fikes,gruber}@google.com

Google, Inc.

#### Abstract

Bigtable is a distributed storage system for managing structured data. It is designed to support very large data sets, petabytes of data across thousands of commodity servers. Many projects at Google store data in Bigtable, including web indexing, Google Earth, and Google Finance. These applications place very different demands on Bigtable, including serving URLs (e.g., URLs to web pages to static imagery) and latency requirements (from batch bulk processing to real-time data service). Despite these differences, Bigtable succeeds in providing a flexible, high-performance solution for all of these Google products. In this paper we describe the simple data model provided by Bigtable, which gives clients dynamic control over the reliability and scalability of the data out of memory or disk.

Section 2 describes the data model in more detail, and Section 3 provides an overview of the client API. Section 4 describes the distributed storage system infrastructure on which Bigtable depends. Section 5 describes the fundamentals of the Bigtable implementation, and Section 6 describes some of the refinements that we made to the original Bigtable prototype. Section 7 provides measurements of Bigtable's performance. We describe several examples of how Bigtable is used at Google and other Google products in the projects, including Google Analytics, Google Finance, Orkut, Personalized Search, Writely, and Google Earth. These products illustrate a wide variety of data models, which range from throughput-oriented batch-processing jobs to latency-sensitive serving of data to our users. The Bigtable clusters used by these products span a wide range of sizes, from a handful to thousands of servers, and store up to several hundred terabytes of data.

#### 2 Data Model

A Bigtable is a sparse, distributed, persistent multi-dimensional sorted map. The map is indexed by a row key, column key, and a timestamp; each value in the map is an uninterpreted array of bytes.

(row:string, column:string, time:int64) → string

## Dynamo: Amazon's Highly Available Key-value Store

### Dynamo: Amazon's Highly Available Key-value Store

Giuseppe DeCandia, Deniz Hasturk, Madan Jannani, Gunavardhan Kakulapati,  
Avinash Lakshman, Andrew Pilchin, Swapna Venkataraman, Peter Vosshall  
and Werner Vogels

Amazon.com

#### ABSTRACT

Reliable, highly available scale is one of the biggest challenges we face at Amazon.com, one of the largest e-commerce operations in the world. One of the slighter easier parts of maintaining financial consistency and serving customers well is the Amazon.com platform, which provides services for many web sites worldwide, is highly distributed, and has many components. The Amazon.com servers and network components located in many datacenters around the world must be highly available and largely continuous, and the way the persistent state is managed in the face of these failures drives the reliability and scalability of the software that runs on top of them.

This paper presents the design and implementation of Dynamo, a highly available key-value storage system that some of Amazon's core web sites use. Dynamo is designed to provide a consistent level of availability and data durability that is hard to achieve with existing systems. Dynamo's schema consistency uses a combination of client-side validation and a set of objects that provide a level of dynamic conflict resolution in a manner that provides a novel trade-off between consistency and performance.

Coresets and Subsample Developers

D4.2 Performance Systems: System Management, D.4.5

(Operating Systems): Reliability, D.4.2 (Operating Systems): Performance

General Terms

Algorithm, Architecture, Measurement, Performance, Design, Reliability

#### 1. INTRODUCTION

Amazon is a world-wide e-commerce platform that serves tens of millions of customers at peak times using tens of thousands of servers. The system must be highly available and meet strict operational requirements on Amazon's platform in terms of performance, reliability and efficiency, and to support continuous growth. One of the most important requirements is consistency, one of the most important requirements because even the slightest change in data can cause significant problems and loss of customer trust. In addition, to support continuous growth, the platform needs to be highly scalable.

There are many services on Amazon's platform that only need to provide a certain level of availability. Some of these services are those that provide best seller lists, shopping carts, customer support, and so on. These services do not require the common pattern of using a relational database would lead to inefficiencies and limit scale and availability. Dynamo employs a single master design, which is a good fit for the needs of these applications.

Dynamo uses a variation of well known techniques to achieve availability and scalability. Data is partitioned and replicated using consistent hashing [16], and consistency is facilitated by object replication. When a client performs a write, the update is maintained by a quorum-like technique and a decentralized replica synchronization protocol. Dynamo employs

**ACID**

# ACID! The King is Dead! Long Live the King!



atomic



consistent



isolated



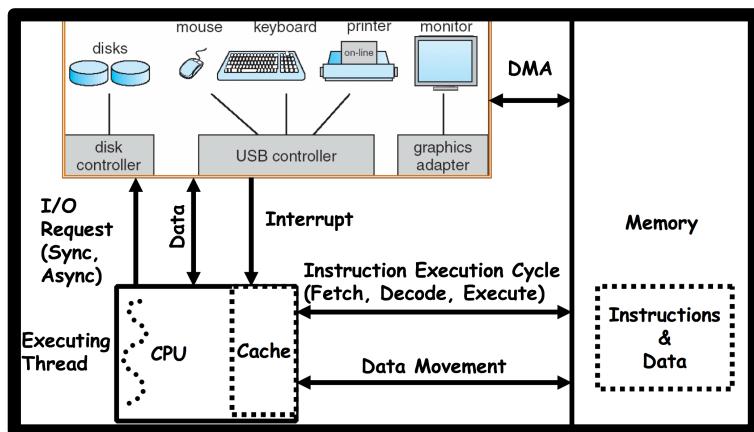
durable

# In order to keep CPU and I/O Busy...

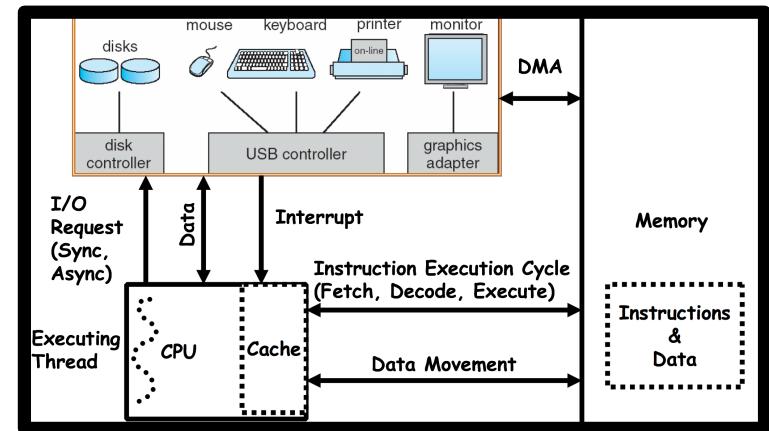
- We invented the concept of “Processes”.
- What are these?

# Process Abstraction...

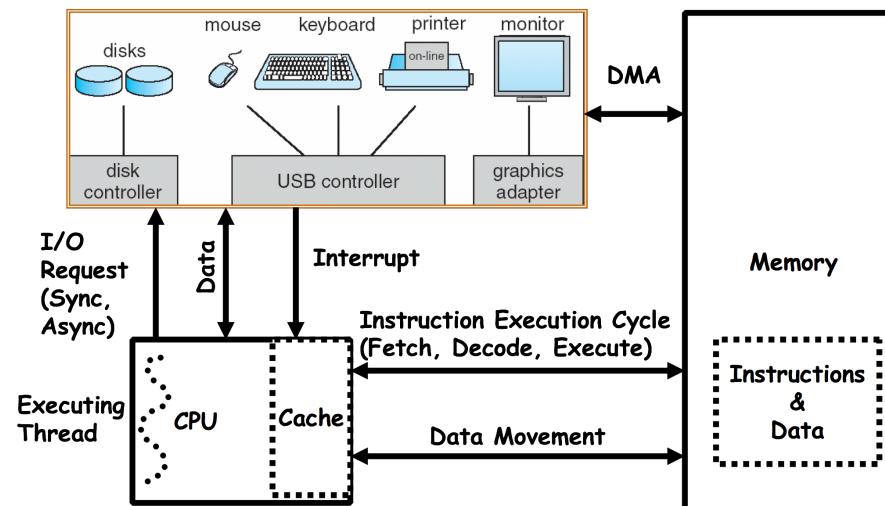
## Virtual Machine & CPU 1



## Virtual Machine & CPU 2



## The Real Machine & CPU



The Context Switch!

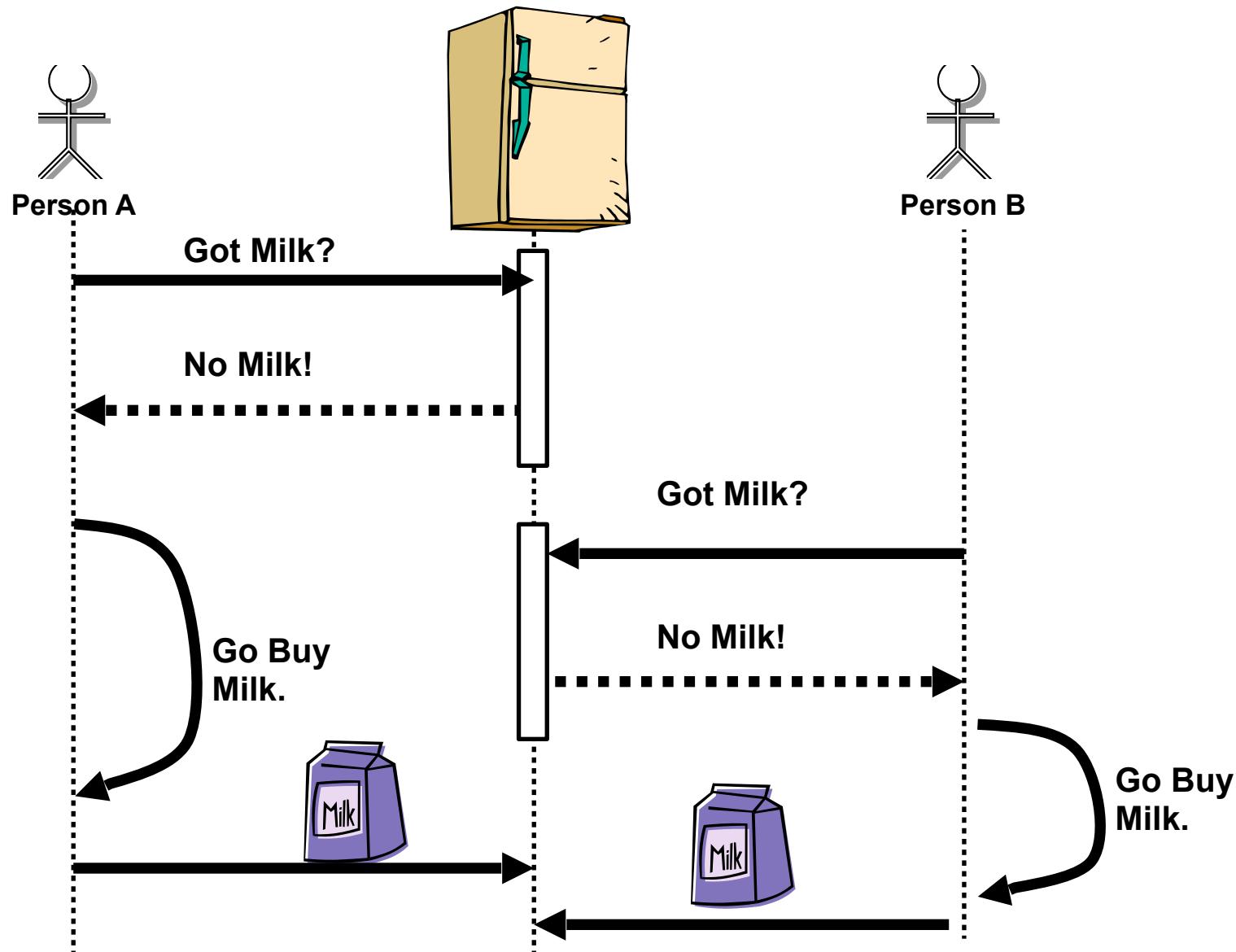
# Illusion of Virtual Machine not perfect!

We have to deal with concurrent access to shared memory.

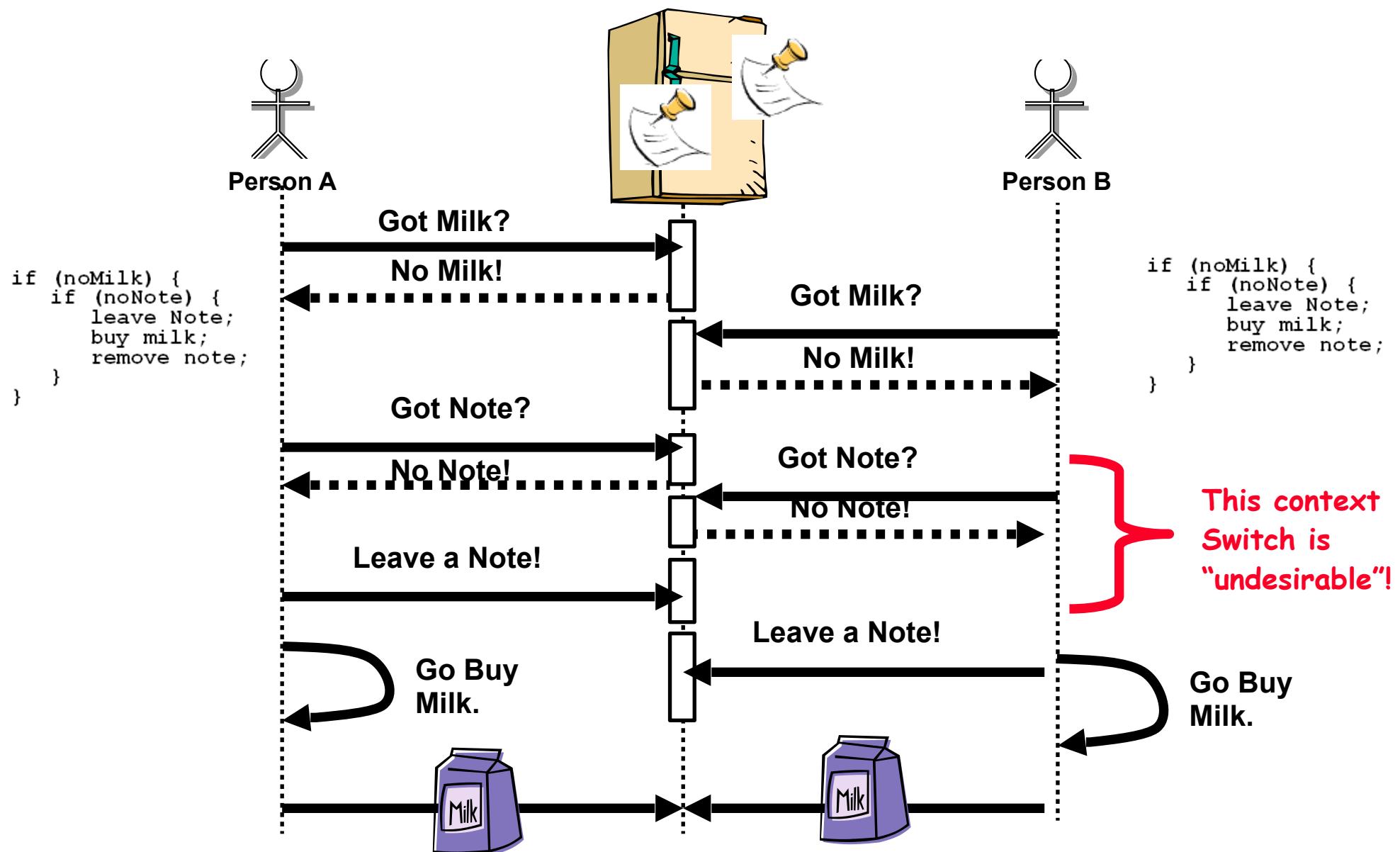
AKA => "Race Condition"

# Too Much Milk! – The Problem.

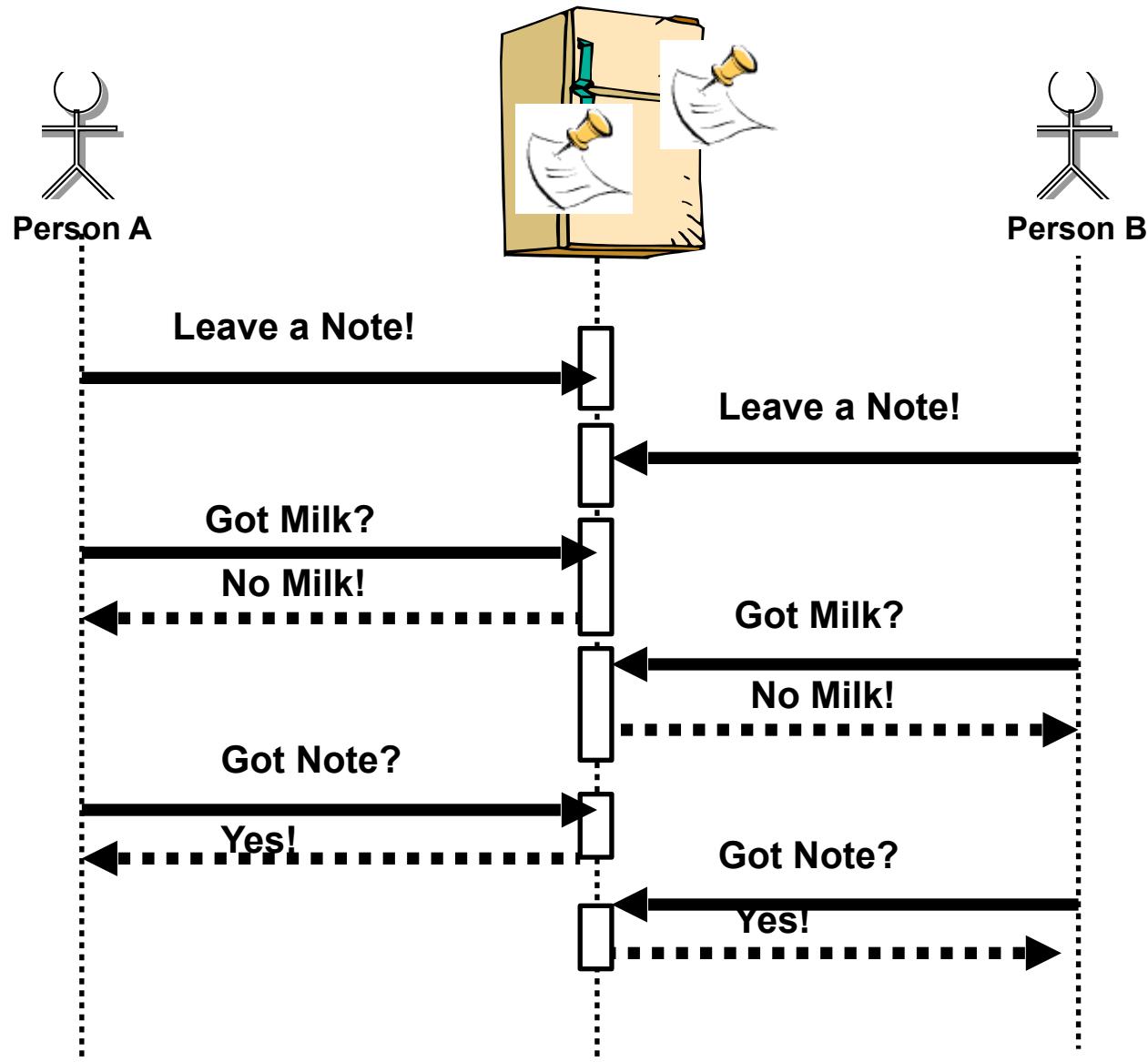
---



# Okay, we can solve this... just leave a note!



# How about leaving a note first?



```
leave Note;  
if (noMilk) {  
    if (noNote) {  
        leave Note;  
        buy milk;  
    }  
}  
remove note;
```

## Solutions to concurrent access/updates are founded on...

The basic idea that we need “**atomic**” operations which can not be interrupted (e.g. context switched in the middle).

