



GWE #3

Database, Data Warehouse, Data Lakes

Database vs Data Warehouse vs Data Lake

Data Warehouse

VS

Data Lake

VS

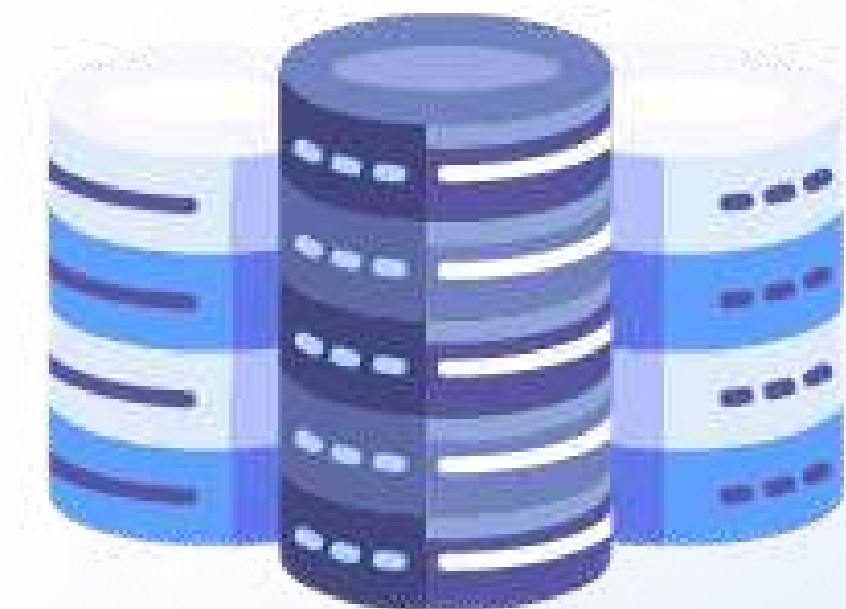
Database



- integrated
