

From RDBMS to NoSQL



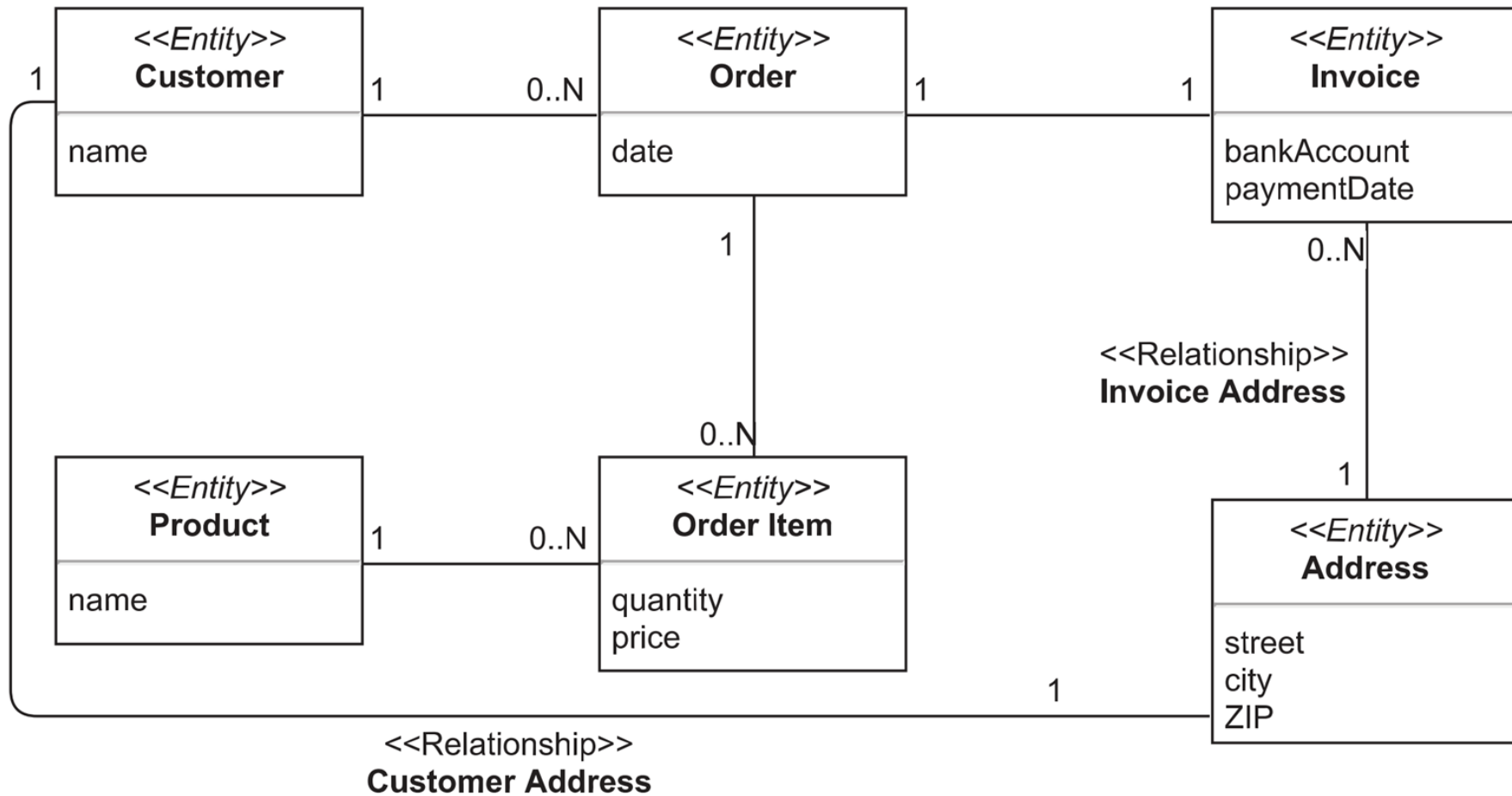
Efficient implementations of table **joins** and of **transactional** processing **require centralized** system.