Q: Are all H/W instructions atomic?

# Major Subsystems of an RDBMS

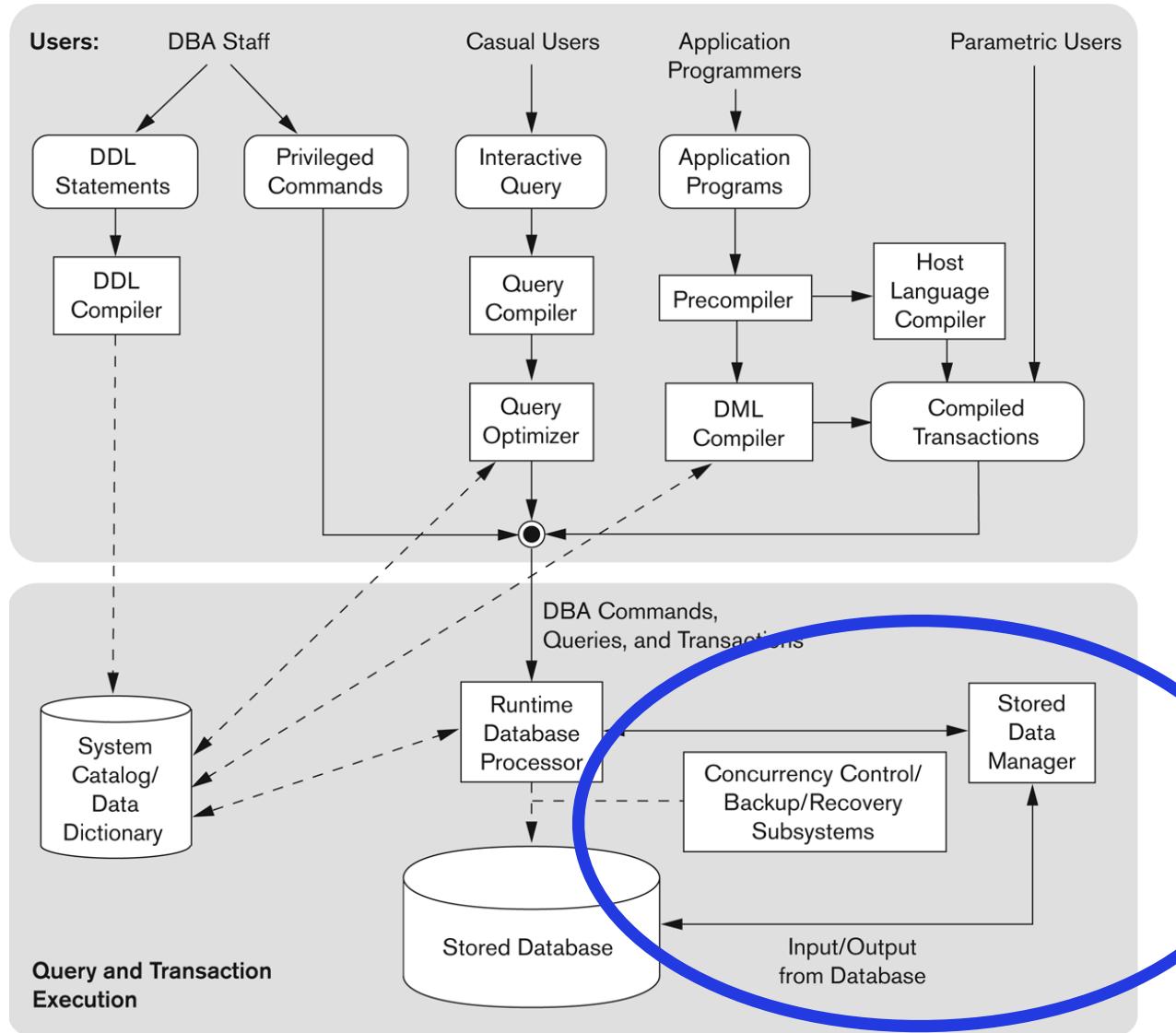


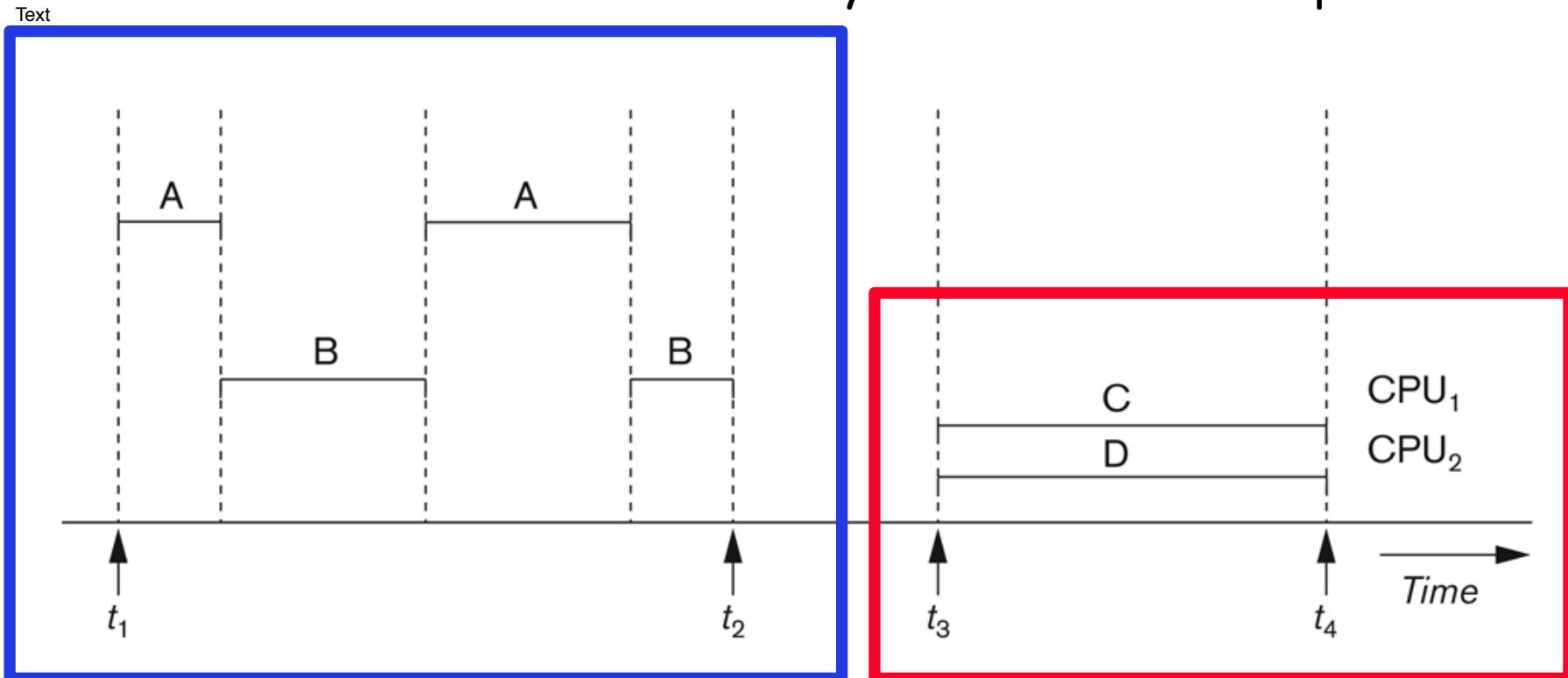
Figure 2.3

Component modules of a DBMS and their interactions.

What can go wrong?

# Transactions Executing Concurrently

- **Interleaved processing**
  - Concurrent execution of processes is interleaved in a single CPU
- **Parallel processing**
  - Processes are concurrently executed in multiple CPUs



## Sally's Program

---

- Sally executes the following two SQL statements, which we call (min) and (max), to help remember what they do.

(max)

```
SELECT MAX(price) FROM Sells  
WHERE bar = 'Joe''s Bar';
```

(min)

```
SELECT MIN(price) FROM Sells  
WHERE bar = 'Joe''s Bar';
```

## Joe's Program

---

- At about the same time, Joe executes the following steps, which have the mnemonic names (del) and (ins).

(del)      DELETE FROM Sells

                WHERE bar = 'Joe''s Bar';

(ins)      INSERT INTO Sells

                VALUES('Joe''s Bar', 'Heineken',  
                      3.50);

## Interleaving of Statements

---

- Although (max) must come before (min) and (del) must come before (ins)
- There are no other constraints on the order of these statements <sup>(1)</sup>.

## Example: Strange Interleaving

---

- Suppose the steps execute in the order

	(max)	(del)	(ins)	(min)	
	2.50, 3.00	2.50, 3.00		3.50	
Joe's Prices:	(max)	(del)	(ins)	(min)	
Statement:	3.00			3.50	
Result:					

- Sally sees MAX < MIN!

## Three Problems with Concurrent Execution

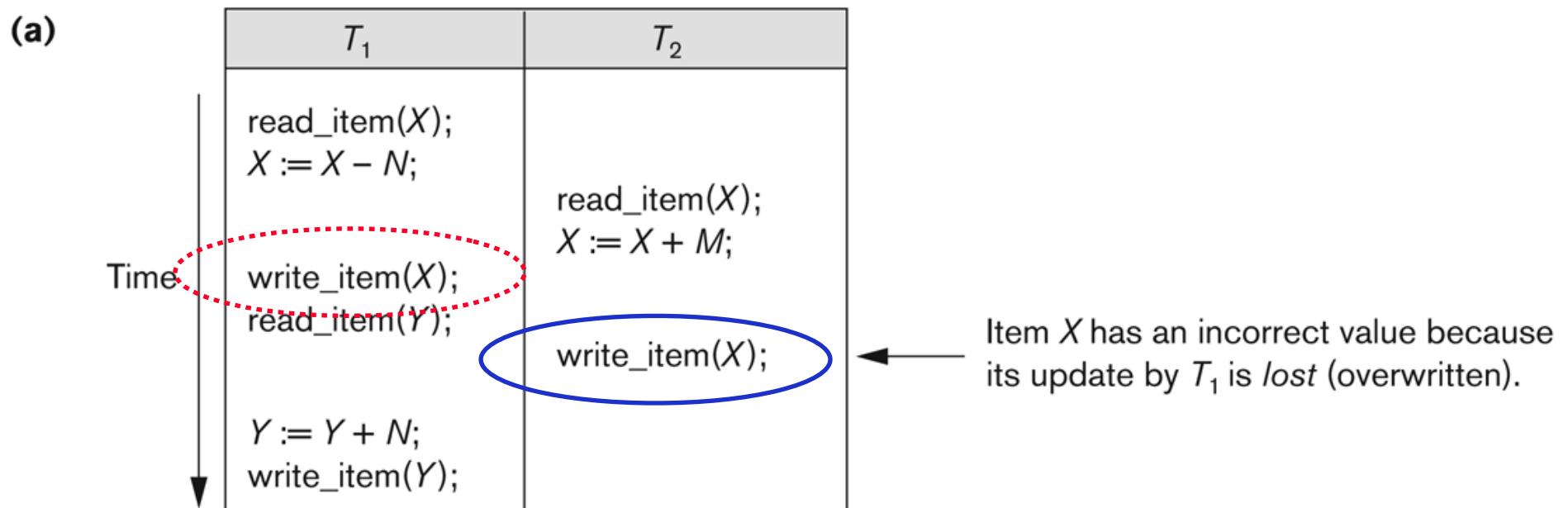
---

- The Lost Update Problem
  - Concurrent writes
- The Temporary Update Problem
  - Dirty read
- The Incorrect Summary Problem
  - Computing results based on inconsistent reads

# The lost update problem.

**Figure 17.3**

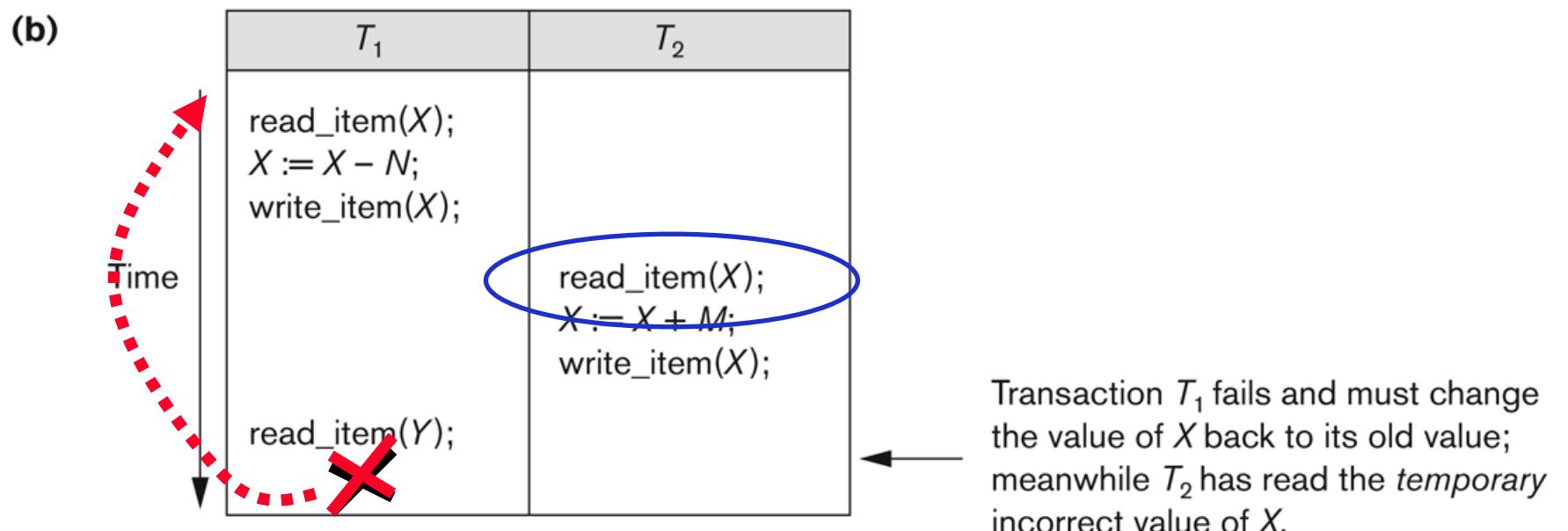
Some problems that occur when concurrent execution is uncontrolled. (a) The lost update problem. (b) The temporary update problem. (c) The incorrect summary problem.



# The temporary update problem.

**Figure 17.3**

Some problems that occur when concurrent execution is uncontrolled. (a) The lost update problem. (b) The temporary update problem. (c) The incorrect summary problem.



# The incorrect summary problem.

**Figure 17.3**

Some problems that occur when concurrent execution is uncontrolled. (a) The lost update problem. (b) The temporary update problem. (c) The incorrect summary problem.

$$A = 1 \quad (c)$$

$$X = 10$$

$$Y = 10$$

$$X = X - 5 \text{ (5)}$$

$$Y = 10 + 5 \text{ (15)}$$

$$A = 1$$

$$X = 5$$

$$Y = 15$$

$T_1$	$T_3$
<pre>read_item(X); X := X - N; write_item(X);</pre>	<pre>sum := 0; read_item(A); sum := sum + A; ⋮</pre> <pre>read_item(X); sum := sum + X; read_item(Y); sum := sum + Y;</pre> <pre>read_item(Y); Y := Y + N; write_item(Y);</pre>

$$\text{SUM} = 1$$

$$\text{SUM} = 1 + 5 + 10 \text{ (16)}$$

$T_3$  reads  $X$  after  $N$  is subtracted and reads  $Y$  before  $N$  is added; a wrong summary is the result (off by  $N$ ).

$$\text{SUM} = 16$$

# ACID



atomic



consistent



isolated



durable

# ACID Properties

---

- **Atomicity**
  - A transaction is an atomic unit of processing; it is either performed in its entirety or not performed at all.
- **Consistency preservation**
  - A correct execution of the transaction must take the database from one consistent state to another.
- **Isolation**
  - A transaction should not make its updates visible to other transactions until it is committed; this property, when enforced strictly, solves the temporary update problem and makes cascading rollbacks of transactions unnecessary (see Chapter 21).
- **Durability or permanency**
  - Once a transaction changes the database and the changes are committed, these changes must never be lost because of subsequent failure.

