# IN3020/IN4020 – Database Systems Spring 2021, Week 18.2

# Data Science:
# The bridge between data & science

M. Naci Akkøk, CEO, In-Virtualis, Assoc. Prof. UiO/Ifi (ASR), Assoc. Prof. OsloMet/CEET (Engineering)

Egor Kostylev, Assoc. Prof, UiO/Ifi, Analytical Solutions & Reasoning (ASR) research group

Sagar Sen, Senior Research Scientist, SINTEF Digital

Renée Wikestad, Principal Cloud Engineer & Specialist, Oracle

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

# Introduction

# Key question – Why/how of newer DBMS´

When new types of applications become main-line, like data science or AI/ML applications or social media applications or IoT etc., they impose new requirements upon the underlying systems, like the infrastructure, and, of course, the DBMS.

How do the new generation DBMS´
answer the needs of data science?

# The program

o Introduction (M. N. Akkøk, UiO/Ifi + OsloMet) – 10 min.

o Graph Neural Networks (E. Kostylev, UiO/Ifi) – 20 min.

o Explainable AI (S. Sen, SINTEF) – 20 min.

o AI/ML and the DBMS in practice (R. Wikestad, ORACLE) – 20 min.


o For discussion:
What does Data Science require from DBMS´? – 20 min.

# A very simple introduction to data science

Insight to the data science techniques, including AI/ML

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

# Statistics – the common base

o Understanding statistical techniques, how they relate to each other and where they are used

    o Mean, variance, standard deviation

    o Curve fitting and simple prediction

    o Anomalies, clustering (K-means), classification

# Statistics – towards data science

o Mean, variance, standard deviation

o Curve fitting and simple prediction

o Anomalies, clustering (K-means), classification

**Mean**. The average of all data-points.
**Variance**. A measure of the degree to which each data point deviates from the
**mean**. A measure of the average distance from the average ☺
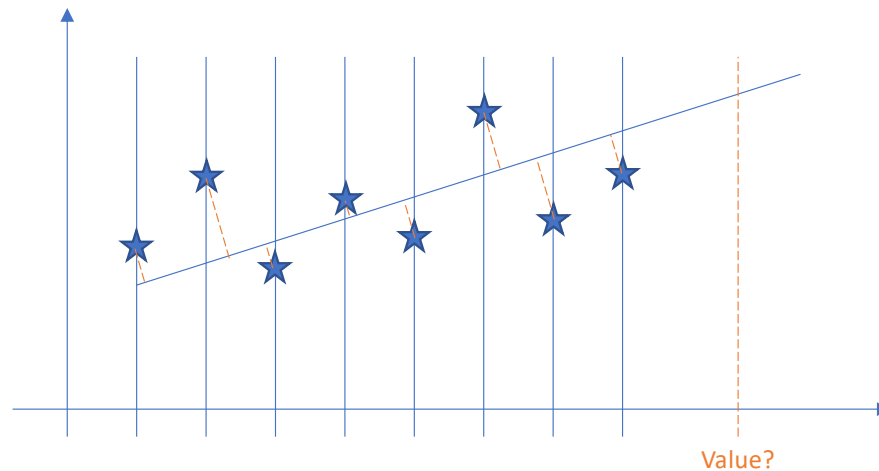**Standard deviation**. Square root of the **variance**.
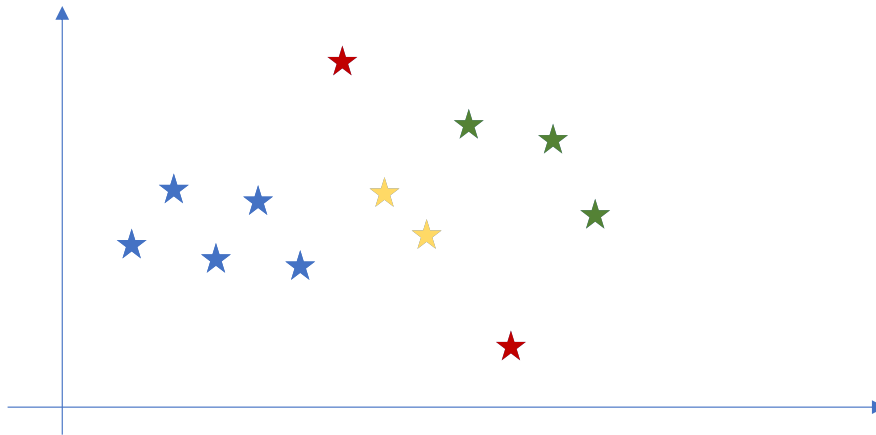An indication the spreading of numbers from the **mean**.

