# IN3020/4020 – Database Systems
# Spring 2021, Week 1.1

# Introduction

Dr. M. Naci Akkøk, Assoc. Prof., Ifi/UiO and CEET/OsloMet, CEO In-Virtualis AS

Dr. Egor Kostylev, Assoc. Prof., Ifi/UiO

Based upon slides by M. Naci Akkøk, Spring 2020 & E. Thorstensen from Spring 2019

UiO **: Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# The course

o Database systems, databases and their usage

- o **Databases** are abstract objects (an entity or model with **data structures** and **a query language**) for storing and retrieving **data**

- o **Database systems** (also referred to as Database Management Systems or **DBMS**) are software systems that allow us to manage these objects, i.e., the databases: They are the "engines" that help manage/run databases.

# Goals:
# At the end of the course, you should be able to…

o Understand what Database Management Systems (DBMS) do, how they work and how they are implemented

o Have a good grasp of the theoretical & practical aspects of database administration

o Have a good grasp of use of SQL and how queries are optimized

o Understand ACID principles & the principles of transaction management

o Know about the TX management types (isolation levels) that modern DBMS offer, and their strengths/weaknesses

UiO **: Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# Inside the database "black box"

o What is apparent to us is that we write an SQL query, the database does something and comes up with a response

o But there is a lot that happens behind the scenes, inside the black box.

o This course is meant to show what happens behind the scenes in the black box, so that **we understand the common challenges** and what to do about them as users and database administrators and understand **how to build database engines** some day.

# Some common challenges

o What can one do if the queries are slow?

o How can one read from and write to the same database at the same time without problems? How would the performance be affected?

o What happens if the database crashes while updating the data?

# Some of the background we will need to build up for the course

o DBMS architecture

o DBMS file systems

o Various data structures and algorithms

o Relational algebra

# The syllabus can be divided into three parts

o Main focus:

   o Relational queries, SQL and query optimization
   o ACID, DBMS characteristics and mechanisms
   o Other DBMS (slightly less than previous years)


o And an additional bit about emerging technologies and database research

UiO **:** **Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# Syllabus: Queries, SQL and query optimization (E. Kostylev)

- o SQL (repetition and a bit more)
- o Query formalization/optimization (including relational algebra), query plans
- o Indexes, index usage and the underlying search

# Syllabus: ACID, DBMS characteristics and Mechanisms (M. N. Akkøk)

o Managing multiple concurrent uses/users and protecting the data
- o ACID – what it stands for

- o Transaction Management
- o Locking, logging, buffer & cache management

- o Relationships between these mechanisms
- o Synchronization/replication challenges

UiO **: Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# Syllabus: Other database systems (briefly) – NoSQL DBMS and newer DBMS technologies

o Categorization of the NoSQL DBMS and the problems they address (their "reason of being")

o A bit more about spatial & graph databases, semantic technologies and databases with built-in intelligence (AI/ML capabilities in and outside the "black box")

o Autonomous databases

o Multimodal databases

# Not really syllabus, but food for thought: Emerging DBMS trends and research areas

o   Performance – how and for what purpose, exactly?

o   Big or fast data – or both? How? How can newer architectures like the Data Mesh help?

o   What is the relation doing in the table? What happens to indexes and keys and search (queries, joins etc.) if it is not part of the data? Can it be as cleanly formalized in an "extended" relational algebra?

o   Intelligence in databases – autonomy, auto-inference (also in RDF semantics) and other built-in mechanisms of learning and intelligence.

o   Extraction-transformation-loading challenges: Intelligence in data quality management and I/O

o   Knowledge Bases: extend data with knowledge, usually in logical form, and efficiently manipulate (e.g., querying) both the explicit data and implicit data that logically follows from the knowledge

UiO **: Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# Dates & practical information

o Two mandatory exercises (*oblig´er*) due Friday 5. March and Friday 30. April 2021 (do follow updates on dates)

> o **Group sessions** dedicated to the mandatory exercises are **the week before the deadlines**

> o Make sure you do NOT co-operate too tightly! Check the rules for the mandatory exercises (*obligreglement*)

o The final written exam is a 4 hours **Inspera Digital Exam** without aids (*uten hjelpemidler*) on the 16th of June 2021 from 09:00 to 13:00

o **Deadline for withdrawing from the exam is 1st of May!**

o ALWAYS CHECK DATES, because things do change ☺

# Dates & practical information

o We may have to be on Zoom throughout the semester but check the announcements on the course´s web site.

o Groups as well.