# Notation for Schedules

(a) 

$T_1$
read_item( $X$ ); $X := X - N$ ; write_item( $X$ ); read_item( $Y$ ); $Y := Y + N$ ; write_item( $Y$ );

(b) 

$T_2$
read_item( $X$ ); $X := X + M$ ; write_item( $X$ );

## Operations:

r = Read

w = Write

c = Commit

a = Abort

(a)  $T1 = r_1(X); w_1(X); r_1(Y); w_1(Y)$

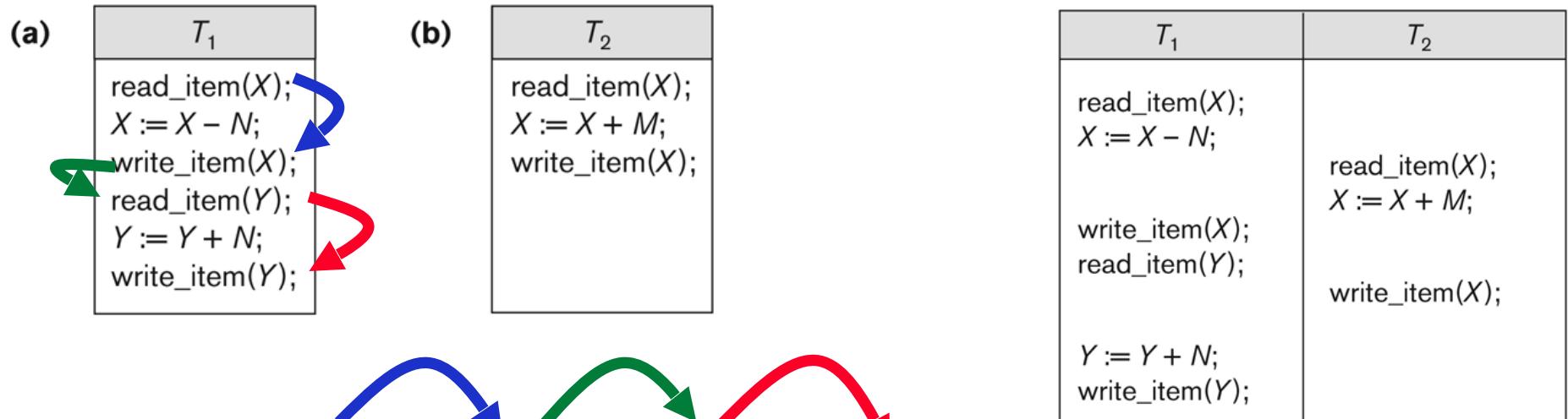
Operation

Transaction ID

Database Item

(b)  $T2 = r_2(X); w_2(X)$

## Notation for Schedules (cont.)



(a)  $T_1 = r_1(X); w_1(X); r_1(Y); w_1(Y)$

(b)  $T_2 = r_2(X); w_2(X)$

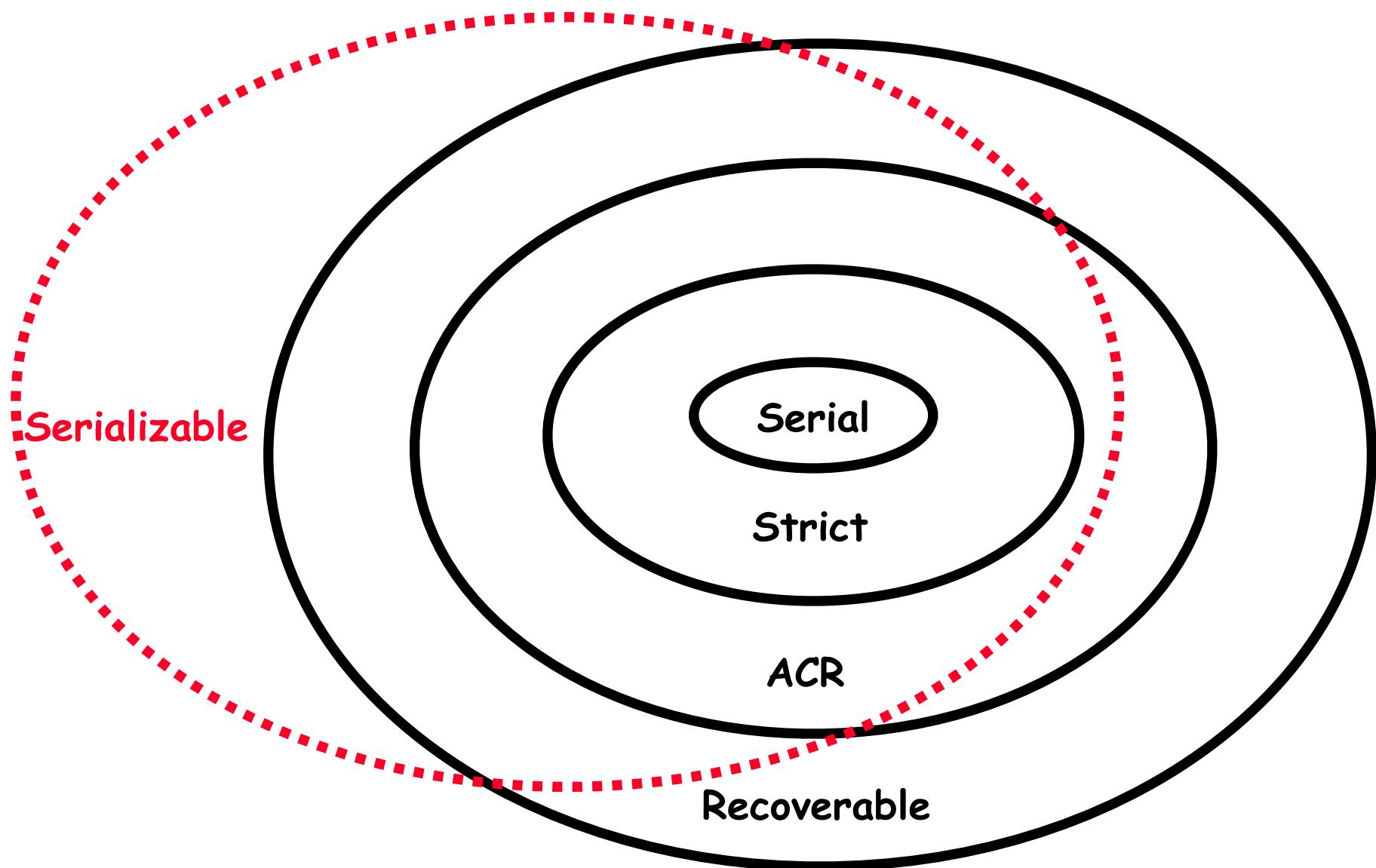
$S = T_1; T_2 = r_1(X); w_1(X); r_1(Y); w_1(Y); r_2(X); w_2(X)$

$S' = r_1(X); r_2(X); w_1(X); r_1(Y); w_2(X); w_1(Y)$

$S$  is a serial schedule

$S'$  is a schedule where operations are interleaved

# Classifying Schedules



**ACR = Avoids Cascading Rollbacks**

# Transaction Support in SQL2 (example)

---

- Sample SQL transaction:

```
EXEC SQL whenever sqlerror go to UNDO;
```

```
EXEC SQL SET TRANSACTION  
    READ WRITE  
    DIAGNOSTICS SIZE 5  
    ISOLATION LEVEL SERIALIZABLE;
```

```
EXEC SQL INSERT  
    INTO EMPLOYEE (FNAME, LNAME, SSN, DNO, SALARY)  
    VALUES ('Robert', 'Smith', '991004321', 2, 35000);
```

```
EXEC SQL UPDATE EMPLOYEE  
    SET SALARY = SALARY * 1.1  
    WHERE DNO = 2;
```

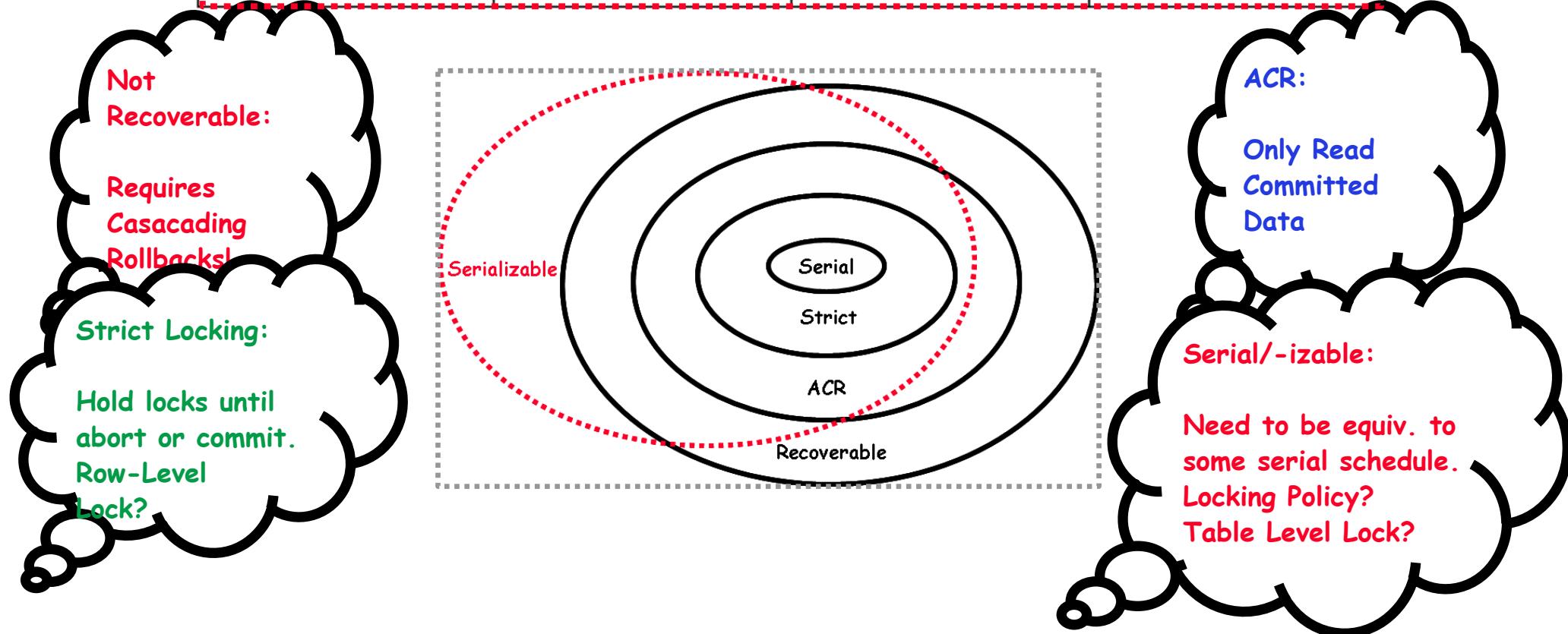
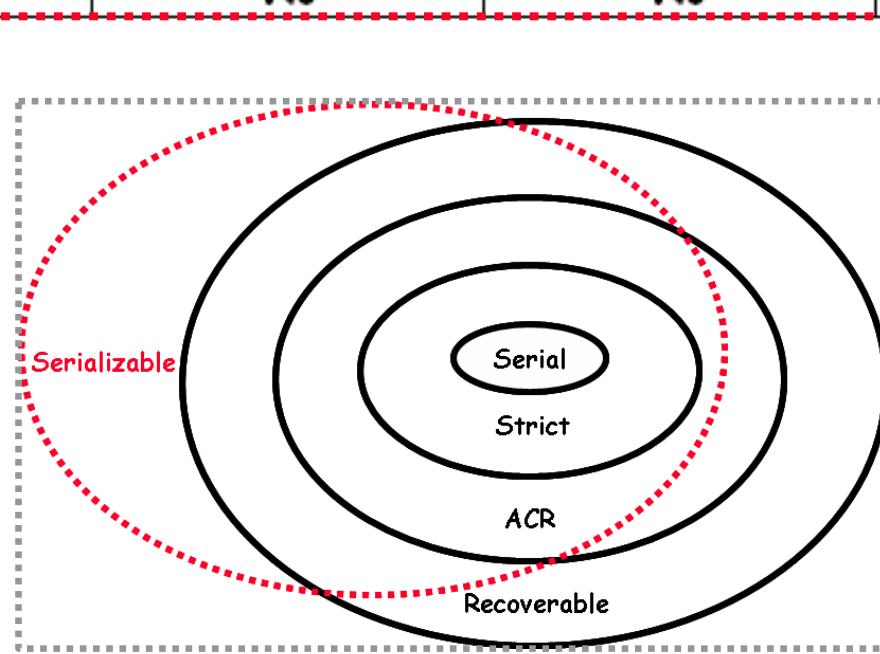
```
EXEC SQL COMMIT;  
GOTO THE_END;
```

```
UNDO: EXEC SQL ROLLBACK;
```

```
THE_END: ...
```

# Possible violation of serializability

Isolation Level	Type of Violation		
	Dirty Read	Non-Repeatable Read	Phantom Read
Read Uncommitted	Yes	Yes	Yes
Read Committed	No	Yes	Yes
Repeatable Read	No	No	Yes
Serializable	No	No	No



# Oracle Example

Transactions

# SET TRANSACTION

## Purpose

Use the SET TRANSACTION statement to establish the current transaction as read-only or read/write.

The operation transaction, issues a COMMIT

### READ WRITE

Specify READ WRITE to establish the current transaction as a read/write transaction. This clause establishes **statement-level read consistency**, which is the default.

## ISOLATION LEVEL Clause

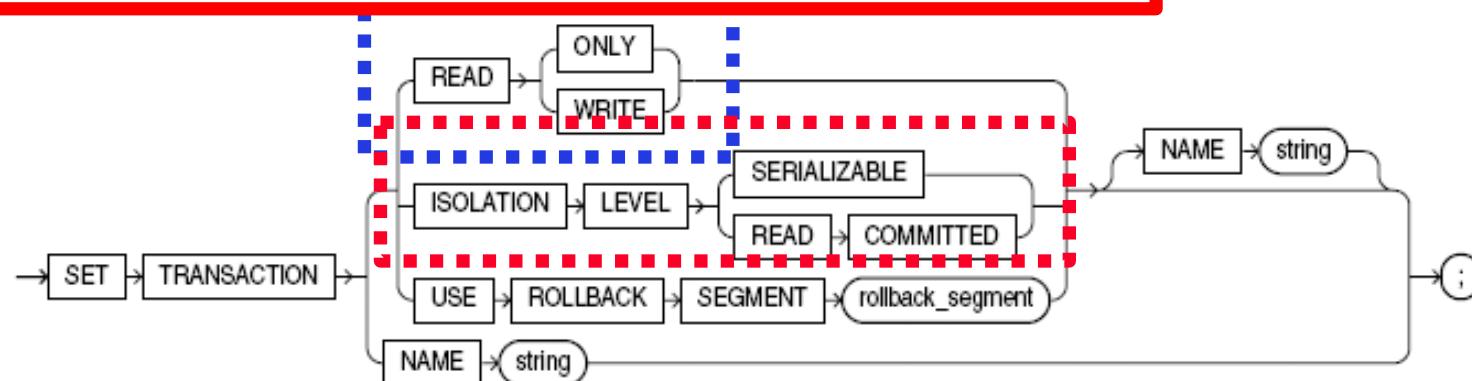
Use the ISOLATION LEVEL clause to specify how transactions containing database modifications are handled.

- The SERIALIZABLE setting specifies serializable transaction isolation mode as defined in the SQL92 standard. If a serializable transaction contains data manipulation language (DML) that attempts to update any resource that may have been updated in a transaction uncommitted at the start of the serializable transaction, then the DML statement fails.
- The READ COMMITTED setting is the default Oracle Database transaction behavior. If the transaction contains DML that requires row locks held by another transaction, then the DML statement waits until the row locks are released.

tion language (DDL)

ent transaction as a read-only transaction. **read consistency**.

ee only changes that were committed nsactions are useful for reports that run while other users update these same



## Transaction Example (What's the Outcome for Test 1 and Test 2?)

---

### Test:

```
Test 1: Default isolation levels for both.  
        S = T1(1);T2(1);T2(2);T1(2);T2(3);T1(3)  
  
Test 2: T1 at serializable. T2 at default.  
        S=T1(1);T2(1);T2(2);T2(3);T1(2);T1(3)
```

### Setup:

```
drop table bank_account ;  
create table bank_account  
( id number, type char(3), balance number(7,2)) ;  
insert into bank_account values ( 1, 'CHK', 1000) ;  
insert into bank_account values ( 2, 'SAV', 1000) ;
```

### T1 (Transfer from Checking to Savings):

1. select \* from bank\_account ;
2. update bank\_account set balance = balance - 500 where id = 1 ;
3. update bank\_account set balance = balance + 500 where id = 2 ;

### T2 (Transfer from Savings to Checking):

1. select \* from bank\_account ;
2. update bank\_account set balance = balance - 500 where id = 2 ;
3. update bank\_account set balance = balance + 500 where id = 1 ;

### Isolation Levels:

```
set transaction isolation level read committed ; // default  
set transaction isolation level serializable ;
```