NoSQL Databases:

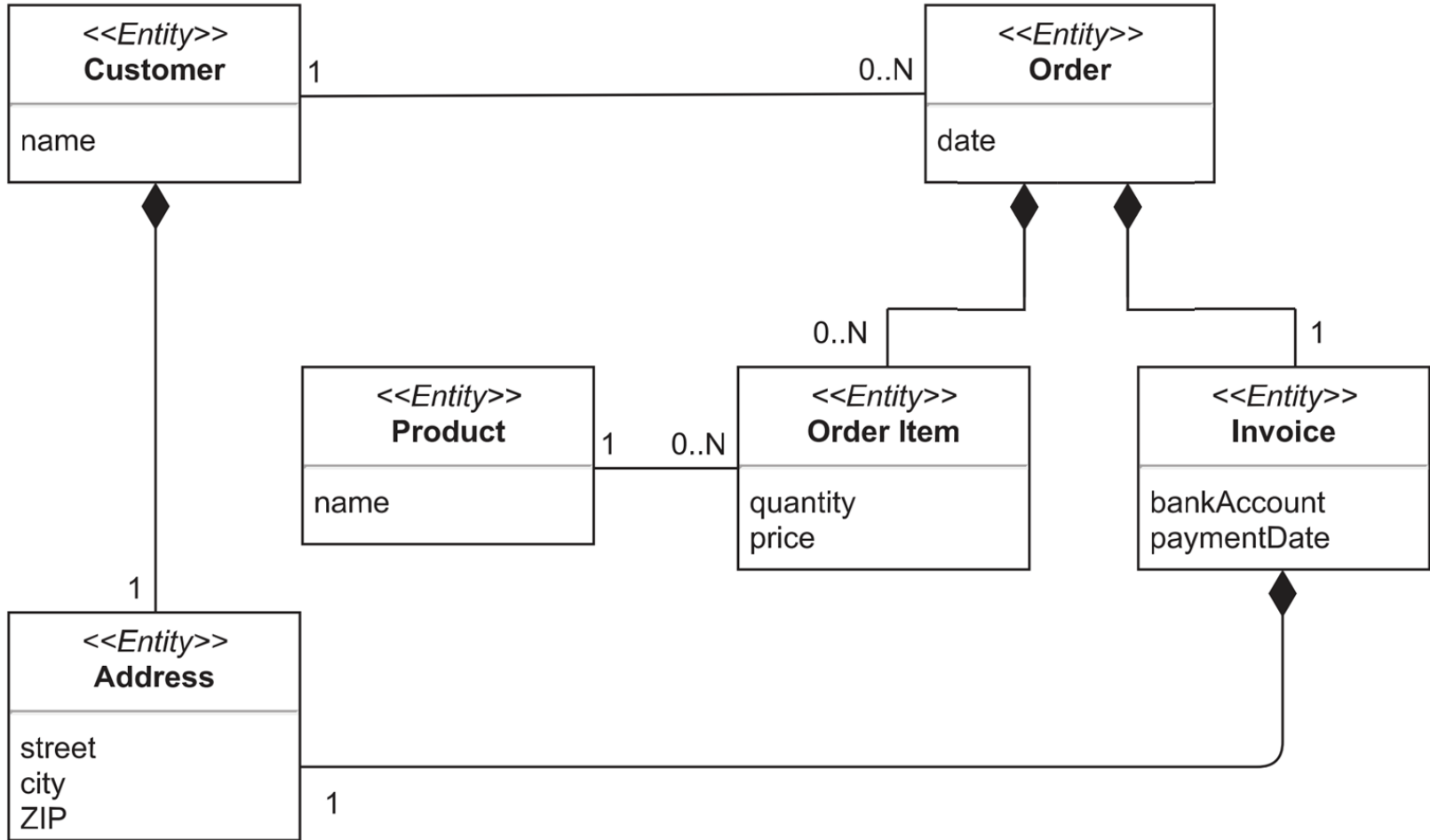
- Database **schema** tailored for a **specific application**
 - **keep together** data pieces that are often accessed together
- Write operations might be slower but **read is fast**
- **Weaker consistency** guarantees

=> **efficiency** and horizontal **scalability**

Example (1): UML Model



Example (3): Aggregates (in UML)



NoSQL Databases: Aggregate-oriented



Many NoSQL stores are aggregate-oriented:

- There is no general strategy to set aggregate boundaries
- Aggregates give the database information about which bits of data will be manipulated together
 - What should be stored on the same node
- Minimize the number of nodes accessed during a search
- Impact on concurrency control:
 - NoSQL databases typically support atomic manipulation of a single aggregate at a time