o Yes, we will record the lectures and put them out in the "Schedule", and we will do our best to publish the slides in advance ☺

o We believe we will use Padlet (<u>not </u>Piazza) as the communication service. We will provide info and link on the Web-site of the course very soon
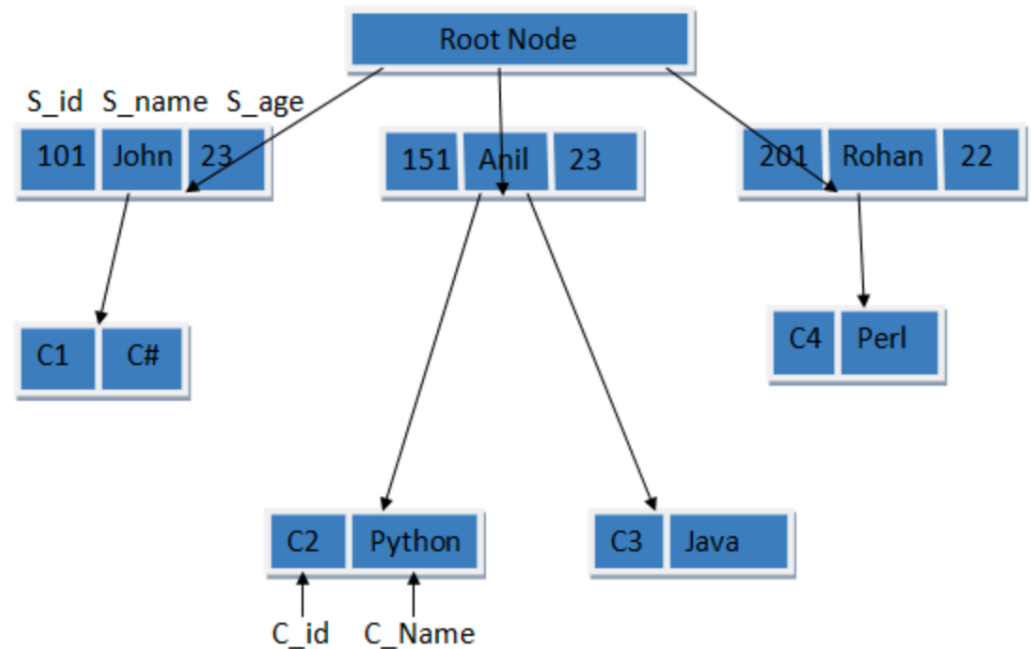
# Brief Database history (look them up)

o 1960-ies:
  o IDS (network DB model) and IMS (hierarchical DB model)
  o Mainframe period. 1 «big» machine for 1 mill. 1964-US$ with 512 Kbyte RAM
  o 50 Mbyte disk
  o Other I/O-equipment magnetic & paper tape, punched-card reader and teletype (which is why unix has a funny abbreviation like tty for the «terminal»)
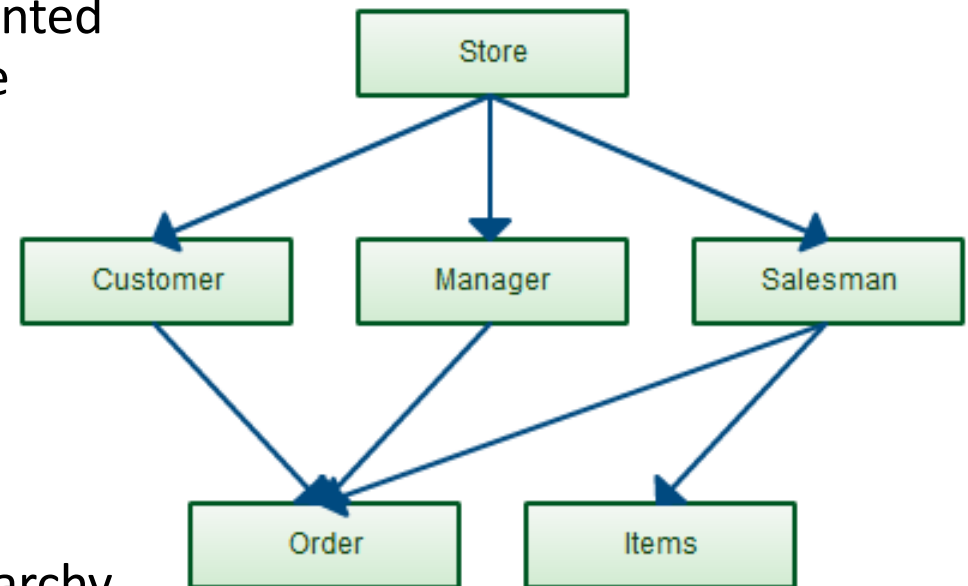o Network and hierarchical databases are somewhat related.