# Statistics – towards data science

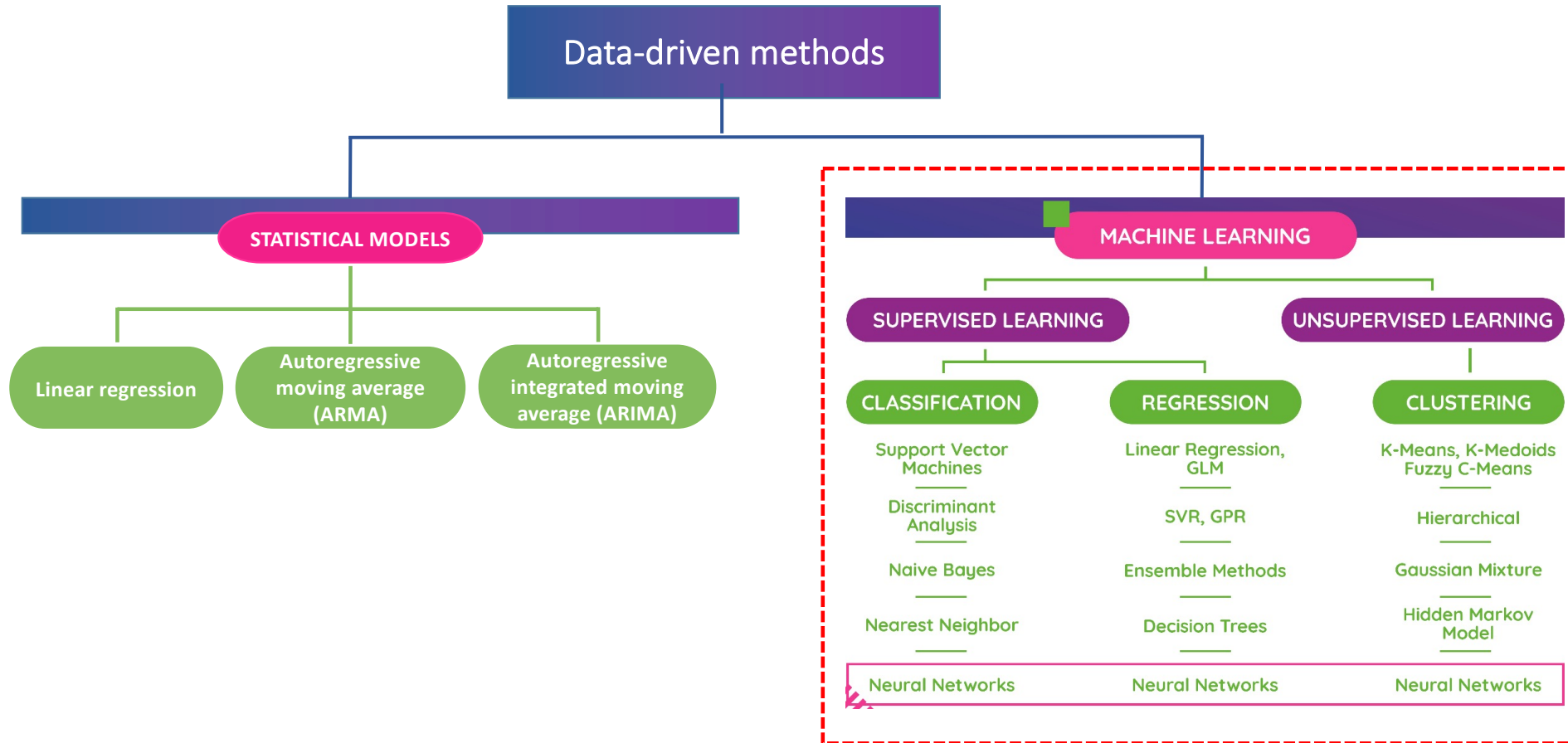o Mean, variance, standard deviation

o Curve fitting and simple prediction

o Anomalies, clustering (K-means), classification



Value?

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

# Statistics – towards data science

o Mean, variance, standard deviation

o Curve fitting and simple prediction

o Anomalies, clustering (K-means), classification

REF.: Dr. Hossain Tabari, Research Scientist, KU Leuven, Department of Civil Engineering

## Supervised learning

**Data:** $(x, y)$

$x$ is data, $y$ is label.