# NoSQL



# Cloud Service Models

# The Long Tail

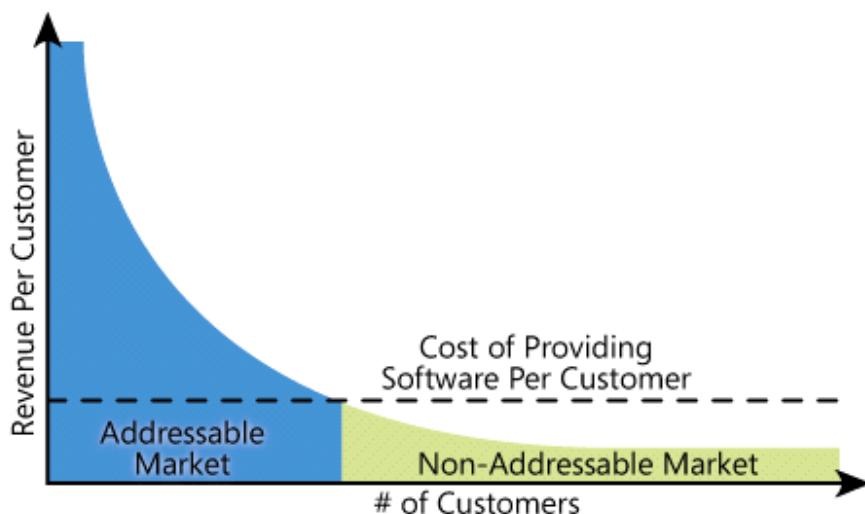
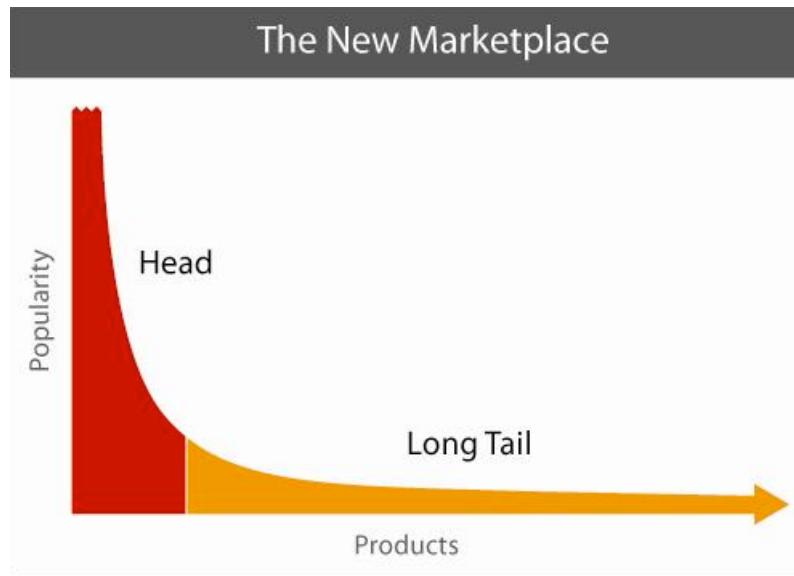


Figure 6. Market curve for LOB software vendors

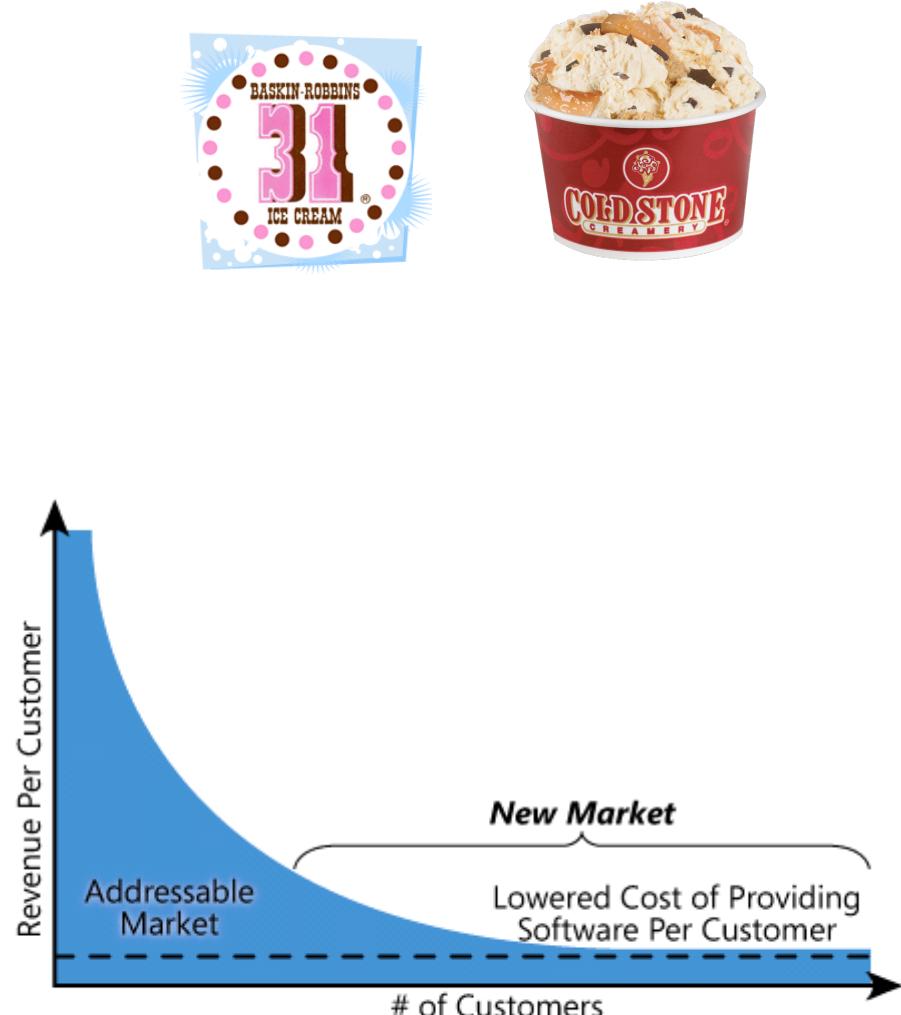


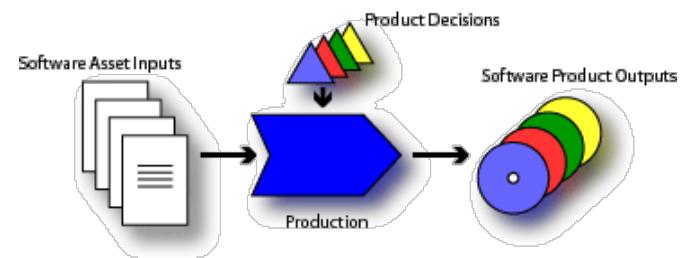
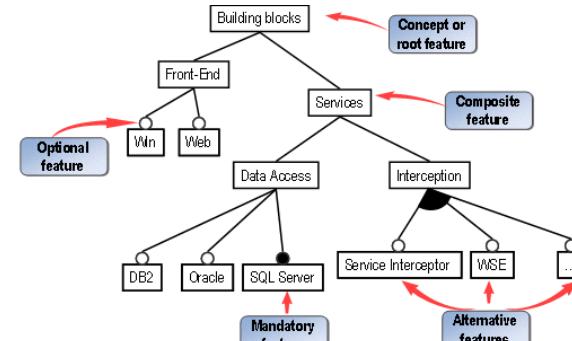
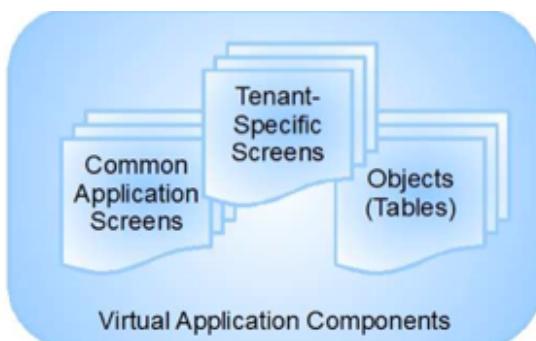
Figure 7. New market opened by lower cost of SaaS

# The Cloud: “*a brave new world*”

	Software	Platform	Infrastructure
Head	Manual Customization	JVM (or .Net)	Dedicated Hardware
Long Tail	Mass Customization	Multi-Paradigm	Virtual Machine
Challenges	Single Instance	Single Runtime	Scaling Data

# Polymorphic Application

- Current Approach
  - *MetaData Driven*
    1. Dedicated Database
    2. Shared DB, Fixed Extension Set
    3. Shared DB, Custom Extensions
- Emerging Approach:
  - *ByteCode Injection*
    1. Beyond Objects
    2. Aspect-Oriented Weaving
    3. Product-Line Engineering



# Polyglot Programming

- Current Approach

- Single Paradigm Design

1. Object-Oriented Analysis
2. Object-Oriented Design
3. Object-Oriented Programming

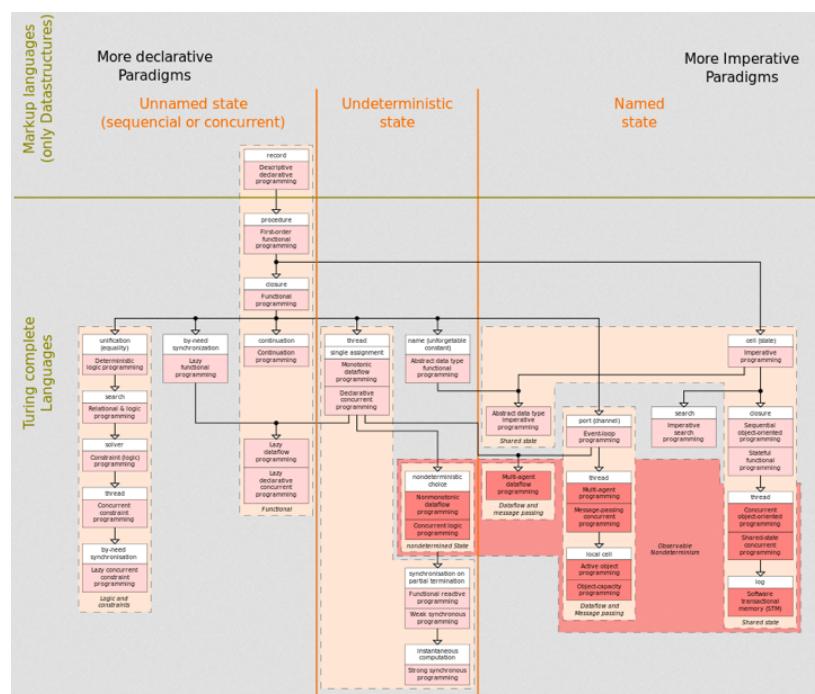


<http://polyglotprogramming.com/>

- Emerging Approach:

- Multi-Paradigm Design

1. Beyond Objects
2. Multiple Languages (GSL + DSL)
3. OOP + Functional + AOP



[http://en.wikipedia.org/wiki/File:Programming\\_paradigms.svg](http://en.wikipedia.org/wiki/File:Programming_paradigms.svg)

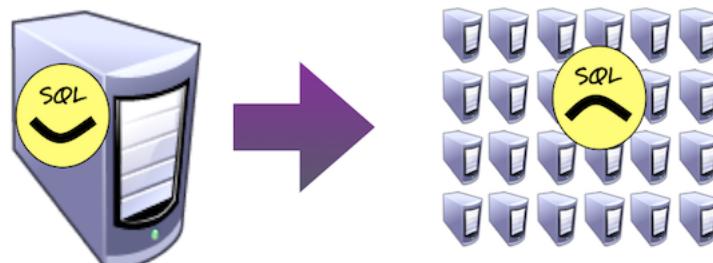
# Polyglot Persistence

- Current Approach
  - *Relational is King!*
  - 1. OLTP (*Transactions*)
  - 2. OLAP (*Analytics*)
  - 3. Centralized DB Server
  - 4. Scale-Up (*add expensive hardware*)
- Emerging Approach:
  - *Not Just Relational*
  - 1. SQL + NoSQL (*Eventual Consistency*)
  - 2. Hadoop + Columnar DB (*Analytics*)
  - 3. Distributed Data Stores
  - 4. Scale-Out (*add cheap computers*)

*SQL's dominance is cracking*

Relational databases are designed to run on a single machine, so to scale, you need buy a bigger machine

But it's cheaper and more effective to *scale horizontally* by buying lots of machines.



# The term “NoSQL”



No SQL Language  
(1998 Definition)

The term “*NoSQL*” was originally used in 1998 to describe a new [lightweight open-source relational DBMS](#) that did not expose an SQL interface. The name was then reprised and used 10 years later by [Eric Evans](#) (of Rackspace) and Johan Oskarsson (of Last.fm) for the [NoSQL meetup/conference](#) focusing on the topic of “*open-source, distributed, non relational databases*“.



Not Just SQL  
(2008 Definition)  
Evans & Oskarsson

**NoSQL**

Your Ultimate Guide to the Non - Relational Universe!

[the **best selected nosql link Archive** in the web]  
...never miss a *conceptual* article again...  
[News Feed](#) covering all changes [here](#) !

**NoSQL DEFINITION:** Next Generation Databases mostly addressing some of the points: being **non-relational, distributed, open-source** and **horizontally scalable**.

The original intention has been **modern web-scale databases**. The movement began early 2009 and is growing rapidly. Often more characteristics apply such as: **schema-free, easy replication support, simple API, eventually consistent / BASE** (not ACID), **a huge amount of data** and more. So the misleading term "*nosql*" (the community now translates it mostly with "**not only sql**") should be seen as an alias to something like the definition above. [based on 7 sources, 14 constructive feedback emails (thanksl!) and 1 disliking comment . Agree / Disagree? [Tell](#) me so! By the way: this is a strong definition and it is out there here since 2009!]

**LIST OF NOSQL DATABASES** [currently 150]

**EVENTS**  
Don't miss these events:  

- 21th May **NoSQL Roadshow Berlin** »
- 26.-27.April **NoSQL Matters Cologne** »

All past NoSQL Conferences »  
register your event here! »

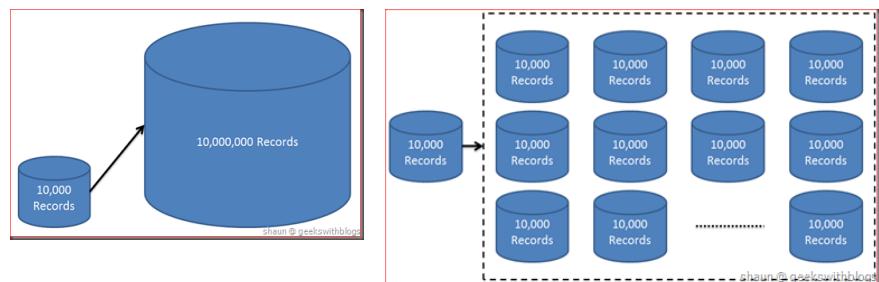
**NoSQL Search**



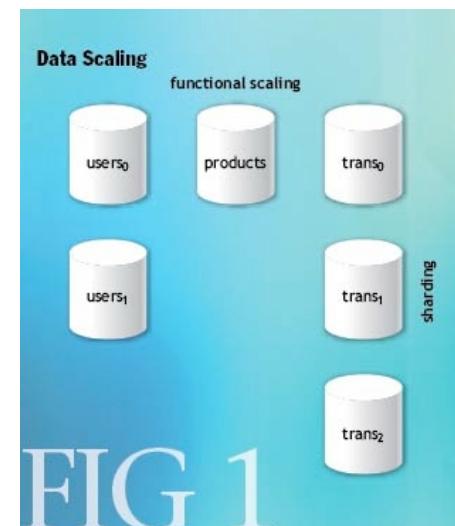
NoSQL is intended as shorthand for “not only SQL.” Complete architectures almost always mix traditional and NoSQL databases.

# The “Scaling” Problem

- **Vertical Scaling** (*aka:* Move everything to a bigger, faster infrastructure) is limited by certain limits (there is a hard-cap upper bound which is the fastest machine available at the time, is expensive, suffers from vendor lock, and is overall inefficient since the resources have to be acquired ahead of the time when they will be needed).
- **Sharding** (*aka:* Horizontal Partitioning) faces its own limitations (the Sharding Principle is not always easy to determine, inter-shard communications are still necessary, rebalancing is difficult and expensive) and enjoys limited support from the traditional RDBMSes vendors.
- **Horizontal Scaling**, on the other hand, seems to be the best solution for scaling a system over time. However, it is a difficult approach to implement when using a traditional RDBMS.



*Web applications have grown in popularity over the past decade. Whether you are building an application for end users or application developers (i.e., services), your hope is most likely that your application will find broad adoption—and with broad adoption will come transactional growth. If your application relies upon persistence, then data storage will probably become your bottleneck.*



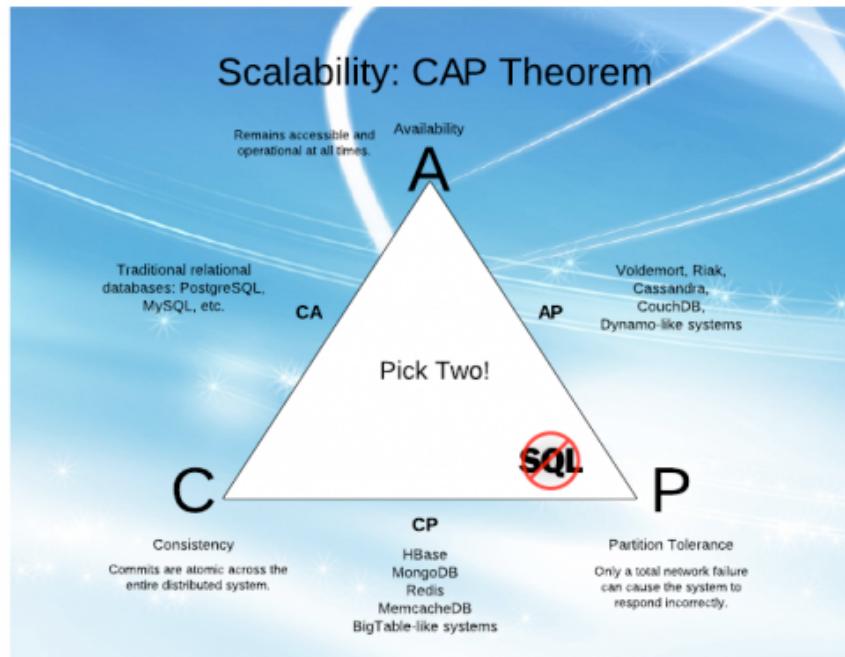
## Functional Partitioning

Functional partitioning is important for achieving high degrees of scalability. Any good database architecture will decompose the schema into tables grouped by functionality. Users, products, transactions, and communication are examples of functional areas. Leveraging database concepts such as foreign keys is a common approach for maintaining consistency across these functional areas.

# Brewer's CAP Theorem

Eric Brewer, a professor at the University of California, Berkeley, and cofounder and chief scientist at Inktomi, made the conjecture that Web services cannot ensure all three of the following properties at once (signified by the acronym CAP):<sup>2</sup> Specifically, a Web application can support, at most, only two of these properties with any database design. Obviously, any horizontal scaling strategy is based on data partitioning; therefore, designers are forced to decide between consistency and availability.