# The Hierarchical Database Model

- o «In the **hierarchical database model** data is represented in a tree-like structure. It represents a parent-child relationship with a single parent for each child» (but allows for more than one child per parent)

- o Own query language.

- o Advantages / disadvantages?

# The Network Database Model

o The **Network Model** is an extension of the hierarchical model. In this data is represented in the form of graphs with more than one parent node for one child node.

o Advantages:
  o Easy to implement.
  o Datamodel follows physical model.
  o Easy to learn?

o Disadvantages:
  o Easy only when query follows the hierarchy.
  o Updates?



UiO : **Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# Relational DB/DBMS (The Relational Model)

o E. F. Codd. 1970. A relational model of data for large shared data banks. Commun. ACM 13, 6 (June 1970), 377-387.

o "Future users of large data banks must be protected from having to know how the data is organized in the machine (the internal representation). A prompting service which supplies such information is not a satisfactory solution. Activities of users at terminals and most application programs should remain unaffected when the internal representation of data is changed and even when some aspects of the external representation are changed. Changes in data representation will often be needed as a result of changes in query, update, and report traffic and natural growth in the types of stored information."

UiO **: Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# What did Codd mean?

o Internal representation (files, pointers) should not be visible to the user

o Users work on an abstract data model

o Codd suggests relations as an abstract data structure

# Relations (the table version)

o Tables with rows (tuples) and columns (attributes).

o A tuple has a value for each attribute

o Table is based upon a variant of the mathematical relation: Set of tuples $R \subseteq d(A1) \times \cdots \times d(Ak)$.

o Note: «Contrary to the usual definition in *mathematics*, there is no ordering to the elements of the tuples of *a relation*!»

o Can be manipulated through operations on the tables and comparison of values

## Mathematical Relations

- Given sets $D_1$, $D_2$, …, $D_n$, not necessarily distinct.
- The *Cartesian product* $D_1 x D_2 x … x D_n$ is the set of all (ordered) n-tuples $<d_1, d_2, …, dn>$ such that $d_1 \in D_1$, $d_2 \in D_2$, …, $d_n \in D_n$
- A *mathematical relation* on $D_1$, $D_2$, …, $D_n$ is a subset of the Cartesian product $D_1 x D_2 x … x D_n$.
- $D_1$, $D_2$, …, $D_n$ are the *domains* of the relation.
- n is the *degree* of the relation.
- The number of n-tuples in a given relation is the *cardinality* of that relation; the cardinality of a relation is always finite.

CSC343 Introduction to Databases  — University of Toronto          The Relational Model — 3

Simple definitions of mathematical relation, function, graph etc. with examples:

https://medium.com/brandons-computer-science-notes/mathematical-relations-5416f027dbbe

Note: In mathematics, a *tuple* is a finite ordered list (sequence) of elements.

# Characteristics of Relations

o Every attribute has its own domain (i.e., set of allowed values)

o All attributes have different (distinct) names. If not, we distinguish them (for example) by giving them sequence numbers

o Relations are sets: Tuples have no order, and no two are alike

o «A table cannot have 0 columns, rows can be duplicated»

o «… and tables cannot be nested»
(careful with newer object-relational databases and databases that allow for abstract data types or ADTs)!

o Attributtes can have the special value of nil (null)

# Relational DB/DBMS (the Relational Model)

o Relation is a an abstract data model. Implementation varies

o We allow for defining redundancy and other characteristics

o We have a very elegant and compact query language (SQL) that builds upon relational algebra

o We allow for ad-hoc queries

o It is considered to be a universal data model

UiO **:** **Institutt for informatikk**
    Det matematisk-naturvitenskapelige fakultet

# Fast forward …

o The relational model and relational databases (DBMS) took over eventually

o Today, everybody seems to talk about Big Data in addition.

  When is Data Big? Or Fast?

