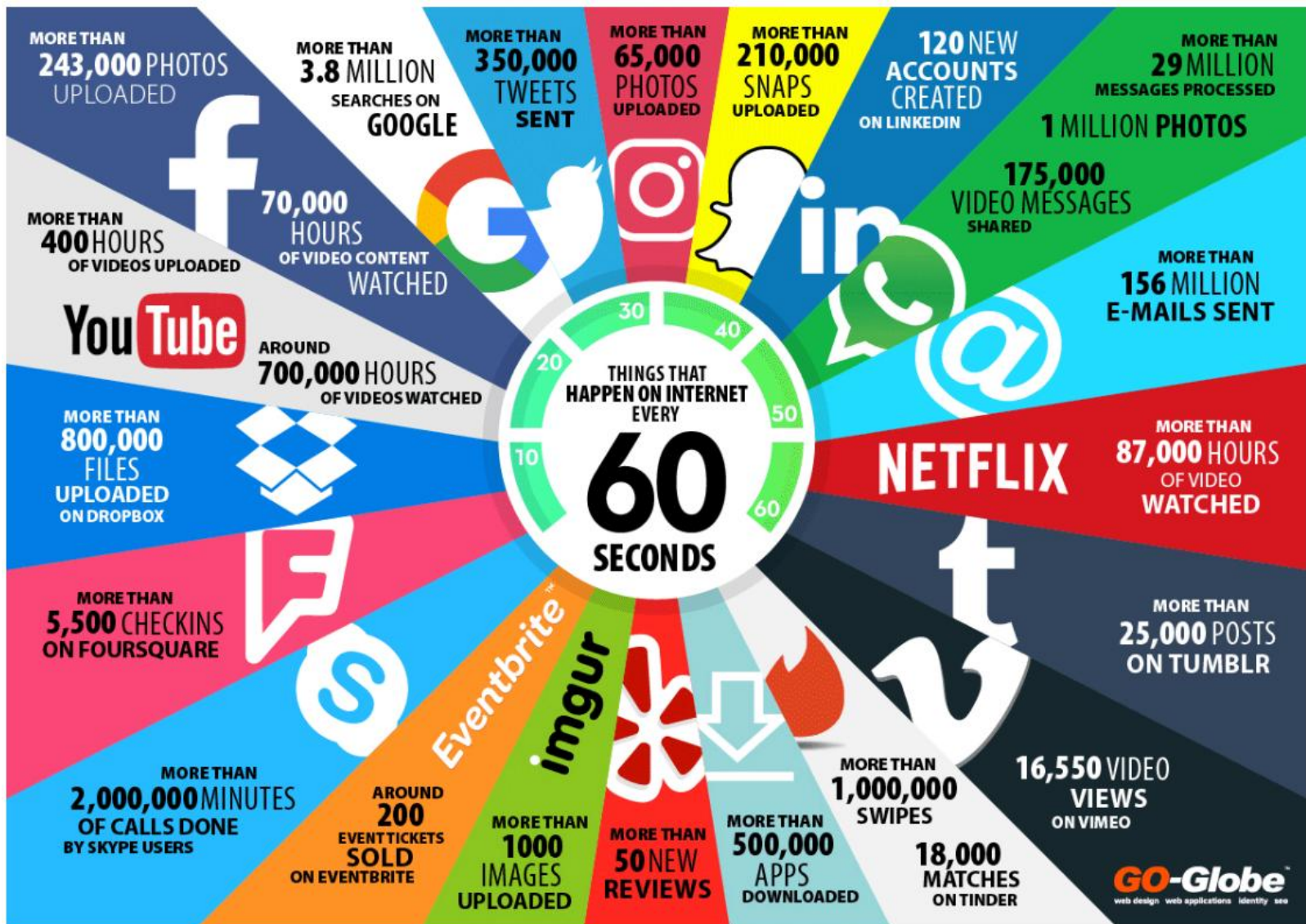


NoSQL

CS460 Databases for Data Scientists



Big Data

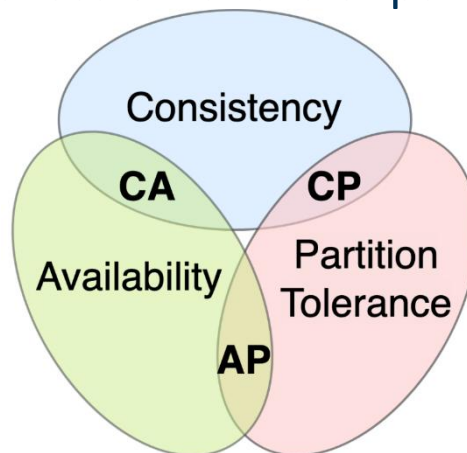


Motivation for NoSQL

- Scalability
- RDBMS usually scale with hardware added to one server
 - very expensive when reaching maximum capacity
- On > 50GB data:
 - MySQL
 - Writes 300 ms avg
 - Reads 350 ms avg
 - Cassandra
 - Writes 0.12 ms avg
 - Reads 15 ms avg
- NoSQL stands for
 - No Relational
 - Not only SQL
 - An umbrella term for a class of products that don't follow RDBMS principles

CAP Theorem

- “Of three properties of shared-data systems — data Consistency, system Availability, and tolerance to network Partitions — only two can be achieved at any given moment in time” — Brewer, 1999
 - **Consistency**: all nodes see the same data at the same time
 - **Availability**: guarantee that every request receives a response about whether it was successful or failed
 - **Partition tolerance**: system continues to operate despite arbitrary message loss or failure of part of the system



NoSQL Categories

- Four groups:
 - Key-value stores