- **Consistency** means that each client always has the same view of the data.
- **Availability** means that all clients can always read and write.
- **Partition tolerance** means that the system works well across physical network partitions.



## **Consistency**



## **Availability**

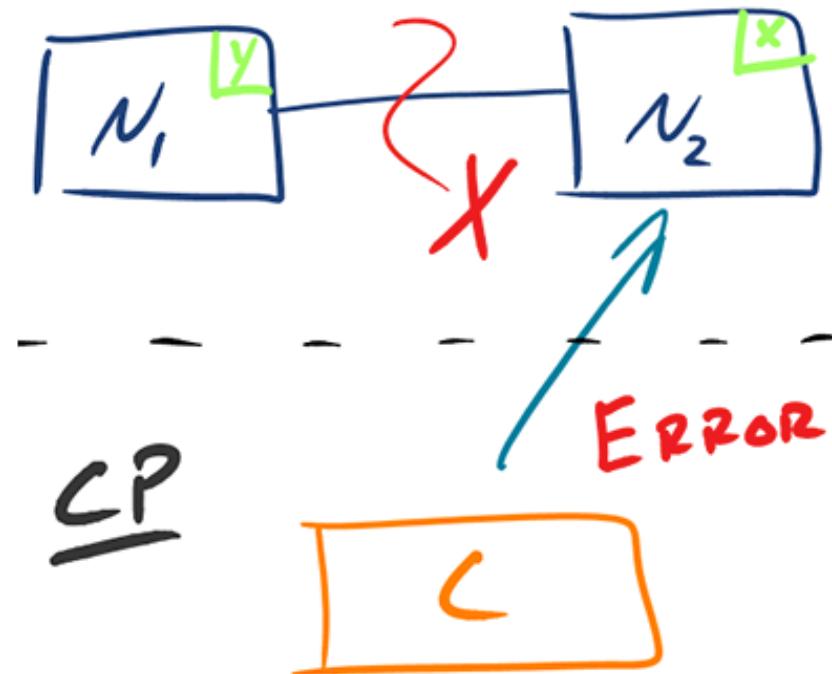


## **Partition Tolerance**

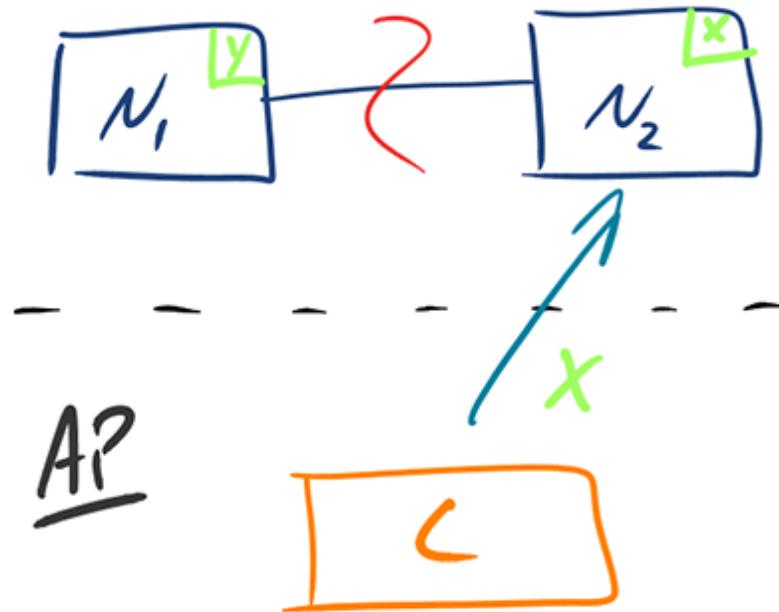


- **Consistency** - A read is guaranteed to return the most recent write for a given client.
- **Availability** - A non-failing node will return a reasonable response within a reasonable amount of time (no error or timeout).
- **Partition Tolerance** - The system will continue to function when network partitions occur.

- **CP** - Consistency/Partition Tolerance - Wait for a response from the partitioned node which could result in a timeout error. The system can also choose to return an error, depending on the scenario you desire. Choose Consistency over Availability when your business requirements dictate atomic reads and writes.



- **AP** - Availability/Partition Tolerance - Return the most recent version of the data you have, which could be stale. This system state will also accept writes that can be processed later when the partition is resolved. Choose Availability over Consistency when your business requirements allow for some flexibility around when the data in the system synchronizes. Availability is also a compelling option when the system needs to continue to function in spite of external errors (shopping carts, etc.)



# ACID vs. BASE

## ACID

In the database world, ACID is an acronym that says, all database transactions should be:

- **atomic**: transaction is either succeed or entirely rolled back
- **consistent**: transaction never invalidates the database state
- **isolated**: transactions should not interfere with each other but done as if they are sequential
- **durable**: completed transactions persist, i.e. what you wrote is what you read

The ACID property is a strict requirement for database. It impairs the availability and performance. For example, ACID makes database making frequent locks on tables to guarantee strong consistency. The lock hurts I/O throughput and make the database not responsive until the lock is over.

## BASE

BASE is described as the other end of a spectrum by Eric Brewer in his 2000 PODC keynote speech. It is:

- **basically available**: system looks to work and responsive at any time
- **soft-state**: no need to be consistent at all time
- **eventually consistent**: but become consistent at some later time

BASE means that, consistency after every transaction is not required, as long as the database is eventually in a consistent state. Therefore, using stale data or providing approximate answers are tolerated. In essence, BASE is a best-effort design that gives up strong consistency for weak consistency.

# ACID vs. BASE

The CAP theorem is given out by Eric Brewer in his PODC keynote speech and proved by Gilbert and Lynch in 2002. It says that it is impossible to achieve all three of:

- consistency: i.e., atomicity in ACID
- availability: i.e., responsiveness of the system
- partition tolerance: i.e., system works (amid partially) even with partial failure

## Examples of systems applying CAP theorem:

- single-site database: Forfeits P for C and A
- distributed database with distributed locking: Forfeits A for C and P
- web caching and DNS: Forfeits C for A and P

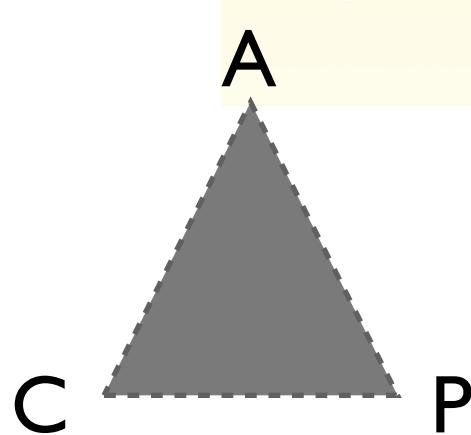
Therefore, in dealing with the CAP theorem, the strategy in designing a distributed system are:

1. Drop partition tolerance:
  - put everything into one machine
  - imposes scaling limits
2. Drop availability
  - upon partition event (i.e. partial system failure), affected service simply wait
  - until data is consistent, system remain unavailable/irresponsive
3. Drop consistency
  - accept things will become eventually consistent

# Examples of CAP Applied

**Consistent, Available (CA) Systems** have trouble with partitions and typically deal with it with replication. Examples of CA systems include:

- Traditional RDBMSs like Postgres, MySQL, etc (relational)
- Vertica (column-oriented)
- Aster Data (relational)
- Greenplum (relational)



**Consistent, Partition-Tolerant (CP) Systems** have trouble with availability while keeping data consistent across partitioned nodes. Examples of CP systems include:

- [BigTable](#) (column-oriented/tabular)
- [Hypertable](#) (column-oriented/tabular)
- [HBase](#) (column-oriented/tabular)
- [MongoDB](#) (document-oriented)
- [Terrastore](#) (document-oriented)
- [Redis](#) (key-value)
- [Scalaris](#) (key-value)
- [MemcacheDB](#) (key-value)
- [Berkeley DB](#) (key-value)

**Available, Partition-Tolerant (AP) Systems** achieve "eventual consistency" through replication and verification. Examples of AP systems include:

- [Dynamo](#) (key-value)
- [Voldemort](#) (key-value)
- [Tokyo Cabinet](#) (key-value)
- [KAI](#) (key-value)
- [Cassandra](#) (column-oriented/tabular)
- [CouchDB](#) (document-oriented)
- [SimpleDB](#) (document-oriented)
- [Riak](#) (document-oriented)

## **Key Value stores**

*Examples: Tokyo Cabinet/Tyrant, Redis, Voldemort, Oracle BDB*  
*Typical applications: Content caching*  
*Strengths: Fast lookups* *Weaknesses: Stored data has no schema*

## **Graph databases**

*Examples: Neo4J, InfoGrid, Infinite Graph* *Typical applications: Social networking, Recommendations*  
*Strengths: Graph algorithms e.g. shortest path, connectedness, n degree relationships, etc.*  
*Weaknesses: Has to traverse the entire graph to achieve a definitive answer. Not easy to cluster.*

## **Distributed Peer Stores**

*Examples: Cassandra, HBase, Riak* *Typical applications: Distributed file systems*  
*Strengths: Fast lookups, good distributed storage of data*  
*Weaknesses: Very low-level API .*

NoSQL is a large and expanding field, for the purposes of this paper the common features of NoSQL data stores are:

- Easy to use in conventional load-balanced clusters
- Persistent data (not just caches)
- Scale to available memory
- Have no fixed schemas and allow schema migration without downtime
- Have individual query systems rather than using a standard query language
- Are ACID within a node of the cluster and eventually consistent across the cluster



Rule of Thumb: NoSQL's primary goal is to achieve horizontal scalability. It attains this by reducing transactional semantics and referential integrity.

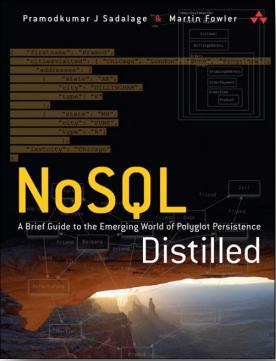
## **Document databases**

*Examples: CouchDB, MongoDb* *Typical applications: Web applications*  
*Strengths: Tolerant of incomplete data*  
*Weaknesses: Query performance, no standard query syntax*

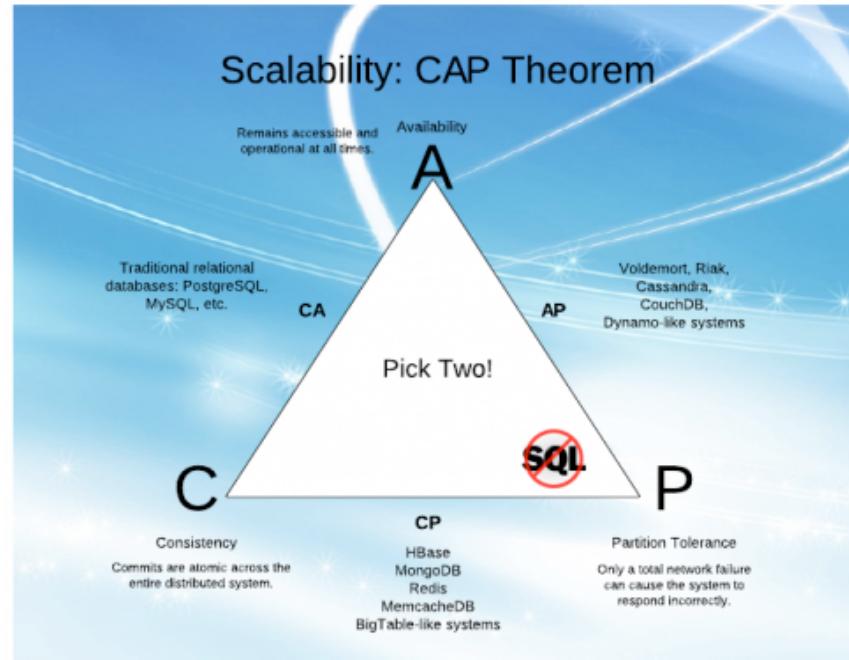
## **Object stores**

*Examples: Oracle Coherence, db4o, ObjectStore, GemStone, Polar* *Typical applications: Finance systems*  
*Strengths: Matches OO development paradigm, low-latency ACID, mature technology* *Weaknesses: Limited querying or batch-update options*

# NoSQL Key Points



- Relational databases have been a successful technology for twenty years, providing persistence, concurrency control, and an integration mechanism.
- Application developers have been frustrated with the impedance mismatch between the relational model and the in-memory data structures.
- There is a movement away from using databases as integration points towards encapsulating databases within applications and integrating through services.
- The vital factor for a change in data storage was the need to support large volumes of data by running on clusters. Relational databases are not designed to run efficiently on clusters.
- NoSQL is an accidental neologism. There is no prescriptive definition—all you can make is an observation of common characteristics.
- The common characteristics of NoSQL databases are
  - Not using the relational model
  - Running well on clusters
  - Open-source
  - Built for the 21st century web estates
  - Schemaless
- The most important result of the rise of NoSQL is Polyglot Persistence.



*The CAP theorem asserts that any networked shared-data system can have only two of three desirable properties.*

*However, by explicitly handling partitions, designers can optimize consistency and availability, thereby achieving some trade-off of all three.*

**Eric Brewer, University of California, Berkeley**

# Modern CAP

CAP prohibits only a tiny part of the design space: perfect availability and consistency in the presence of partitions, which are rare.

When partitions are present, there is an incredible range of flexibility for handling partitions and recovering from them.

The modern CAP goal should be to maximize combinations of consistency and availability that make sense for the specific application.

1. Partition Mode Operation
2. Partition Recovery

# 2 out of 3

The easiest way to understand CAP is to think of two nodes on opposite sides of a partition.

Allowing at least one node to update state will cause the nodes to become inconsistent, thus forfeiting C.

Likewise, if the choice is to preserve consistency, one side of the partition must act as if it is unavailable, thus forfeiting A.

Only when nodes communicate is it possible to preserve both consistency and availability, thereby forfeiting P.

# CAP, BASE, ACID

The general belief is that for wide-area systems, designers cannot forfeit P and therefore have a difficult choice between C and A.

The NoSQL movement is about creating choices that focus on availability first and consistency second.

Databases that adhere to ACID properties (atomicity, consistency, isolation, and durability) do the opposite.

# CAP Confusion

“2 of 3” view is misleading on several fronts.

First, because **partitions are rare**, there is little reason to forfeit C or A when the system is not partitioned.

Second, the choice between C and A can occur many times within the same system at very fine granularity; not only can **subsystems** make different choices, but the choice can change according to the **operation** or even the specific **data** or user involved.

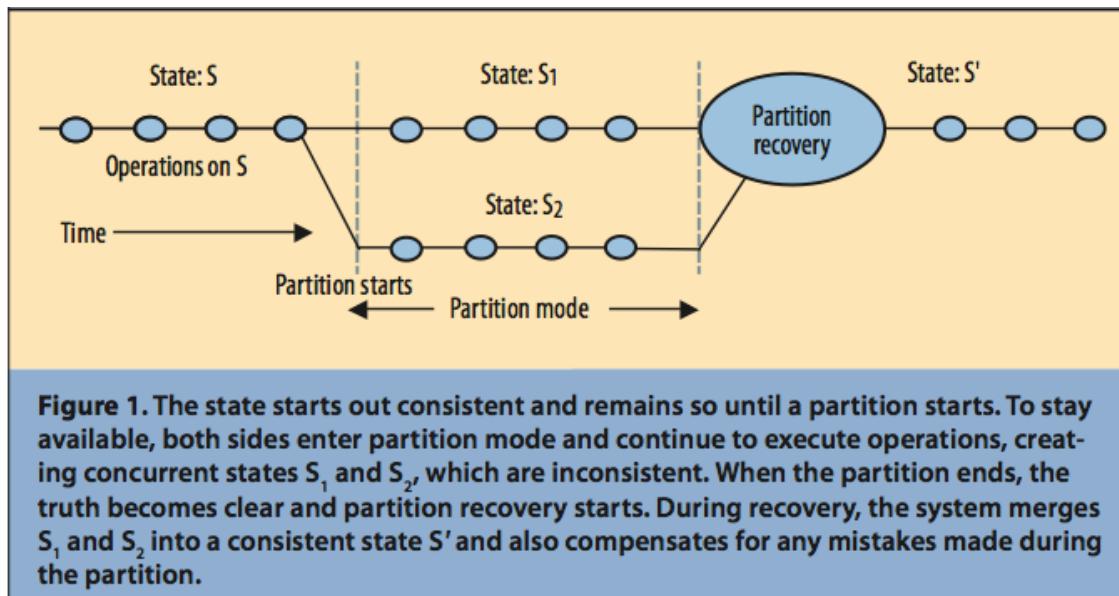
Finally, **all three properties are more continuous than binary**. Availability is obviously continuous from 0 to 100 percent, but there are also many levels of consistency, and even partitions have nuances, including disagreement within the system about whether a partition exists

# Manage the Partitions

The challenging case for designers is to mitigate a partition's effects on consistency and availability. The key idea is to manage partitions very explicitly, including not only detection, but also a specific recovery process and a plan for all of the invariants that might be violated during a partition.

Three Steps:

1. Detect the start of a partition
2. Enter an explicit partition mode that may limit some operations
3. Initiate partition recovery when communication is restored



# Partition Recovery

The designer must solve two hard problems during recovery:

- The state on both sides must become consistent, and
- There must be compensation for the mistakes made during partition mode.

Using **commutative operations** is the closest approach to a general framework for automatic state convergence.

The system concatenates logs, sorts them into some order, and then executes them. Commutativity implies the ability to rearrange operations into a preferred consistent global order. Unfortunately, using only commutative operations is harder than it appears; for example, addition is commutative, but addition with a bounds check is not (a zero balance, for example)

The ATM system designer could choose to prohibit withdrawals during a partition, since it is impossible to know the true balance at that time, but that would compromise availability. Instead, using stand-in mode (partition mode), modern ATMs limit the net withdrawal to at most  $k$ , where  $k$  might be \$200. Below this limit, withdrawals work completely; when the balance reaches the limit, the system denies withdrawals. Thus, the ATM chooses a sophisticated limit on availability that permits withdrawals but bounds the risk.