**In 2019**, the following server could be rented for NOK900/month:
- Cores: 10 x 2.20 GHz (Single 10 Core)
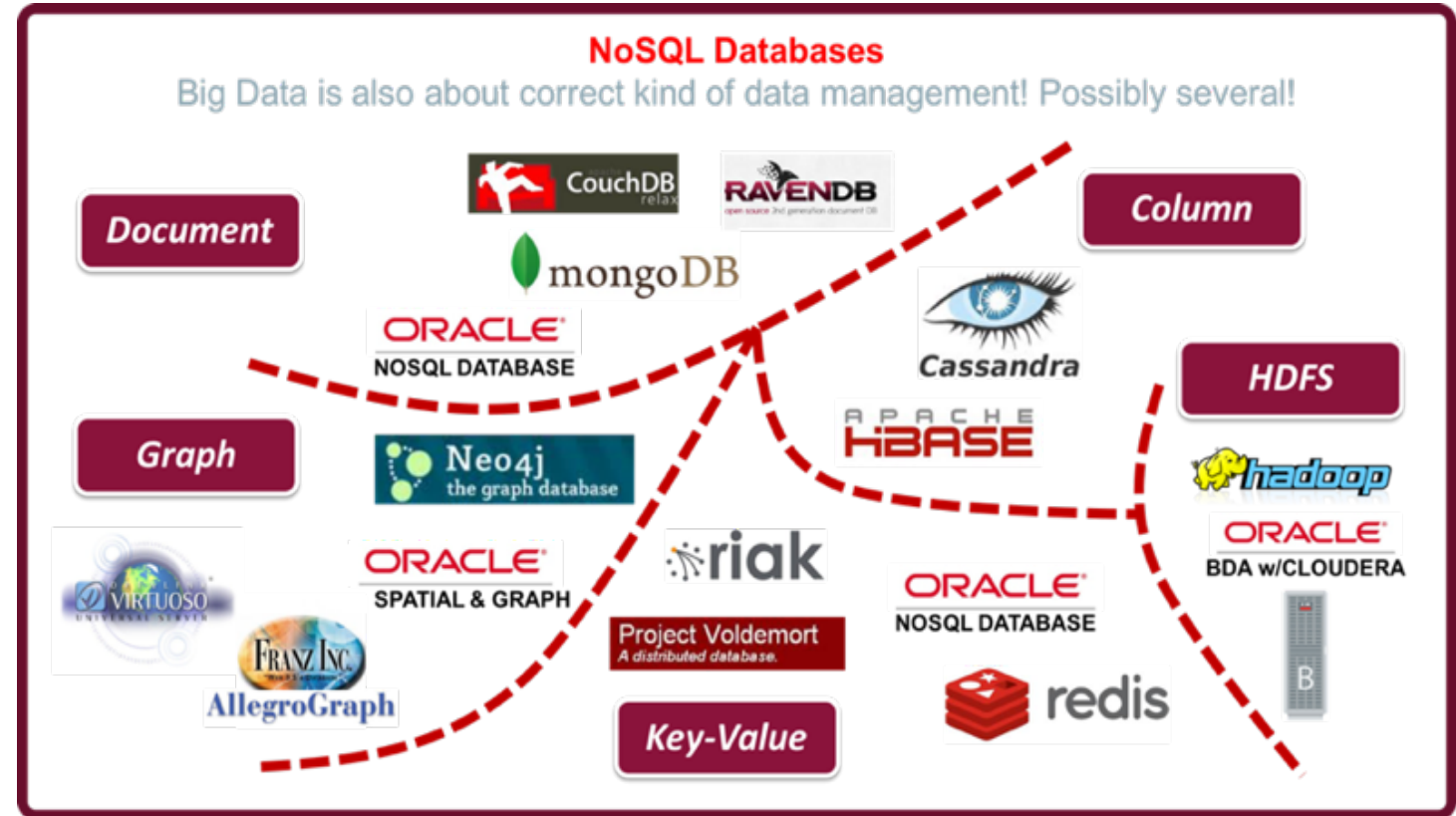- RAM: 64 GB DDR4 ECC
- HDDs: 2 x 1TB SATA 7.2k RPM

**In 2020**, the following cloud infrastructure can be «rented» for NOK630/month
- 64 cores
- 512 GB memory
- 2 x 25 GB Gbps network
- Up to 1 PB of block volumes (5TB/month free)
- Can be reduced or increased as needed, giving an average usage for slightly better price and with better performance

# NoSQL DBMS – a collective concept

Relational took over, but both hierarchical and network models are still in use, and there are several DB technologies introduced in addition

# Databases, Database Systems & Database Managament Systems

o IMPORTANT TO UNDERSTAND THE DIFFERENCE

o Database System (DS) =
Database Management System (DBMS):

o The engines that «run» the databases

o Run the queries, optimize, manage transactions, help manipulate the tables or other DBMS structures etc.

o What else?

UiO **: Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet

# The reality of database usage

o Concurrent multi-user (many users use the database at the same time)

o Usage is random and needs to be controlled, orchetsrated.

o Machines and storage can crash.

o For distributed databases: The network can be broken, messages can be lost

# ACID

o Four very important characteristics

o Atomicity

o Consistency

A and C are for correctness

o Isolation

I for multi-user concurrency issues

o Durability

D for avoiding data loss

# Example DBMS architecture



INF3100 - xx.1.2017 – Ellen Munthe-Kaas

UiO : Institutt for informatikk
Det matematisk-naturvitenskapelige fakultet