**Goal:** Learn a function to map $x \longrightarrow y$

**Use cases:** Classification, regression,
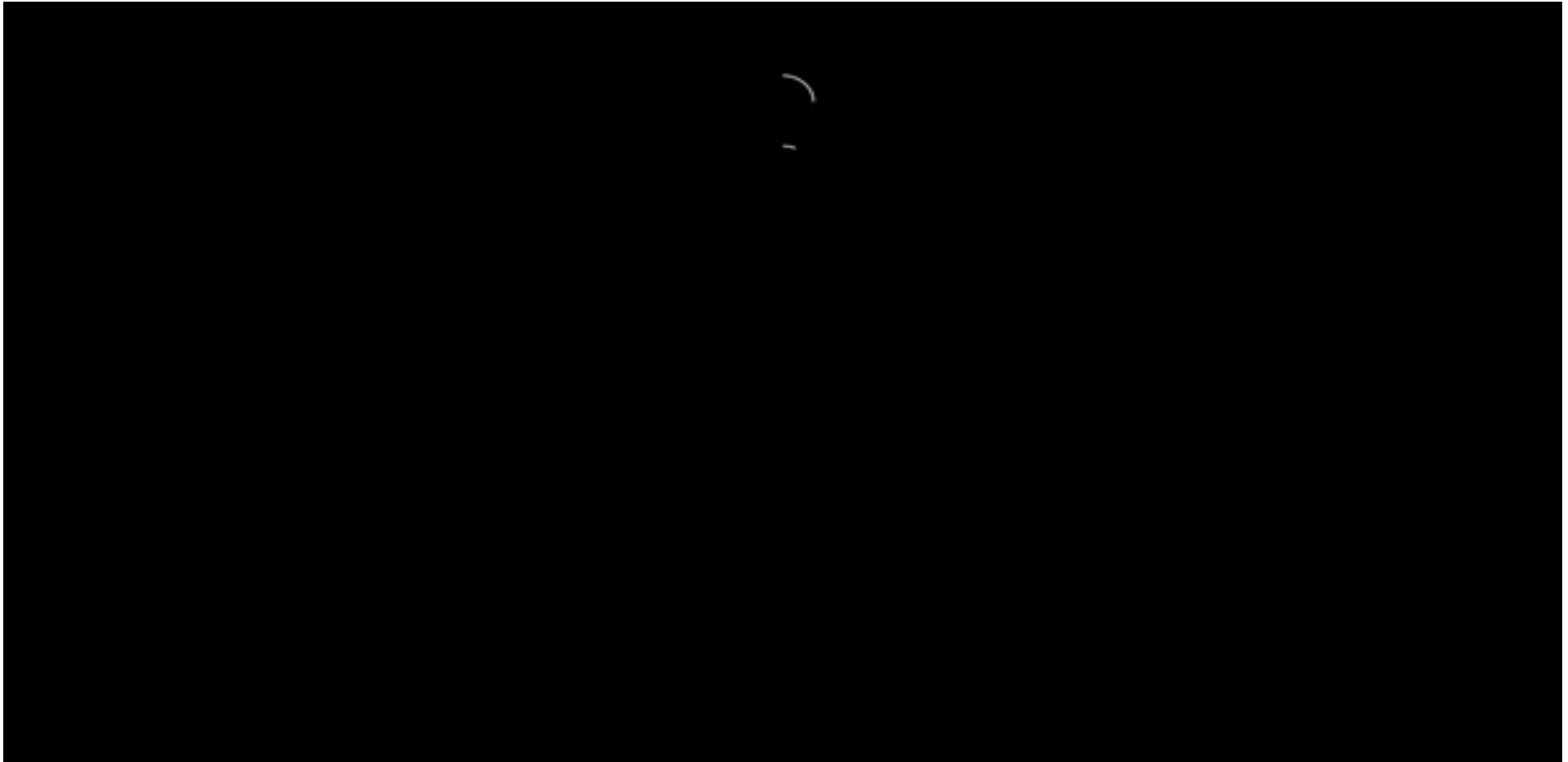
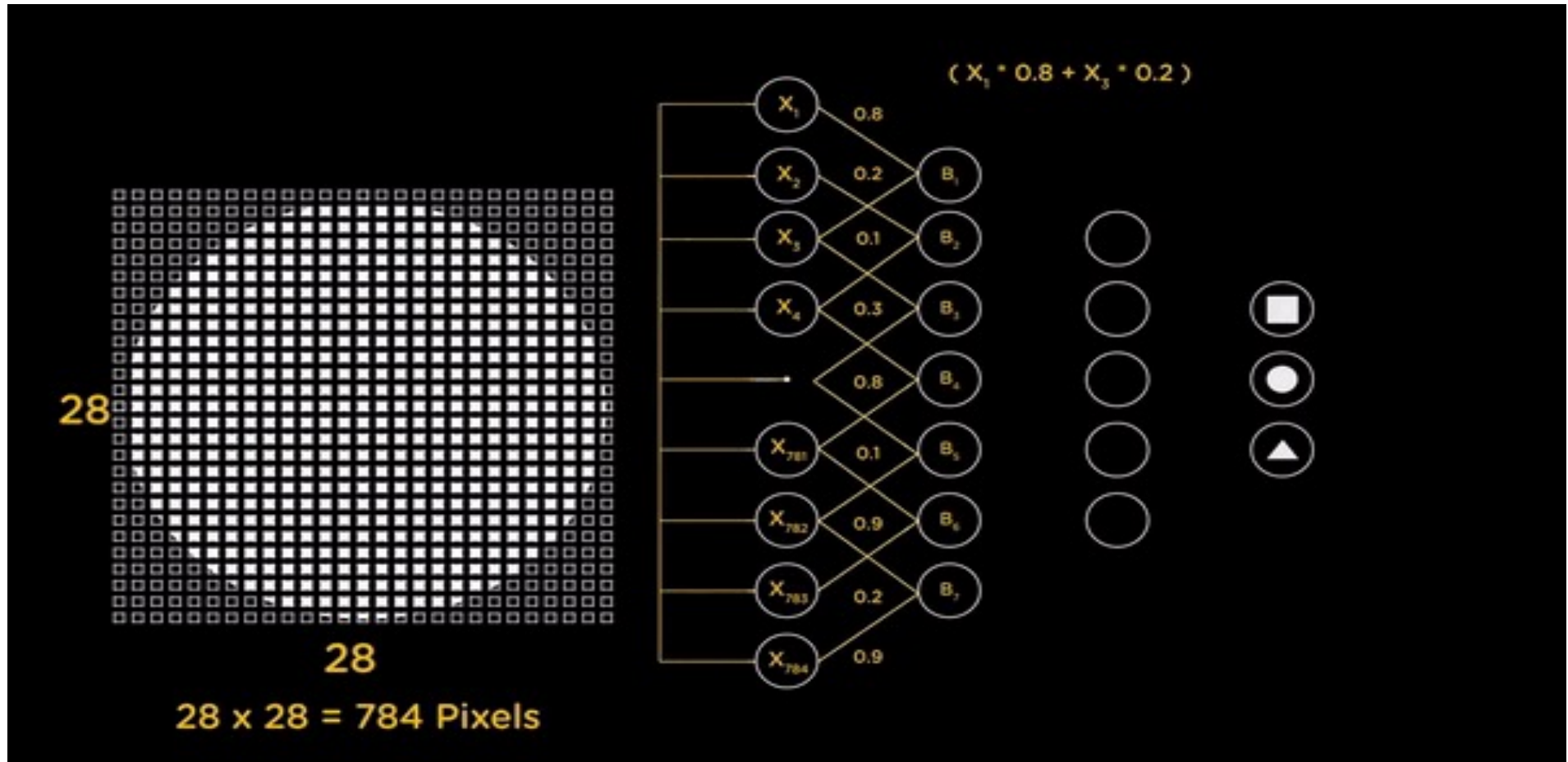object detection, etc.