 - Column-based families or wide column systems
 - Document stores
 - Graph databases

Key-value stores

- Simplest NoSQL databases
 - collection of key, value pairs
- Queries are limited to query by key
- Examples:

- Redis
- Riak
- Voldermort
- DynamoDB
- MemcacheDB

Key	Value
K1	AAA,BBB,CCC
K2	AAA,BBB
K3	AAA,DDD
K4	AAA,2,01/01/2015
K5	3,ZZZ,5623

Column-Based Families

- Data is stored in a big table
- Access control, disk storage, and memory accounting performed on column families

- Example:

- Cassandra
- Hbase
- DuckDB
- Hypertable

id	user_name	current_balance	number_of_transactions
1	'taggy@yahoo.com'	1059298	1045
2	'fra@hotmail.com'	3910	194
...

Column storage better

```
SELECT sum(current_balance)
      AS total_transactions
FROM table
WHERE user_id > 2
```

Row storage better

```
SELECT user_id, user_name, current_balance
FROM table
WHERE user_id = 1
```

Document Databases

- Collections of similar documents
- Each document can resemble a complex model
- Good for single view or data hub applications

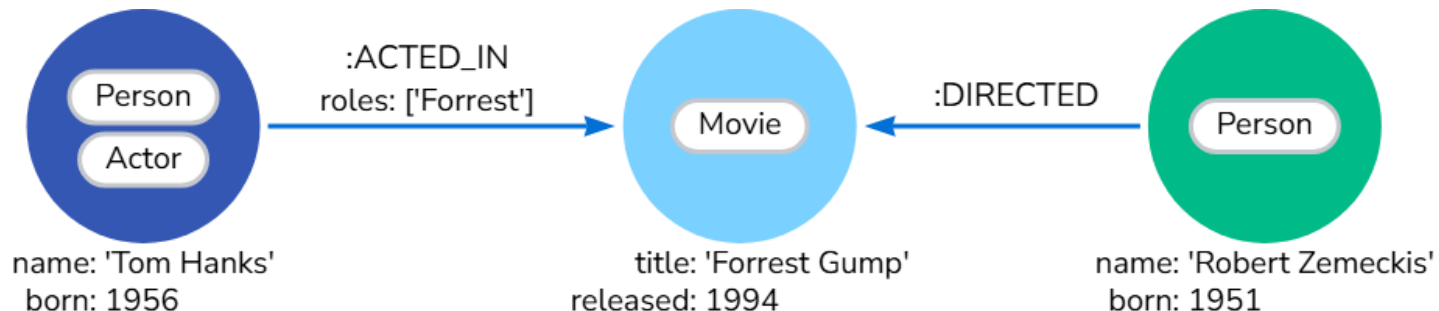
- Examples:

- MongoDB
- ArangoDB
- CouchDB

```
{  
  "firstName": "Bob",  
  "lastName": "Smith",  
  "address": {  
    "type": "Home",  
    "street1": "5 Oak St.",  
    "city": "Boys",  
    "state": "AR",  
    "zip": "32225",  
    "country": "US"  
  },  
  "hobby": "sailing",  
  "phone": {  
    "type": "Cell",  
    "number": "(555)-123-4567"  
  }  
}
```

Graph Database

- Collection of vertices (nodes) and edges (relations) and their properties
- Example:
 - Neo4j
 - Amazon Neptune
 - AllegroGraph
 - VertexDB



What's Next?

- Database Systems II (FSS, Moerkotte)
 - e.g., distributed DBMS, object-relational DBs, deductive DBs
- Query Optimization (FSS, Moerkotte)
 - more sophisticated query optimization
- Large-Scale Data Management (HWS, Gemulla)
 - e.g., parallel & distributed databases, MapReduce, SPARQL, NoSQL
- Knowledge Graphs (Hertling)
 - Graph Representation and Inference
 - Knowledge Modeling and Integration

What's Next?

- Data Security and Privacy (FSS, Armknecht)
 - also covers aspects such as encryption
- Web Data Integration (HWS, Bizer)
 - dealing with multiple databases
 - automatically integrating them into a single one
 - can be accompanied with a practical project
- Data Mining (FSS/HWS, Bizer/Hertling)
 - finding patterns in data
 - entry point to more specific lectures in the data analytics field
 - includes a practical project

Exam

- Date: Monday, 01.06.2026
- Time: 12:00-13:00
- Location:
 - Room: B 243
 - Building: A 5, 6 Part B

Exam

- 6 questions
 - Each question 10 points
- Prepare for similar questions as the exercise
- Focus:
 - Be able to apply your knowledge to a specific database / example at hand

Exam

- SQL
 - Be able to write SQL queries (know the syntax)
 - **No** views, no roles
 - **No** temporary relations using WITH
- ER Models
 - Write and understand ER models
 - Be able to reduce a ER model to Relation Schema
- Normalization
 - Be able to normalize a given database

Exam

- Indexing/Hashing
 - Understand differences between different indices
 - Create and modify given indices
 - **Excluding** special cases like „merge siblings“ or „redistribute pointers“ for B+ tree
- Database Architectures
 - General understanding of concepts (no details)

Exam

- Query Processing
 - Know the different join possibilities
 - Know how indices can be used for selections
- Query Optimization
 - Be able to apply equivalence rules (don't need to remember them)
 - Understand execution/evaluation plans
 - Be able to decide which one is better given a concrete example
 - Be able to compute the number of tuples for a given selection

Exam

- Transactions & Concurrency
 - Understand and Write schedules
 - Know what Conflicts are (detect conflicts)
 - Know the Two-Phase Locking Protocol
 - Know/detect deadlocks
 - **No** Timestamp-based Scheduling or Validation Based Protocol
- Recovery
 - Know and apply log-based recovery mechanism (undo, redo operations)
 - Including checkpoints

Exam

- Applications
 - Know the security implications
 - SQL Injection
 - Cross Site Scripting
 - Password storage
 - **No** HTML or servlet code etc
- NoSQL
 - **Not** exam relevant

Course Evaluation

- See Message Board in Ilias

Q & A