Agenda



- Fundamentals of RDBMs and NoSQL Databases
- Data Model of Aggregates
- Models of **Data Distribution**
 - scalability
 - **sharding** vs. **replication**: master-slave, peer-to-peer
 - combination
- Consistency
 - write-write vs. read-write conflict
 - strategies and techniques
 - relaxing consistency

Vertical Scalability (Scaling up)



- Involve **larger** and more **powerful** machines
 - large disk storage using **disk arrays**
 - massively **parallel** architectures
 - large **main memories**
- Traditional choice
 - in favour of **strong consistency**
 - very **simple** to realize (**no** handling of data **distribution**)
- Works in many cases **but...**

Vertical Scalability: Drawbacks



- Higher **costs**
 - Large machines **cost more** than equivalent commodity HW
- Data growth **limit**
 - Large machine works well until the data grows to fill it
 - Even the **largest** of machines **has a limit**
- **Proactive** provisioning
 - In the beginning, **no idea** of the **final scale** of the application
 - An **upfront budget** is needed when scaling vertically
- **Vendor** lock-in
 - Large machines are produced by **a few vendors**
 - Customer is **dependent on** a single vendor (proprietary HW)

Horizontal Scalability (Scaling out)



System is distributed across **multiple machines/nodes**

- **Commodity** machines, cost effective
- Provides **higher scalability** than vertical approach
 - Data is partitioned over **many disks**
 - Application can use **main memory** of all machines
 - **Distribution computational** model
- Introduces new **problems**:
 - synchronization, consistency, handling partial-failures, etc.



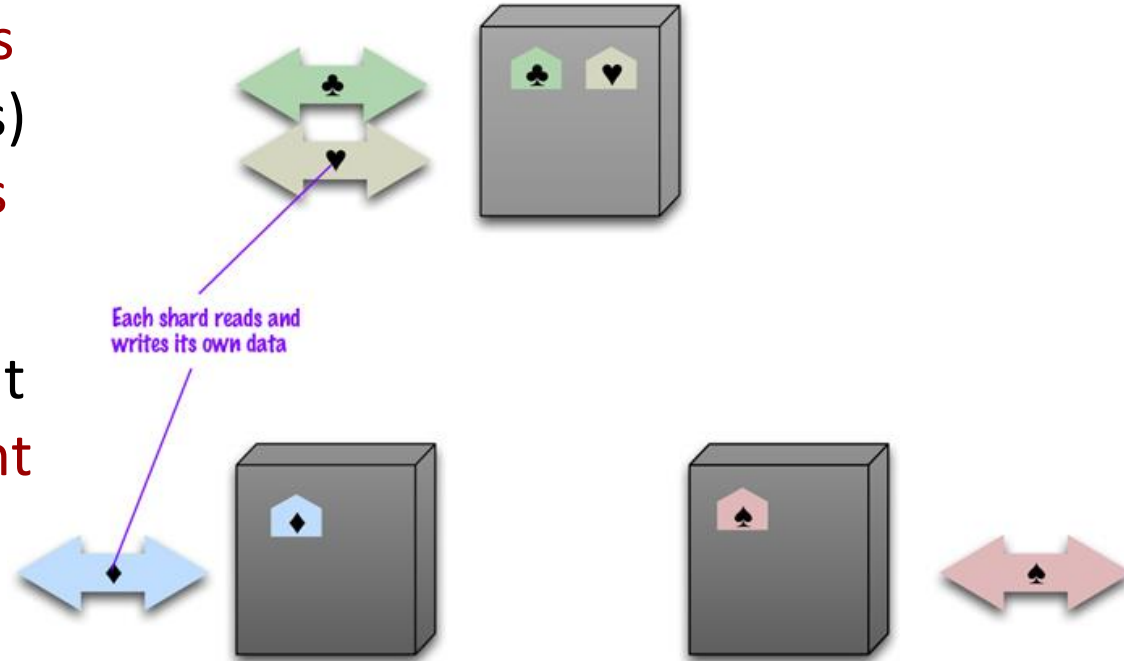
Distribution Models: Overview

- for horizontal scalability
- **Two generic ways of data distribution:**
 - **Replication** – the same data is copied over multiple nodes
 - Master-slave vs. peer-to-peer
 - **Sharding** – different data chunks are put on different nodes (data partitioning)
 - Master-master
- We can use either or **combine them**
 - **Distribution models** = specific ways to do **sharding**, **replication** or combination of both