## Unsupervised learning

**Data:** $(x)$

$x$ is data, no labels!

**Goal:** Learn some hidden or underlying structure of the data

**Use cases:** Clustering, feature or dimensionality reduction, etc.

REF.: Dr. Hossain Tabari, Research Scientist, KU Leuven, Department of Civil Engineering

## Artificial neural networks



REF.: Dr. Hossain Tabari, Research Scientist, KU Leuven, Department of Civil Engineering

REF.: Dr. Hossain Tabari, Research Scientist, KU Leuven, Department of Civil Engineering

# Topic for discussion:
# What does data science require of database management systems?

What are the characteristics of data science applications?

# Characteristics of Data Science

o Data science applications are <span style="color:red">data-driven</span> or <span style="color:red">data-intensive</span>

o Remember Big Data characterization?

| | |
|---|---|
| Volume | Large amounts of data (maybe not always?) |
| Velocity | Fast! Real time? Ingestion, access, processing? |
| Variety | Different sources, formats, granularity |
| Value | Quality, completeness, fitness-to-purpose |

UiO **:** **Institutt for informatikk**
Det matematisk-naturvitenskapelige fakultet
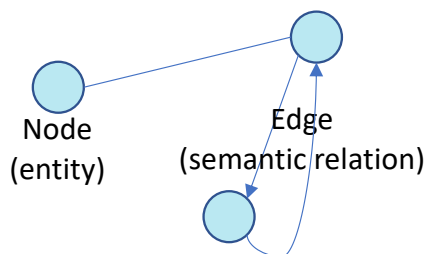
# Example: Computer aided molecular design (CAMD)

o Computational molecular modeling & simulation

o Very complex and time consuming

o Some optimization solutions, directions:
  o Efficient data representation: Representing molecular structure & free-energy data using (extended) property graphs
  o Using graph structure in machine learning (Egor & Sagar)
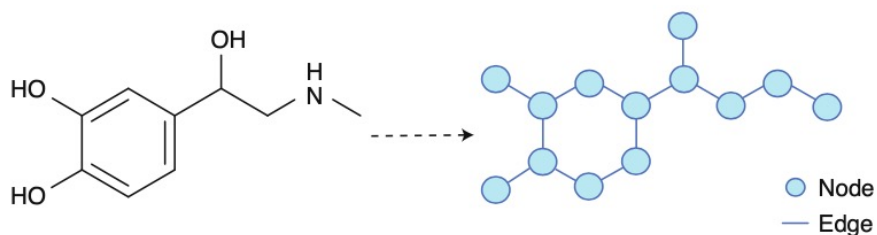  o Processing (search/retrieval) optimization

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

# Technology #03 (Using Extended Property Graph)

ANVIL
by
In-Virtualis

- Graph

Node
(entity)