When the partition ends, there must be some way to both restore consistency and compensate for mistakes made while the system was partitioned. Restoring state is easy because the operations are commutative, but **compensation** can take several forms. A final balance below zero violates the invariant. In the normal case, the ATM dispensed the money, which caused the mistake to become external. The bank compensates by charging a fee and expecting repayment.

# Technology Radar

Prepared by the ThoughtWorks Technology Advisory Board - May 2013

## Platforms

### Adopt

- 29. Elastic Search
- 30. MongoDB
- 31. Neo4j
- 32. Redis
- 33. SMS and USSD as a UI

### Trial

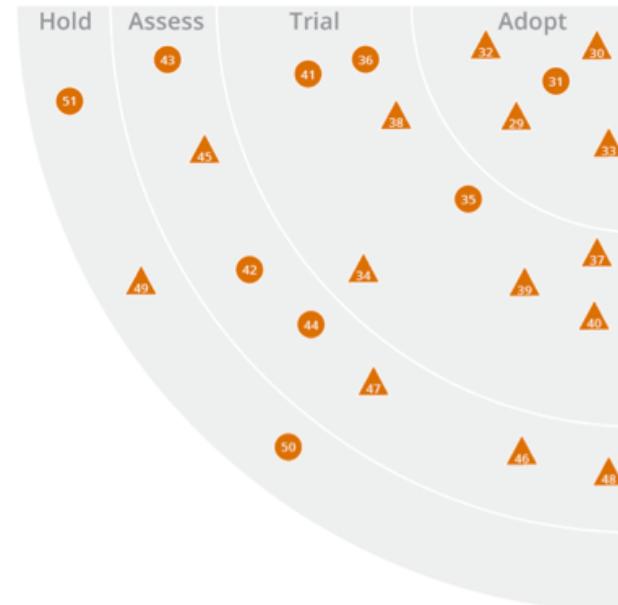
- 34. BigQuery
- 35. Continuous integration in the cloud
- 36. Couchbase
- 37. Hadoop 2.0
- 38. Node.js
- 39. OpenStack
- 40. Rackspace Cloud
- 41. Riak

### Assess

- 42. Azure
- 43. Calatrava
- 44. Datomic
- 45. PhoneGap/Apache Cordova
- 46. PostgreSQL for NoSQL
- 47. Vumi
- 48. Zepto.js

### Hold

- 49. Big enterprise solutions
- 50. Singleton infrastructure
- 51. WS-\*



*"For problems that fit the document database model, **MongoDB** is now the most popular choice"*

<http://www.thoughtworks.com/radar>

*"PostgreSQL is expanding to become the NoSQL choice of SQL databases"*

## Non-relational

### Analytic

Hadoop  
Teradata  
Aster  
Piccolo  
HPCC

## Relational

Netezza  
Infobright  
Teradata  
Calpont  
ParAccel  
EMC Greenplum  
Actian  
SAP Sybase IQ  
IBM InfoSphere  
VectorWise  
HP Vertica

### NoSQL

#### Graph

Neo4J  
InfiniteGraph  
OrientDB  
DEX  
NuvolaBase

#### Big tables

DataStax Enterprise  
Castle  
Acunu  
Citrusleaf  
BerkeleyDB  
Cassandra  
HBase  
Oracle NoSQL  
RethinkDB  
HandlerSocket\*  
Riak  
Redis-to-go  
SimpleDB  
LevelDB  
Redis  
Membrain  
Voldemort  
Couchbase

#### -as-a-Service

### SAP HANA

Oracle  
SkySQL  
MySQL  
PostgreSQL  
SQL Server

#### -as-a-Service

FathomDB  
Amazon RDS  
Postgres Plus Cloud  
Rackspace MySQL Cloud  
Google Cloud SQL

### NewSQL

#### -as-a-Service

NuoDB  
VoltDB  
MemSQL  
Drizzle  
StormDB  
Xeround  
GenieDB

#### Storage engines

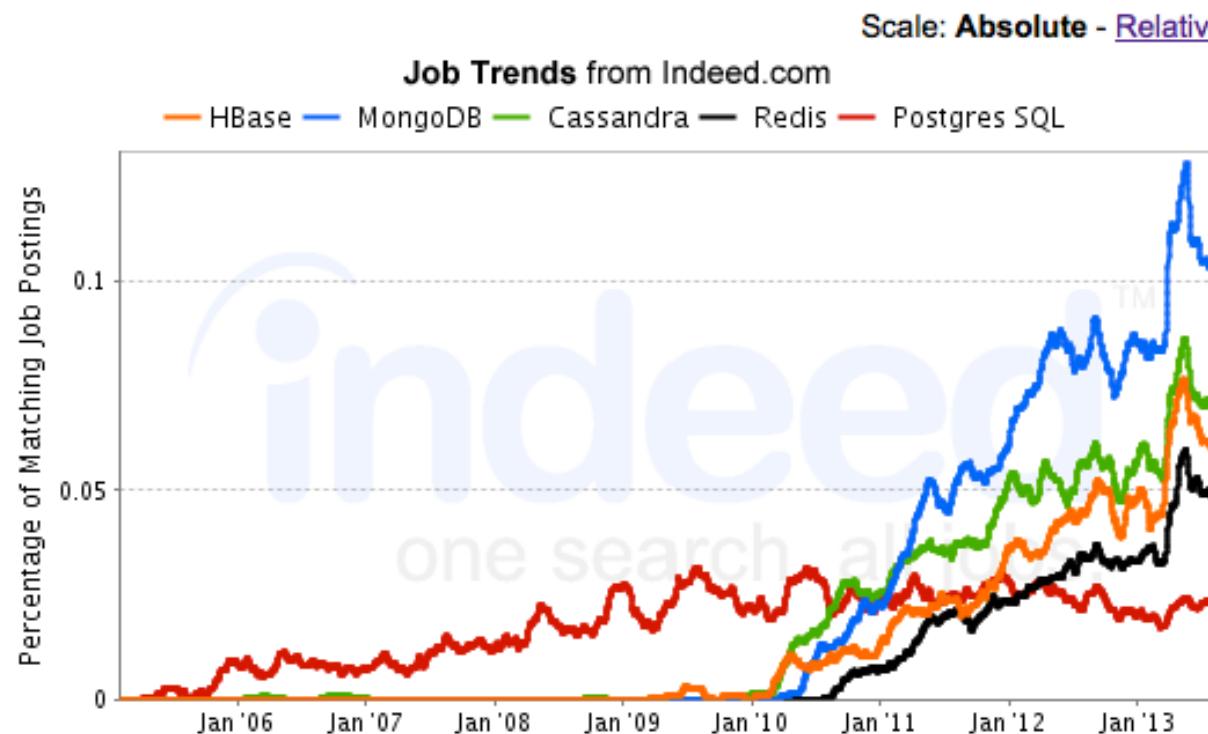
VoltDB  
JustOneDB  
Translattice  
Akiban  
SchoonerSQL  
Clustrix  
ScaleArc  
ParElastic  
Scale  
Continuent  
Galera  
CodeFutures  
MySQL Cluster  
ScaleBase

## Operational

Starcounter

InterSystems

# HBase, MongoDB, Cassandra, Redis, Postgres SQL Job Trends



Indeed.com searches millions of jobs from thousands of job sites.  
This job trends graph shows the percentage of jobs we find that contain your search terms.

Find [Hbase jobs](#), [Mongodb jobs](#), [Cassandra jobs](#), [Redis jobs](#), [Postgres SQL jobs](#)

Feel free to [▶ share this graph](#)

▶ [Email to a friend](#)

▶ [Post on your blog/website](#)

## Top Job Trends

1. [HTML5](#)
2. [MongoDB](#)
3. [iOS](#)
4. [Android](#)
5. [Mobile app](#)
6. [Puppet](#)
7. [Hadoop](#)
8. [jQuery](#)
9. [PaaS](#)
10. [Social Media](#)

<http://www.indeed.com/jobtrends>

# Aggregate Data Models

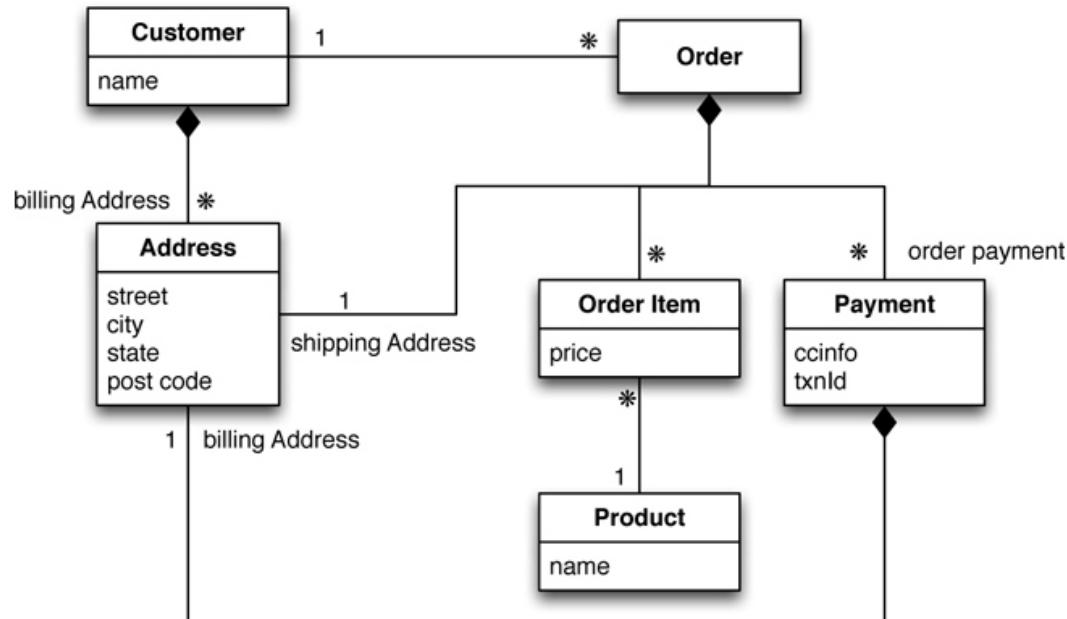
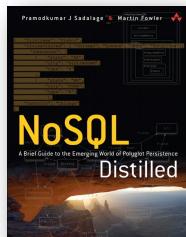


Figure 2.3. An aggregate data model

```
// in customers
{
  "id":1,
  "name":"Martin",
  "billingAddress":{ "city":"Chicago"}}

// in orders
{
  "id":99,
  "customerId":1,
  "orderItems":[
    {
      "productId":27,
      "price": 32.45,
      "productName": "NoSQL Distilled"
    }
  ],
  "shippingAddress":{ "city":"Chicago"}}
  "orderPayment":[
    {
      "ccinfo":"1000-1000-1000-1000",
      "txnid":"abelif879rft",
      "billingAddress": { "city": "Chicago"}
    }
  ],
} }
```



*"NoSQL Databases support atomic manipulation of a single aggregate at a time. This means that if we need to manipulate multiple aggregates in an atomic way, we have to manage that ourselves in the application code"*

# NoSQL Data Models

**Example of a standard structured  
(relational) database table for user  
records**

ID	First Name	Last Name
1	Samantha	Jones
2	James	Bond
3	James	Kirk

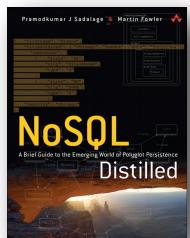
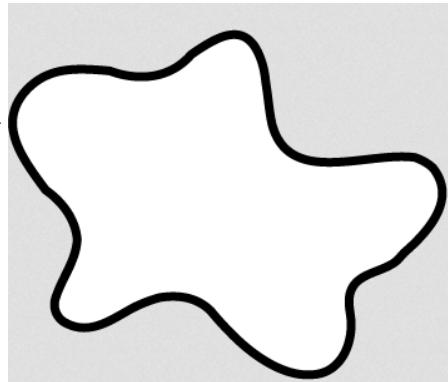
**Example of unstructured data for user records**

Key: 1	ID: sj	First Name: Sam
-----------	--------	-----------------

Key: 2	Email: jb@gmail.com	Location: London	Age: 37
-----------	------------------------	---------------------	------------

Key: 3	Facebook ID: jkirk	Password: xxx	Name: James
-----------	-----------------------	------------------	----------------

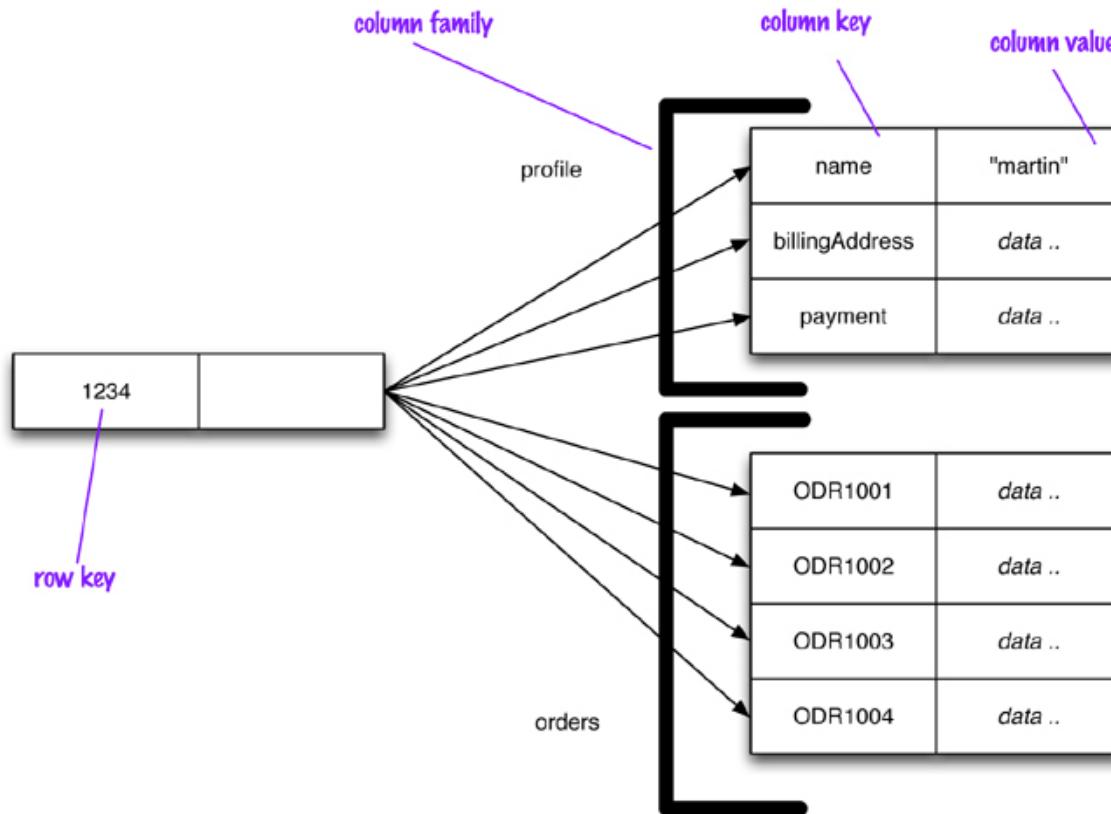
Key



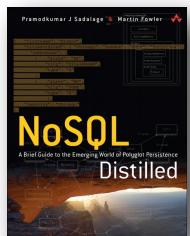
```
// in customers
{
  "id":1,
  "name":"Martin",
  "billingAddress": [ { "city":"Chicago"} ]
}

// in orders
{
  "id":99,
  "customerId":1,
  "orderItems":[
    {
      "productId":27,
      "price": 32.45,
      "productName": "NoSQL Distilled"
    }
  ],
  "shippingAddress": [ { "city":"Chicago"} ]
  "orderPayment": [
    {
      "ccinfo":"1000-1000-1000-1000",
      "txnid":"abelif879rft",
      "billingAddress": { "city": "Chicago"}
    }
  ],
}
```

# NoSQL Data Models

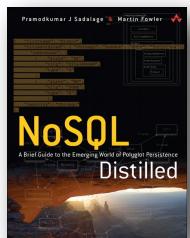
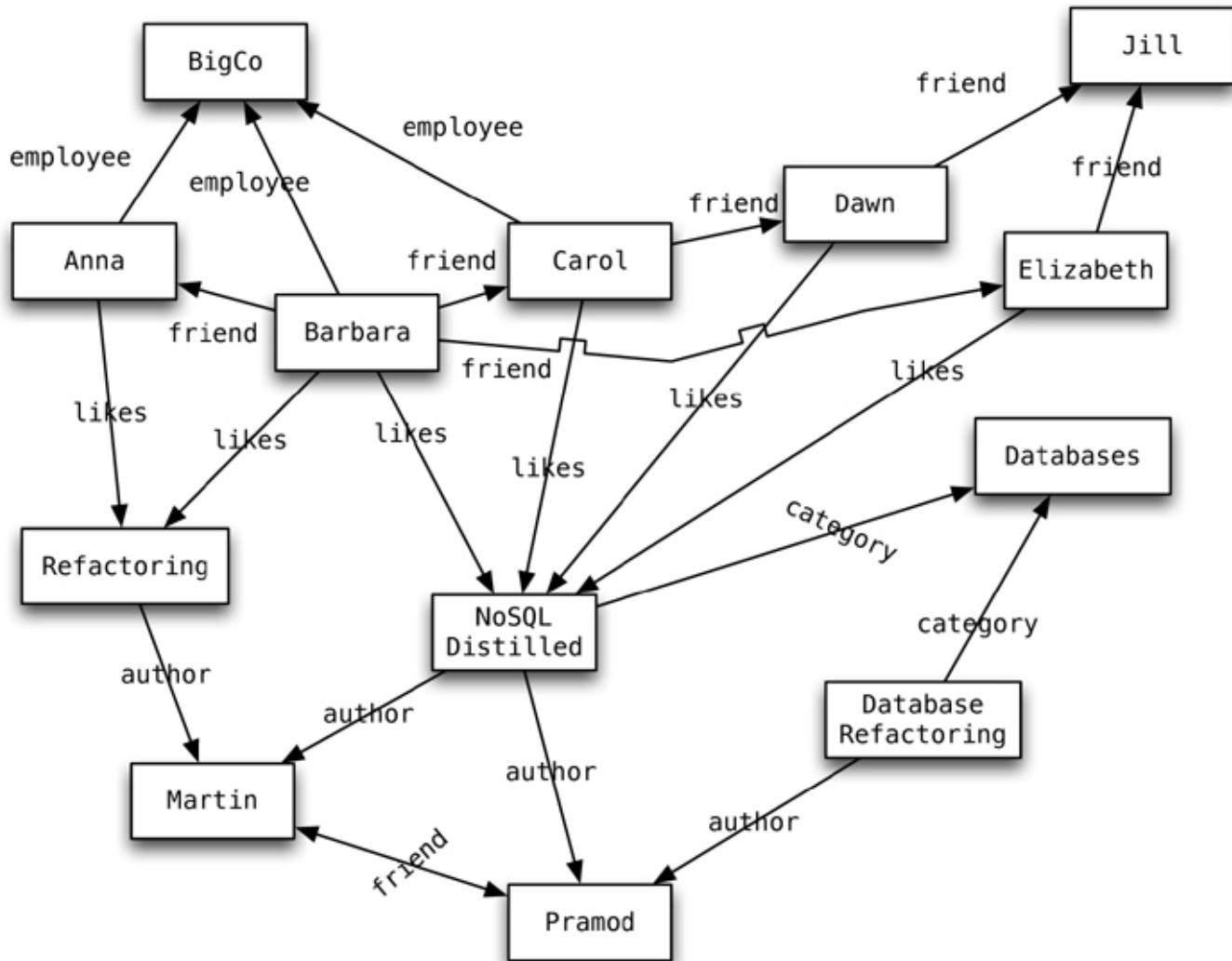