Sharding (Data Partitioning)



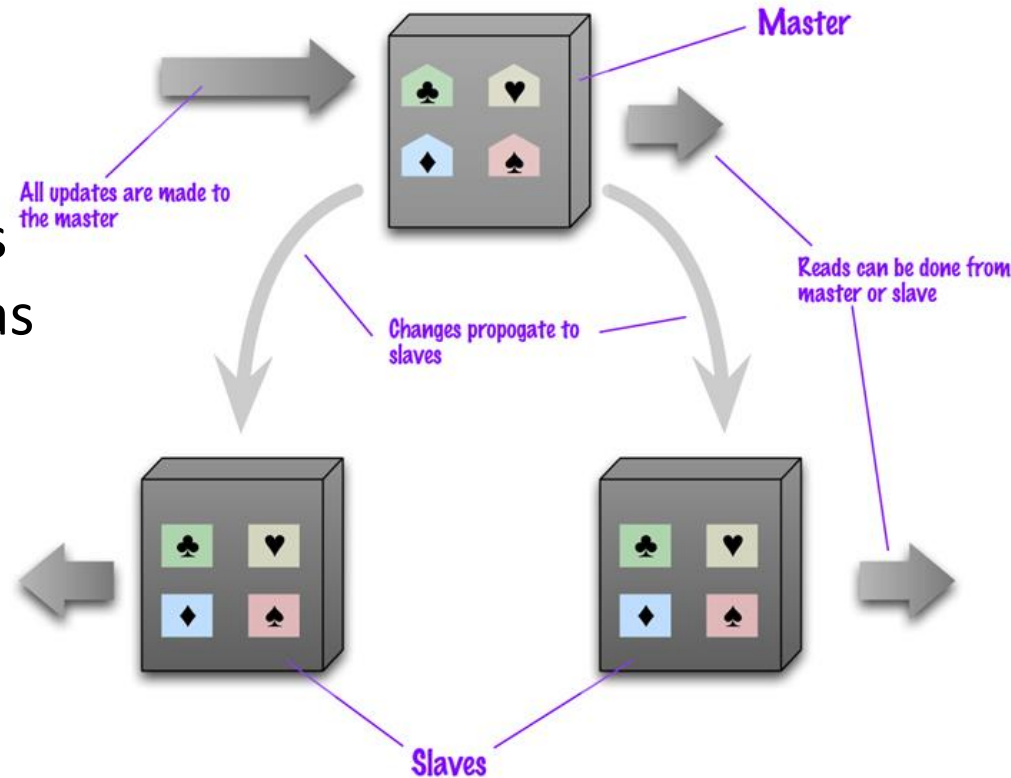
- Placing **different parts** of the data (card suits) onto **different servers**
- Applicability: Different clients **access different** parts of the dataset



Master-slave Replication



- We replicate data across multiple nodes
- One node is designated as primary (**master**), others as secondary (**slaves**)
- **Master** is responsible for processing **all updates** to the data
- **Reads** from **any** node



Master-slave Replication (2)



- For scaling a **read-intensive** application
 - More read requests → more slave nodes
 - The **master fails** → the slaves can still **handle read** requests
 - A slave can become a new master quickly (it is a replica)
- **Limited** by ability of the master to process updates
- Masters are **selected** manually or automatically
 - User-defined vs. cluster-elected

Peer-to-peer Replication (2)



- **Problem:** consistency
 - Users can **write simultaneously** at two different nodes
- **Solution:**
 - When writing, the **peers coordinate** to avoid conflict
 - At the cost of **network traffic**
 - The **write** operation **waits** till the coordination process is finished
 - Not all replicas need to **agree on** the write, just a majority (details below)

Consistency in Databases



- “Consistency is the lack of contradiction in the DB”
- Centralized RDBMS ensure **strong consistency**
- Distributed NoSQL databases typically **relax consistency** (and/or durability)
 - Strong consistency → **eventual** consistency
 - BASE (basically available, soft state, eventual consistency)
 - **CAP** theorem
 - **tradeoff** between consistency and availability