Edge
(semantic relation)

- Graph representation of molecular structure is more natural & effective

○ Node
— Edge

- Property Graph

$N(1)$

$h_4^{(l)}$  $\alpha_{14}^{(l)}$  $h_2^{(l)}$

$h_1^{(l+1)}$  $\alpha_{12}^{(l)}$

$\alpha_{13}^{(l)}$

$h_3^{(l)}$

- Extenden Property Graph
  - Add quantum/mechanical/stochastic energy and interaction info
  - Enable dynamic edges
  - Also adding folding info (3D geometry),
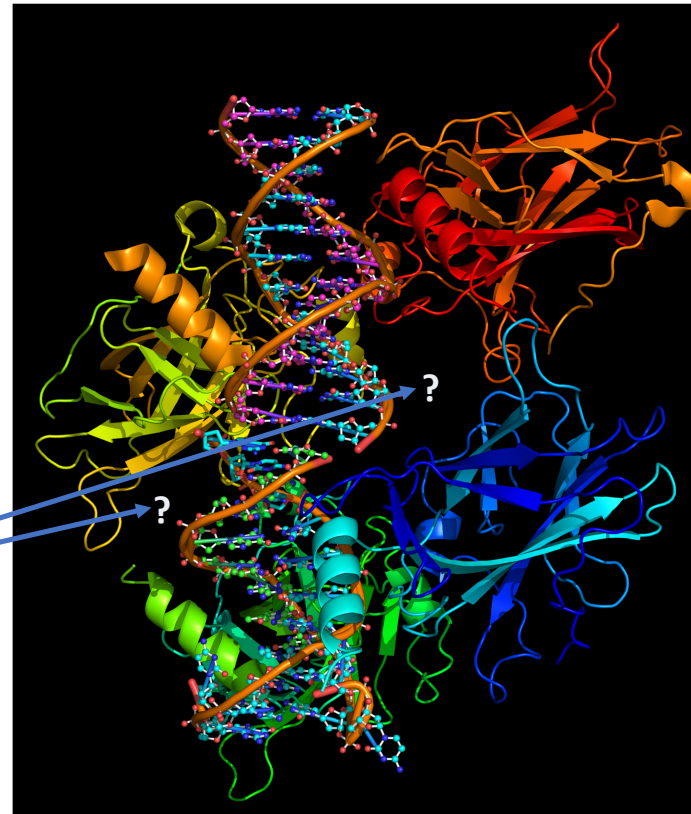  - and affinity info (see next slide)

UiO : Institutt for informatikk
Det matematisk-naturvitenskapelige fakultet

# Technology #04 (Folding, Affinity & VR)



Molecules "fold" due to the distribution of pull-and-push forces.

When folded, some areas with specific surpluss "charge" (positive or negative) are exposed,

where other structures (with proper "charges") may "dock".

Much easier to see & refine in VR.

UiO : Institutt for informatikk
Det matematisk-naturvitenskapelige fakultet

REF.: Dr. M. Naci Akkøk, CEO, In-Virtualis AS

# Example: Computer aided molecular design (CAMD)

o Computational molecular modeling & simulation

o Very complex and time consuming **– obvious need for HPC**

o Some optimization solutions, directions:

  o Efficient data representation: Representing molecular structure & free-energy data using (extended) property graphs

  o Using graph structure in machine learning (Egor & Sagar)

  o Processing (search/retrieval) optimization

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

# Topic for discussion:
This is what we briefly saw last time

# Date Quality Management (DQM)

o **Data Quality Management** is for ensuring the quality of the data (the information) we are pulling into the database or data we already have in the database.

o Any application that relies on data
  o Will need data of high quality
    (high enough with respect to its purpose)
  o And will be negatively affected by low quality data (may fail or give wrong information)

# Data Quality?

Data is of high quality, if the data is **fit for the intended purpose** of use and if the data **correctly represent the real-world construct** that the data describes.