**Figure 2.5. Representing customer information in a column-family structure**



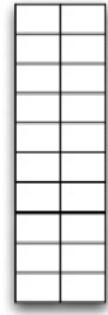
*"Column-family databases organize their columns into column families. Each column has to be part of a single column family, and the column acts as unit for access, with the assumption that data for a particular column family will be usually accessed together."*

# Graph Databases



# NoSQL Data Models

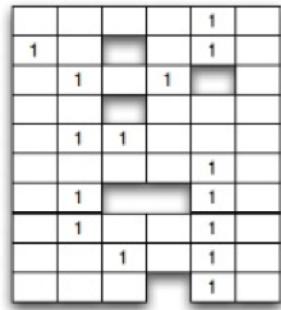
## Key-Value



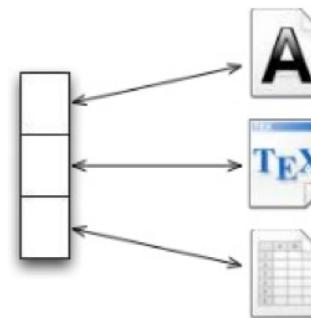
## Graph DB



## BigTable



## Document



# Distributions Models

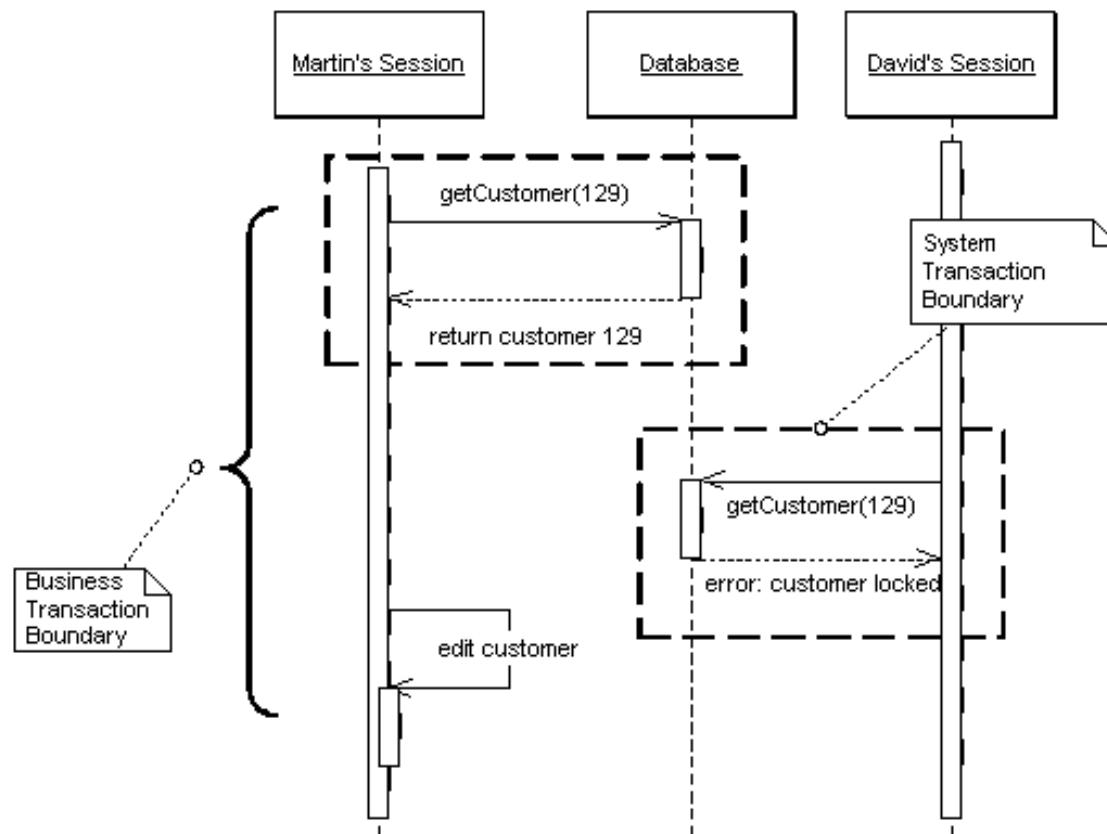
- Single-Server
  - Using NoSQL for Data Model
- Master-Slave Replication
  - Same Data is replicated from Master to Slaves. Master services all Writes. Reads services from Masters or Slaves.
- Sharding
  - Different Data is on separate nodes, each of which does its own Reads and Writes.
- Peer-To-Peer Replication
  - All Nodes services Reads and Writes to All Data.

# Pessimistic Offline Lock

by David Rice

*Prevents conflicts between concurrent business transactions by allowing only one business transaction at a time to access data.*

For a full description see P of EAA page 426

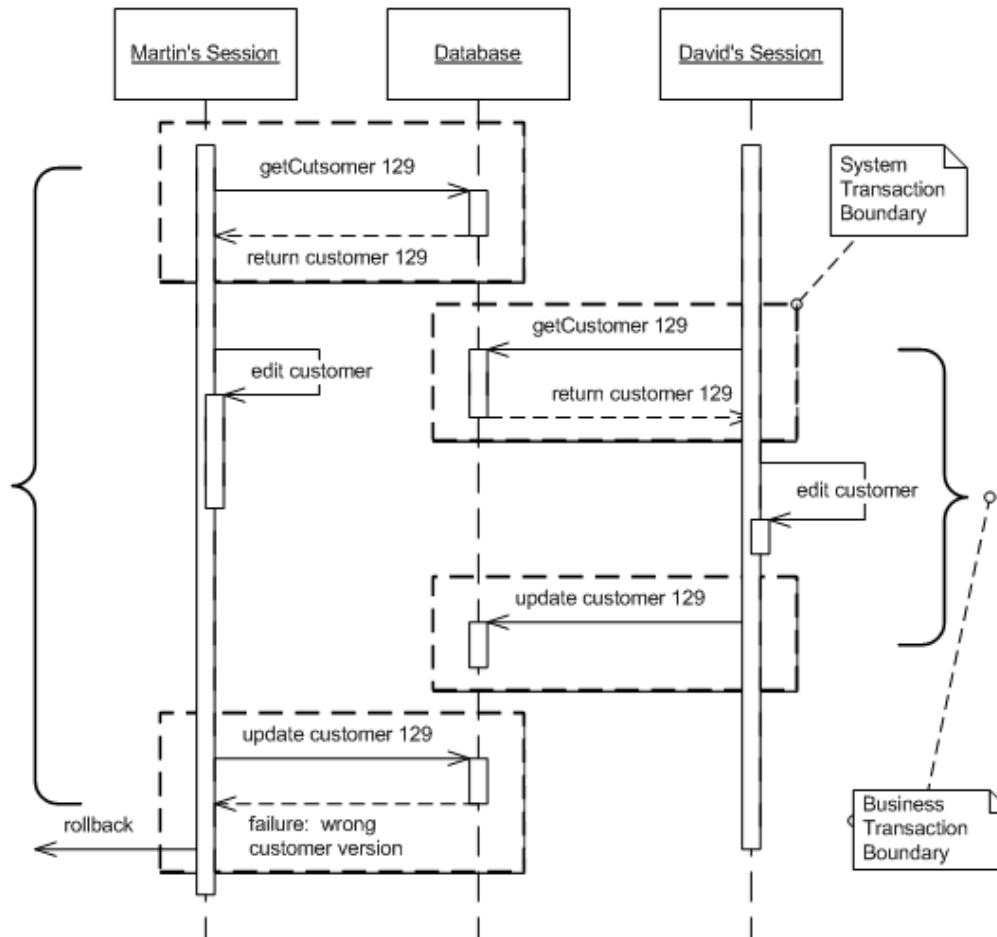


# Optimistic Offline Lock

by David Rice

*Prevents conflicts between concurrent business transactions by detecting a conflict and rolling back the transaction.*

For a full description see P of EAA page 416



<http://martinfowler.com/eaaCatalog/optimisticOfflineLock.html>

# Version Stamps

A field that changes each time data is updated

Examples:

*Counter*      Requires a Single Master to generate unique values.

*GUID*      Anyone can generate. But values are large and hard to compare.

*Hash*      Hash of content with a big enough hash key size for global uniqueness.

*Timestamp*      Nodes have to synchronize clocks.

*Vector Stamps*      A vector of version stamps from each node.  
Example: [Node-A: 2, Node-B: 5, Node-C: 10]

# Dynamo

Optimized for  
Availability

Key-Value  
(Opaque)

Consistent  
Hashing  
(Ring System)

Peer-To-Peer  
Replication

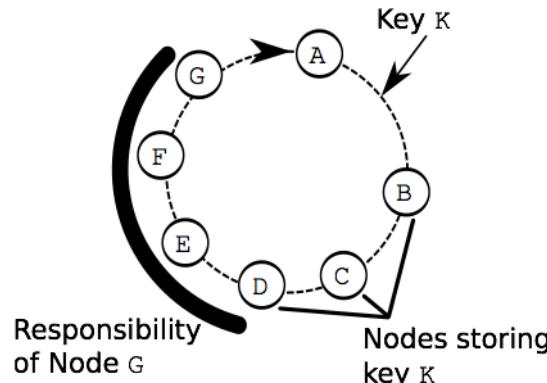


Figure 2.2: Key K is stored with a replication count of 3. It is stored on node B, while copies are kept on nodes C and D. Node G is responsible for all keys from D to G (left black arc).

Table 1: Summary of techniques used in *Dynamo* and their advantages.

Problem	Technique	Advantage
Partitioning	Consistent Hashing	Incremental Scalability
High Availability for writes	Vector clocks with reconciliation during reads	Version size is decoupled from update rates.
Handling temporary failures	Sloppy Quorum and hinted handoff	Provides high availability and durability guarantee when some of the replicas are not available.
Recovering from permanent failures	Anti-entropy using Merkle trees	Synchronizes divergent replicas in the background.
Membership and failure detection	Gossip-based membership protocol and failure detection.	Preserves symmetry and avoids having a centralized registry for storing membership and node liveness information.

# Big Table

Optimized for  
Consistency

Key/Record  
(Column-  
Oriented)

Master-Slave  
Replication

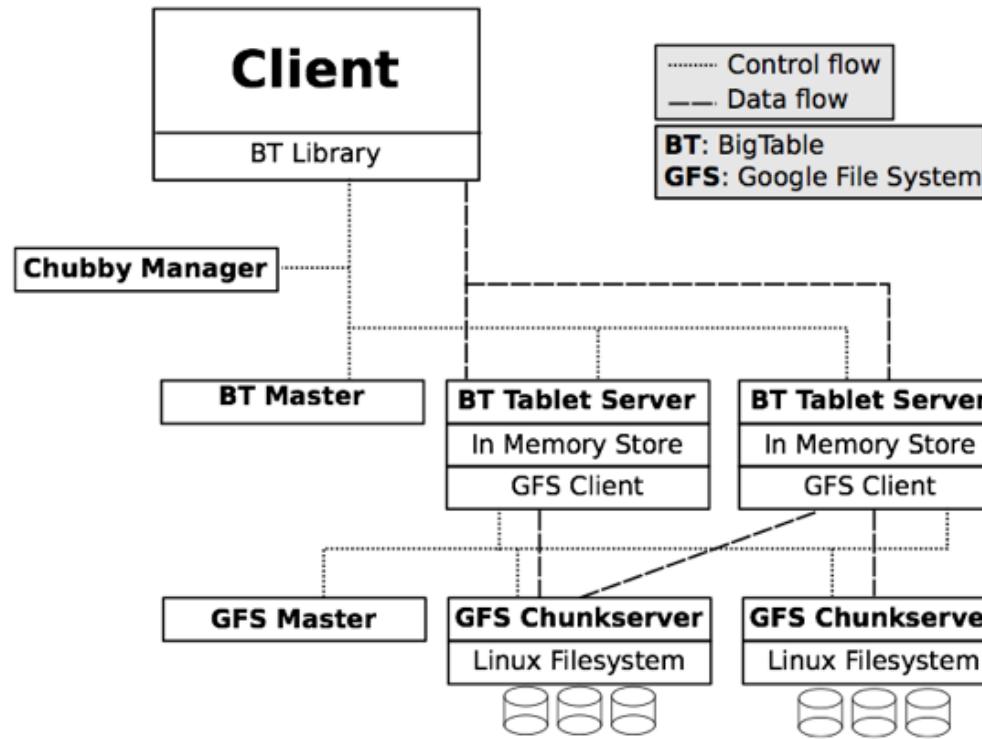


Figure 2.3: Architecture of Google BigTable: BigTable (BT) is created on top of two existing components: Google File System (GFS) and the Chubby lock service. While the control flow (dotted lines) is mainly between clients and the managers, the data (dashed lines) is transferred between the clients and the providers without the master.

# Comparison of Distribution Technologies in Different NoSQL Database Systems

	Programming Language	License	First Release
Apache Cassandra	Java	Apache License 2	July 2008
Basho Riak	Erlang, JavaScript	Apache License 2	2009
Project Voldemort	Java	Apache License 2	February 2009
Apache HBase	Java	Apache License 2	2007
Redis	C	New BSD	March 2009
CouchDB	Erlang, JavaScript	Apache License 2	2005
MongoDB	C++	GNU AGPL v3.0	Summer 2008
Membase	C++, Erlang	Apache License 2	2009

Table 4.1: Overview of the system which were examined

	Replication	Fragmentation	Category
Cassandra	✓	✓	
Riak	✓	✓	Ring (4.2)
Project Voldemort	✓	✓	
HBase	<i>by underlying DFS</i>	✓	
MongoDB	✓	✓	Master-Slave (4.3)
Membase	✓	✓	
Redis	✓ ( <i>unidirectional</i> )	<i>using external tools</i>	Replication based (4.4)
CouchDB	✓ ( <i>bidirectional</i> )	<i>using external tools</i>	

Table 4.2: Categorization and features of the systems

	Column-Oriented	Document-Oriented	Schemaless
Ring	Cassandra		Riak, Project Voldemort
Master-Slave	HBase	MongoDB	Membase
Replication based		CouchDB	Redis

Table 5.1: Systems compared in regard to data models and distribution categories

	Ring	Master-Slave	Replication Based
Single Point Of Failure	none	master and slaves	unidirectional: the master otherwise: none
Consistency Model	tunable consistency	strict consistency	eventual consistency
Availability	write: always available read: maybe unavailable	maybe unavailable	unidirectional: Maybe unavailable bidirectional: always available
Data Access	range scans make no sense (if randomly placed)	range scans fast (same slave)	range scans fast (all data available)

Table 4.3: Differences of the categories which have been discussed above

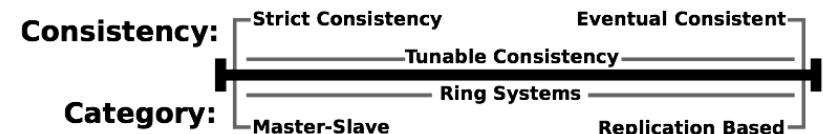


Figure 4.1: The categories introduced directly correlate with the consistency models they implement if put on a scale between strict consistency on the one side and eventual consistency on the other.