Transaction Processing in NoSQL



- Basically, no problem if the DB is **centralized**
 - ACID can be implemented
 - Various **levels of isolation** (details later in the course)
 - read uncommitted
 - read committed
 - repeatable reads
 - serializable
- **Distributed transactions** (details later in the course)
 - X/Open Distributed Trans. Processing Model (X/Open XA)
 - Two-phase Commit Protocol (**2PC**)
 - Strong Strict Two-phase Locking (SS2PL)

CAP Theorem: Formulation

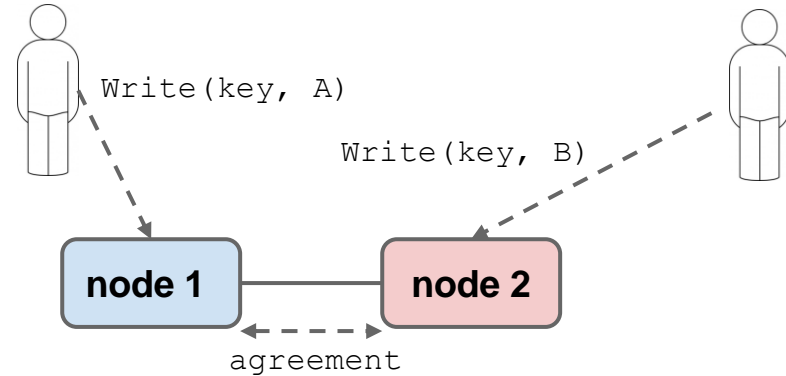


- **CAP Theorem:** A “shared-data” system **cannot** have **all three** CAP properties
 - Or: only **two of** the **three** CAP properties are possible
 - This is the common version of the theorem
- First **formulated** in 2000: prof. Eric Brewer
 - ACM PODC Conference Keynote speech
 - www.cs.berkeley.edu/~brewer/cs262b-2004/PODC-keynote.pdf
- **Proven** in 2002: Seth Gilbert & Nancy Lynch
 - SIGACT News 33(2) <http://dl.acm.org/citation.cfm?id=564601>

PC: Partition Tolerance & Consistency

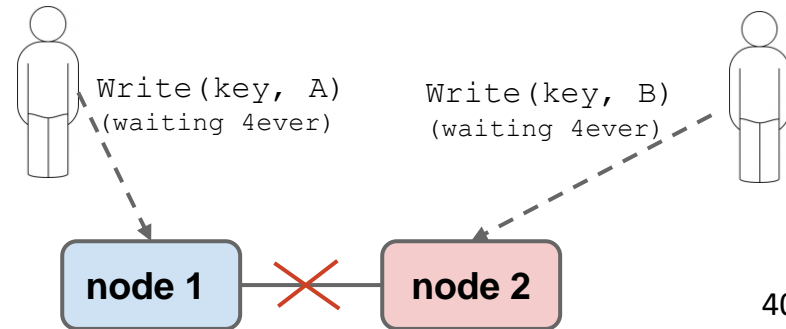


Example: two **users**, two **masters**, two **write** attempts



- **Strong** consistency:
 - Before the write is committed, **both** nodes have to **agree** on the order of the writes

- If the nodes are **partitioned**, we are **losing Availability**
 - (but reads are still available)



Summary of the Lesson



- **Aggregate**-oriented data modeling
- **Sharding vs. replication**
 - **Master-slave vs. peer-to-peer** replication
 - Combination of sharding & replication
- Database **consistency**:
 - **Write/Read** consistency (write-write & write-read **conflict**)
 - **Replication** consistency (also, read-your-own-writes)
- **Relaxing** consistency:
 - CAP (**C**onsistency, **A**vailability, Tolerance to **P**artitions),
 - Eventual consistency
 - **Quora** (write/read quorum)
 - **can** ensure **strong** replication consistency; wide range of settings

References



- I. Holubová, J. Kosek, K. Minařík, D. Novák. Big Data a NoSQL databáze. Praha: Grada Publishing, 2015. 288 p.
- Sadalage, P. J., & Fowler, M. (2012). NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence. Addison-Wesley Professional, 192 p.
- doc. RNDr. Irena Holubova, Ph.D. MMF UK course NDBI040: Big Data Management and NoSQL Databases
- Eric Brewer: Towards Robust Distributed Systems.
www.cs.berkeley.edu/~brewer/cs262b-2004/PODC-keynote.pdf