*Ref. Profisee*

| Characteristic | How It's Measured |
|---|---|
| Accuracy | Is the information correct in every detail? |
| Completeness | How comprehensive is the information? |
| Reliability | Does the information contradict other trusted resources? |
| Relevance | Do you really need this information? |
| Timeliness | How up- to-date is information? Can it be used for real-time reporting? |

*Ref. Syncsort*

Content & format **CORRECTNESS**:

Key fields and other relevant fields are non-empty and are in the right format, and of the right type.

Content of field makes sense with respect to its format and expected use.

*Ref. practice @ Oracle*

**FUNCTIONS**: Profiling, auditing, visualization, parsing & standardization, matching & merging , case-based clean-up, address/format verification.
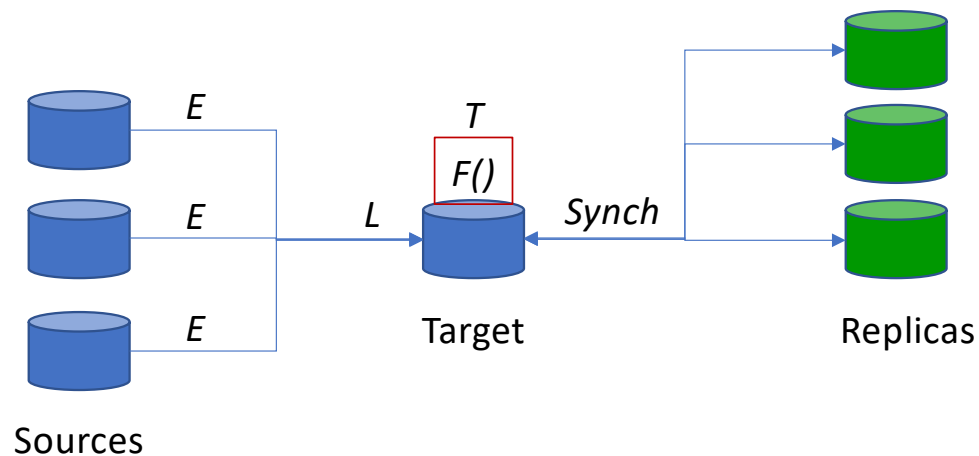
*Ref. EDQ sheet @ Oracle*

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

# DQM is related to ETL/ELT and DRM tools

o **Data Quality Management** is practically always related to **Extract-Transform-Load** or **Extract-Load-Transform**$^{(*)}$ tools and **Data Replication Management** tools.

o Typical process

(*) ELT usually faster

# AI, ML and Data Science
## New areas of applications
## that demand new types of DBMS

Yet another category of data intensive applications!

# AI/ML in DQM

o One main area is introducing AI/ML to automate and improve the DQM process
  o Learn and automate profiling
  o Do the necessary corrections
  o Learn and automate transformations
  o Learn and improve ingestion performance
  o …

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

# AI/ML in DBMS

o One main area is introducing some AI/ML into the DBMS
- o Self-install
- o Self-tune (through learned and optimized indexing etc.)
- o Self-repair

- o ...

# A DBMS for AI/ML

o How would you best represent a deep learning (neural) network in a DBMS?

   o Suggestion: Graph Neural Networks (Explainable AI)

o How would you design your DBMS to improve performance in the case of AI/ML and Data Science applications that require large amounts of data fast?

   o Suggestion: Embed the algorithms into the DB and let them run in the database (so that you don´t have to pull the data out first)

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

# In-Virtualis

Contact:     naci@in-virtualis.com
             nacia@ifi.uio.no
             mehmetna@oslomet.no

             +47 47026879

enabled by

The Life Science Cluster

# Join us
## for the first MoMS/CAMD Seminar
## June 9th at 09:00-11:00

Register at https://www.in-virtualis.com/events
(on-line, soon to be published)

UiO : Institutt for informatikk
Det matematisk-naturvitenskapelige fakultet