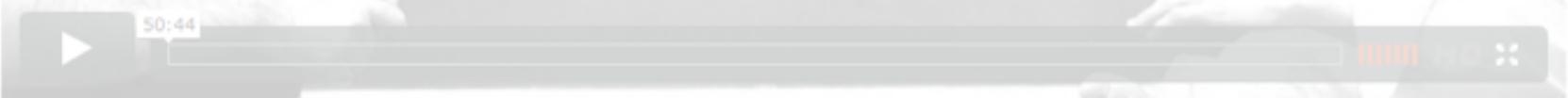
# Map Reduce



THE NOSQL TAPES WERE MADE POSSIBLE BY THE SUPPORT OF

**SCALITY**

YOUR VIDEO WILL START IN 4 SECONDS



AUTOPLAY IS ON

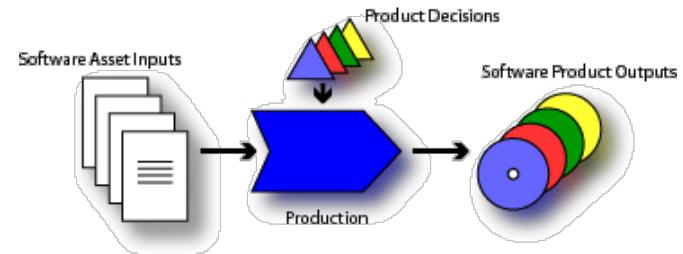
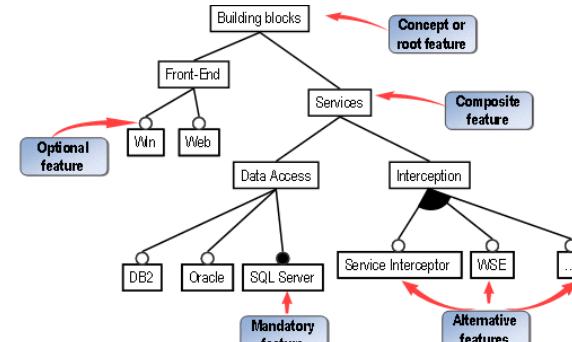
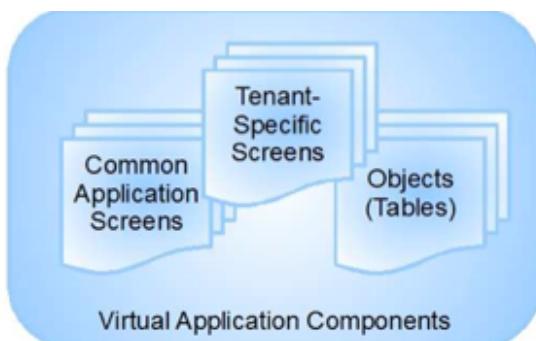
VOLUME 8

## UNDERSTANDING MAPREDUCE – WITH MIKE MILLER

RECORDED ON SEPTEMBER 26, 2010 IN SEATTLE, WA.

# Polymorphic Application

- Current Approach
  - *MetaData Driven*
    1. Dedicated Database
    2. Shared DB, Fixed Extension Set
    3. Shared DB, Custom Extensions
- Emerging Approach:
  - *ByteCode Injection*
    1. Beyond Objects
    2. Aspect-Oriented Weaving
    3. Product-Line Engineering



# Polyglot Programming

- Current Approach

- Single Paradigm Design

1. Object-Oriented Analysis
2. Object-Oriented Design
3. Object-Oriented Programming

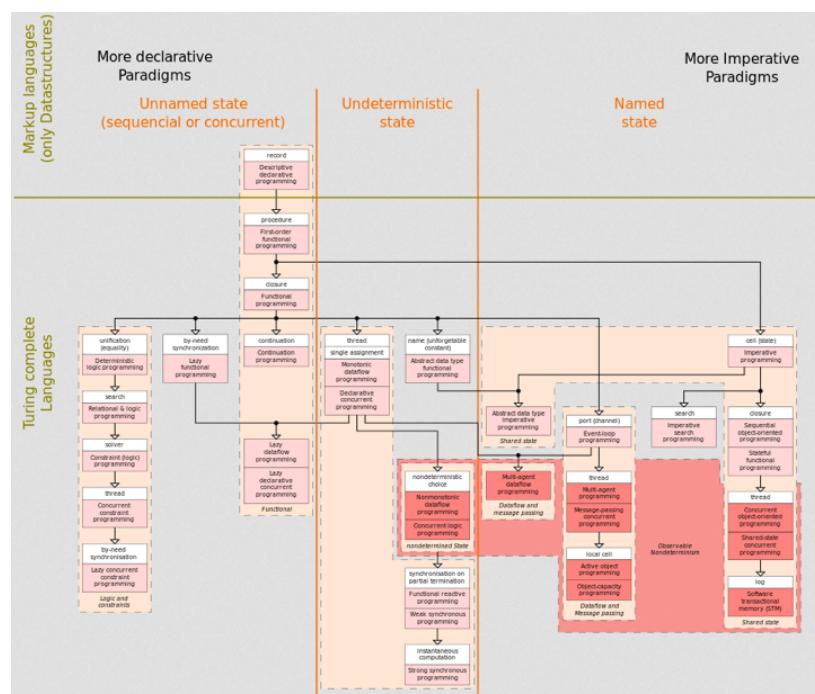


<http://polyglotprogramming.com/>

- Emerging Approach:

- Multi-Paradigm Design

1. Beyond Objects
2. Multiple Languages (GSL + DSL)
3. OOP + Functional + AOP



[http://en.wikipedia.org/wiki/File:Programming\\_paradigms.svg](http://en.wikipedia.org/wiki/File:Programming_paradigms.svg)

# Polyglot Persistence

- Current Approach
  - *Relational is King!*
  - 1. OLTP (*Transactions*)
  - 2. OLAP (*Analytics*)
  - 3. Centralized DB Server
  - 4. Scale-Up (*add expensive hardware*)
- Emerging Approach:
  - *Not Just Relational*
  - 1. SQL + NoSQL (*Eventual Consistency*)
  - 2. Hadoop + Columnar DB (*Analytics*)
  - 3. Distributed Data Stores
  - 4. Scale-Out (*add cheap computers*)

*SQL's dominance is cracking*

Relational databases are designed to run on a single machine, so to scale, you need buy a bigger machine



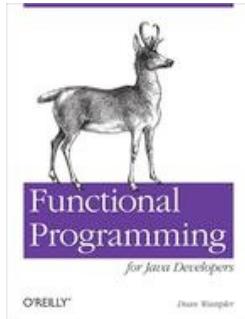
But it's cheaper and more effective to *scale horizontally* by buying lots of machines.



# Map & Reduce in Functional Programming

Two important concepts in functional programming:

1. **Map**: do something to everything in a list
2. **Fold**: combine results of a list in some way



## Combinator Functions: The Collection Power Tools

You already think of lists, maps, etc. as “collections,” all with a set of common methods. Most collections support Java `Iterators`, too. In functional programming, there are three core operations that are the basis of almost all work you do with collections:

### Filter

Create a new collection, keeping only the elements for which a filter method returns `true`. The size of the new collection will be less than or equal to the size of the original collection.

### Map

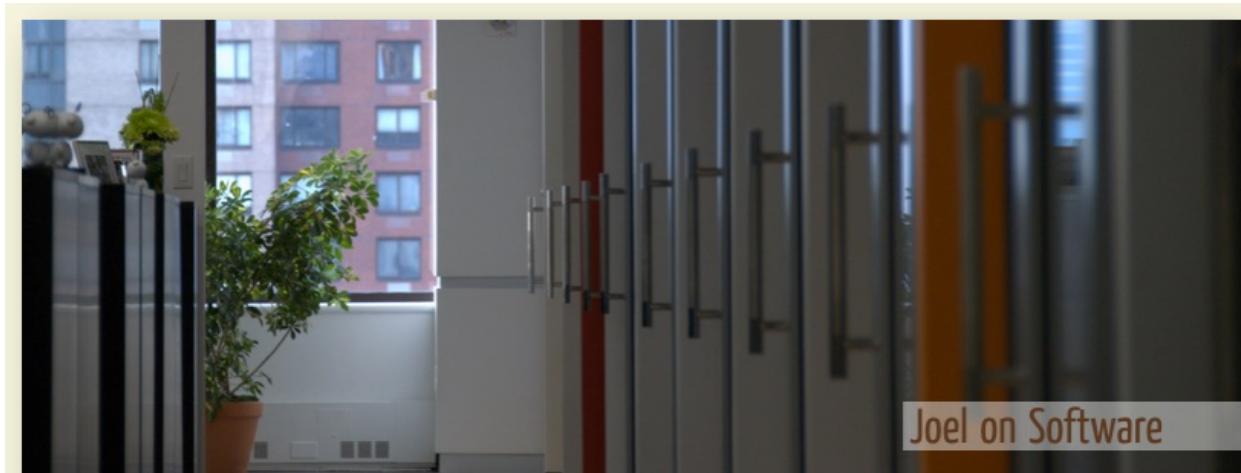
Create a new collection where each element from the original collection is transformed into a new value. Both the original collection and the new collection will have the same size. (Not to be confused with the `Map` data structure.)

### Fold

Starting with a “seed” value, traverse through the collection and use each element to build up a new final value where each element from the original collection “contributes” to the final value. An example is summing a list of integers.

Many other common operations can be built on top of these three. Together they are the basis for implementing concise and *composable* behaviors. Let’s see how.

<http://www.joelonsoftware.com/items/2006/08/01.html>



Joel on Software

Joel on Software

## Can Your Programming Language Do This?

by Joel Spolsky

Tuesday, August 01, 2006

One day, you're browsing through your code, and you notice two big blocks that look almost exactly the same. In fact, they're exactly the same, except that one block refers to "Spaghetti" and one block refers to "Chocolate Moose."

```
// A trivial example:  
  
alert("I'd like some Spaghetti!");  
alert("I'd like some Chocolate Moose!");
```

These examples happen to be in JavaScript, but even if you don't know JavaScript, you should be able to follow along.

The repeated code looks wrong, of course, so you create a function:

Wanted: [Software Engineer – Service Developer \(Video Service\)](#) at [Samsung R&D](#) (San Jose, CA).

See this and other great job listings on [the jobs page](#).

 CAREERS 2.0  
by stackoverflow

# SQL vs. Map Reduce

Below is a selection from the "Orders" table:

OrderID	CustomerID	EmployeeID	OrderDate	ShipperID
10248	90	5	1996-07-04	3
10249	81	6	1996-07-05	1
10250	34	4	1996-07-08	2

And a selection from the "Shippers" table:

ShipperID	ShipperName	Phone
1	Speedy Express	(503) 555-9831
2	United Package	(503) 555-3199
3	Federal Shipping	(503) 555-9931

And a selection from the "Employees" table:

EmployeeID	LastName	FirstName	BirthDate	Photo	Notes
1	Davolio	Nancy	1968-12-08	EmpID1.pic	Education includes a BA....
2	Fuller	Andrew	1952-02-19	EmpID2.pic	Andrew received his BTS....
3	Leverling	Janet	1963-08-30	EmpID3.pic	Janet has a BS degree....

# SQL vs. Map Reduce

SQL Statement: Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName,COUNT(Orders.OrderID) AS NumberOfOrders FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID  
GROUP BY ShipperName;
```

**Run SQL »**

Result:

Number of Records: 3

ShipperName	NumberOfOrders
Federal Shipping	68
Speedy Express	54
United Package	74

*How does a SQL  
Query Engine  
Process this Query?*

*How would this be  
done in Map  
Reduce?*

# Breaking Down Query Processing

SQL Statement: Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName, Orders.OrderID, Orders.CustomerID FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID
```

**Run SQL »**

Result:

Number of Records: 196

ShipperName	OrderID	CustomerID
Federal Shipping	10248	90
Speedy Express	10249	81
United Package	10250	34
Speedy Express	10251	84
United Package	10252	76
United Package	10253	34
United Package	10254	14
Federal Shipping	10255	68
United Package	10256	88
Federal Shipping	10257	35
Speedy Express	10258	20
Federal Shipping	10259	13
Speedy Express	10260	55
United Package	10261	61

SQL Statement: Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName, 1 FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID
```

**Run SQL »**

Result:

Number of Records: 196

ShipperName	Expr1001
Federal Shipping	1
Speedy Express	1
United Package	1
Speedy Express	1
United Package	1
United Package	1
United Package	1
Federal Shipping	1
United Package	1
Federal Shipping	1
Speedy Express	1
Federal Shipping	1
Speedy Express	1
United Package	1

# Breaking Down Query Processing

SQL Statement: Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName, SUM(1) FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID  
GROUP BY ShipperName;
```

[Run SQL »](#)

Result:

Number of Records: 3

ShipperName	Expr1001
Federal Shipping	68
Speedy Express	54
United Package	74

*These are Equivalent!*

SQL Statement: Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName,COUNT(Orders.OrderID) AS NumberOfOrders FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID  
GROUP BY ShipperName;
```

[Run SQL »](#)

Result:

Number of Records: 3

ShipperName	NumberOfOrders
Federal Shipping	68
Speedy Express	54
United Package	74

[http://www.w3schools.com/sql/sql\\_groupby.asp](http://www.w3schools.com/sql/sql_groupby.asp)

# What's Going on Here?

SQL Statement:

Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName, 1 FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID
```

Number of Records: 196

ShipperName	Expr1001
Federal Shipping	1
Speedy Express	1
United Package	1
Speedy Express	1
United Package	1
United Package	1
United Package	1
Federal Shipping	1
United Package	1
Federal Shipping	1
Speedy Express	1
Federal Shipping	1
Speedy Express	1
United Package	1

SQL Statement:

Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT Shippers.ShipperName, SUM(1) FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID  
GROUP BY ShipperName|
```

ShipperName

Expr1001

Federal Shipping

68

Speedy Express

54

United Package

74

*Relational  
Projections  
on Base Tables*

*In Memory  
Temp Table*

*Relational  
Aggregations*

*Result (Output)*

# The Map Reduce Way

```
OrderID,ShipperName,CustomerID,OrderDate  
10248, Federal Shipping, 90, 7/4/1996  
10249, Speedy Express, 81, 7/5/1996  
10250, United Package, 34, 7/8/1996
```

**Map Input**  
*(Key = OrderID,  
Values = Line)*

Number of Records: 196	
ShipperName	Expr1001
Federal Shipping	1
Speedy Express	1
United Package	1
Speedy Express	1
United Package	1
Federal Shipping	1
United Package	1
Federal Shipping	1
Speedy Express	1
Federal Shipping	1
Speedy Express	1
United Package	1

**Intermediate Results**  
*(Key = Shipper, Value = 1)*

**Reduce Intermediate Results**  
*Add all the 1's for each Shipper*  
*(Key = Shipper, Value = Sum of  
Ones')*

ShipperName	Expr1001
Federal Shipping	68
Speedy Express	54
United Package	74

**Result (Output)**

# Which is Easier for Developers?

*Write SQL Query*



SQL Statement:

```
SELECT Shippers.ShipperName, SUM(1) FROM Orders  
LEFT JOIN Shippers  
ON Orders.ShipperID=Shippers.ShipperID  
GROUP BY ShipperName|
```

Run SQL >

Result:

ShipperName	Expr1001
Federal Shipping	68
Speedy Express	54
United Package	74

```
OrderID,ShipperName,CustomerID,OrderDate  
10248, Federal Shipping, 90, 7/4/1996  
10249, Speedy Express, 81, 7/5/1996  
10250, United Package, 34, 7/8/1996
```



*Write Map Function (Java?)*

Number of Records: 196	
ShipperName	Expr1001
Federal Shipping	1
Speedy Express	1
United Package	1
Speedy Express	1
United Package	1
United Package	1
United Package	1
Federal Shipping	1
United Package	1
Federal Shipping	1
Speedy Express	1
Federal Shipping	1
Speedy Express	1
United Package	1



*Get Results!*

ShipperName	Expr1001
Federal Shipping	68
Speedy Express	54
United Package	74



*Write Reduce Function (Java?)*

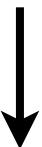
# Which can Scale?

Terabytes of Data across  
Thousands of Servers



OrderID,ShipperName,CustomerID,OrderDate  
10248, Federal Shipping, 90, 7/4/1996  
10249, Speedy Express, 81, 7/5/1996  
10250, United Package, 34, 7/8/1996

1 Million Rows



SQL

SQL Statement:

```
SELECT Shippers.ShipperName, SUM(1) FROM Orders
LEFT JOIN Shippers
ON Orders.ShipperID=Shippers.ShipperID
GROUP BY ShipperName|
```

Run SQL »

Result:

ShipperName	Expr1001
Federal Shipping	68
Speedy Express	54
United Package	74



Map

Number of Records: 196	
ShipperName	Expr1001
Federal Shipping	1
Speedy Express	1
United Package	1
Speedy Express	1
United Package	1
United Package	1
United Package	1
Federal Shipping	1
United Package	1
Federal Shipping	1
Speedy Express	1
Federal Shipping	1
Speedy Express	1
United Package	1



Reduce

ShipperName	Expr1001
Federal Shipping	68
Speedy Express	54
United Package	74

# Examples from Paper

```
map(String key, String value):
    // key: document name
    // value: document contents
    for each word w in value:
        EmitIntermediate(w, "1");

reduce(String key, Iterator values):
    // key: a word
    // values: a list of counts
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```

**Inverted Index:** The map function parses each document, and emits a sequence of  $\langle \text{word}, \text{document ID} \rangle$  pairs. The reduce function accepts all pairs for a given word, sorts the corresponding document IDs and emits a  $\langle \text{word}, \text{list(document ID)} \rangle$  pair. The set of all output pairs forms a simple inverted index. It is easy to augment this computation to keep track of word positions.

**Distributed Sort:** The map function extracts the key from each record, and emits a  $\langle \text{key}, \text{record} \rangle$  pair. The reduce function emits all pairs unchanged. This computation depends on the partitioning facilities described in Section 4.1 and the ordering properties described in Section 4.2.

**Distributed Grep:** The map function emits a line if it matches a supplied pattern. The reduce function is an identity function that just copies the supplied intermediate data to the output.

**Count of URL Access Frequency:** The map function processes logs of web page requests and outputs  $\langle \text{URL}, 1 \rangle$ . The reduce function adds together all values for the same URL and emits a  $\langle \text{URL}, \text{total count} \rangle$  pair.

**Reverse Web-Link Graph:** The map function outputs  $\langle \text{target}, \text{source} \rangle$  pairs for each link to a target URL found in a page named source. The reduce function concatenates the list of all source URLs associated with a given target URL and emits the pair:  $\langle \text{target}, \text{list(source)} \rangle$

**Term-Vector per Host:** A term vector summarizes the most important words that occur in a document or a set of documents as a list of  $\langle \text{word}, \text{frequency} \rangle$  pairs. The map function emits a  $\langle \text{hostname}, \text{term vector} \rangle$  pair for each input document (where the hostname is extracted from the URL of the document). The reduce function is passed all per-document term vectors for a given host. It adds these term vectors together, throwing away infrequent terms, and then emits a final  $\langle \text{hostname}, \text{term vector} \rangle$  pair.

# Map Reduce Execution Model