- for reporting purpose

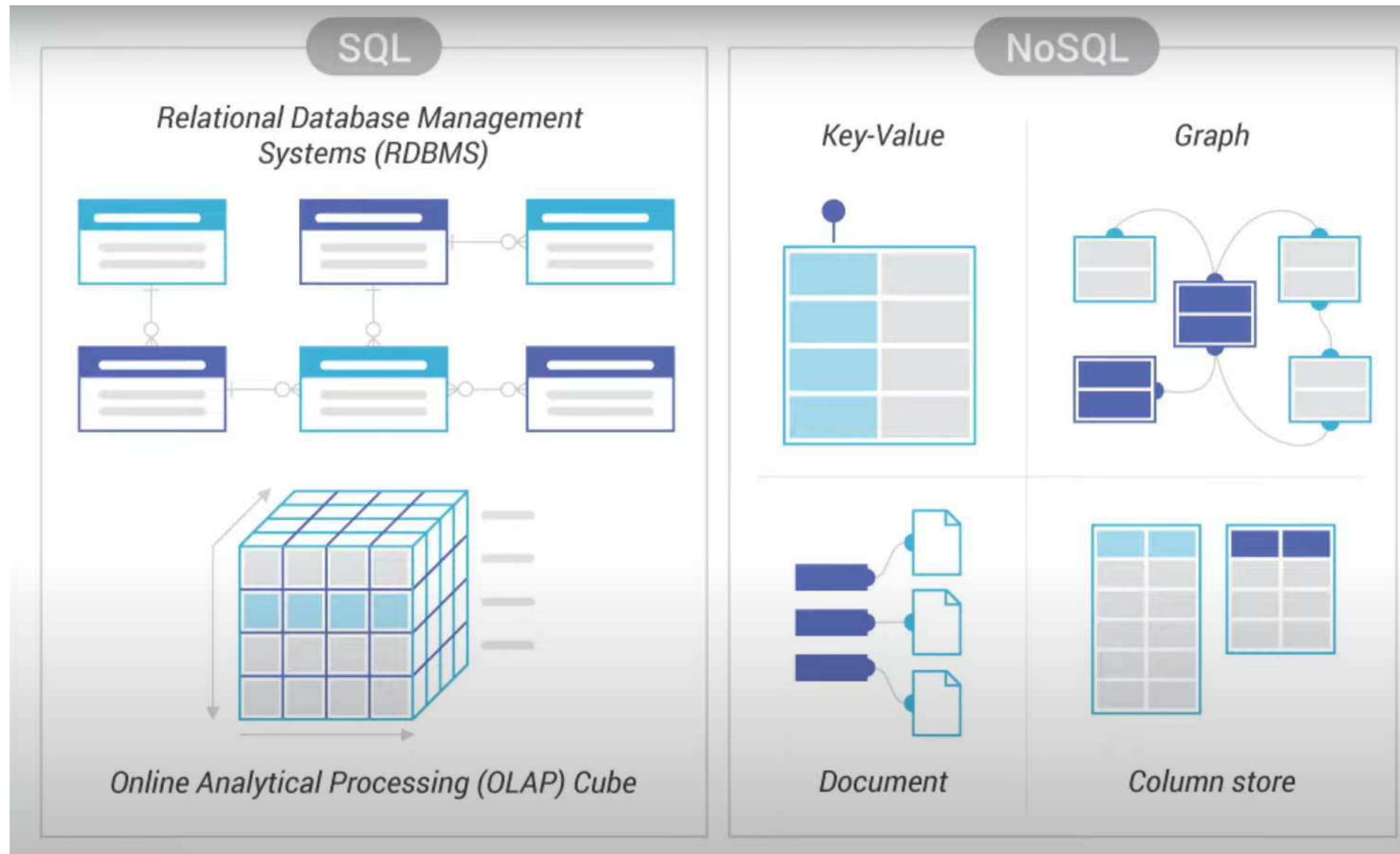


- Schema-on-read
- for ML purpose

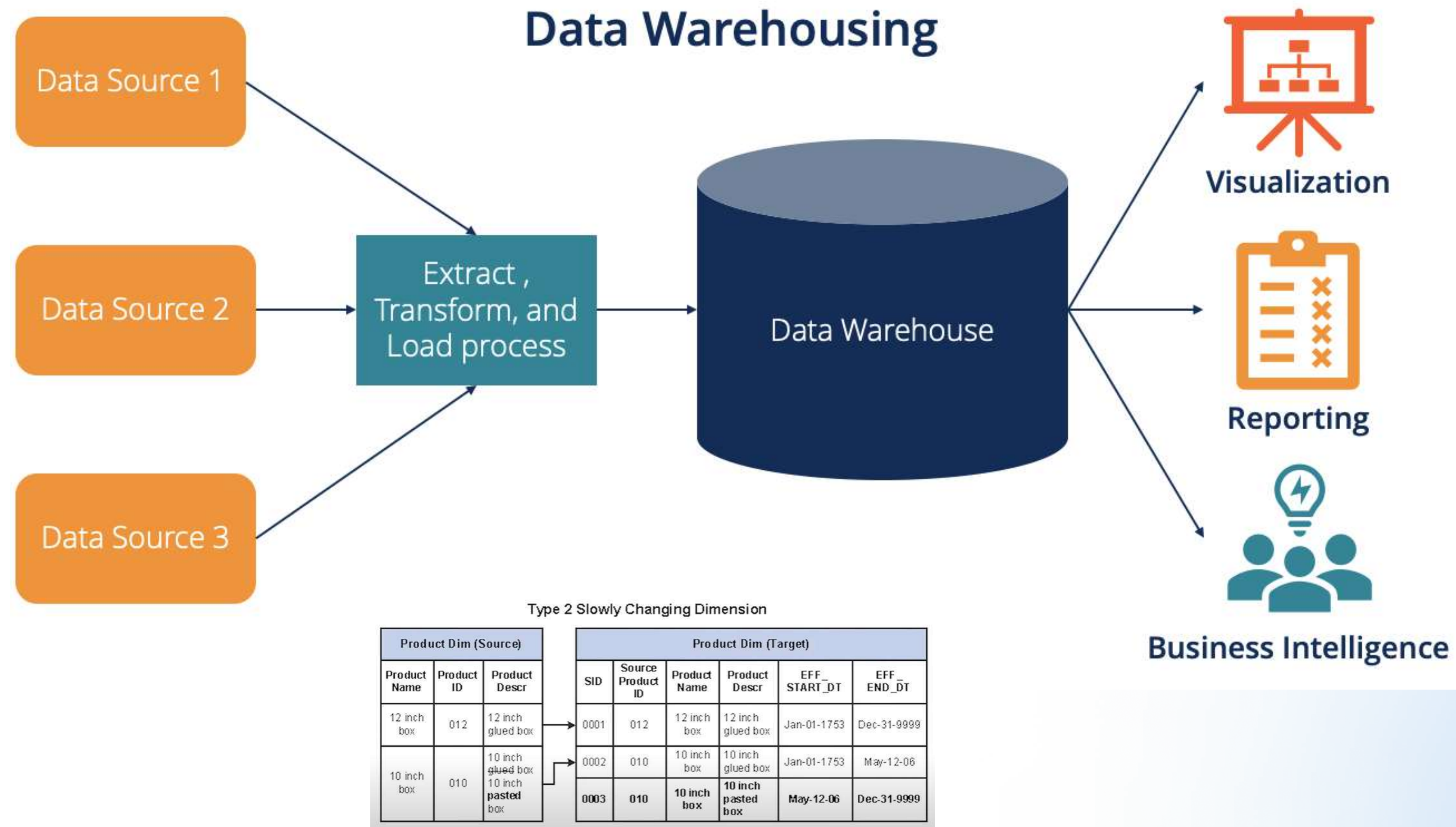


- CRUD
- Operational purpose (software)

Database

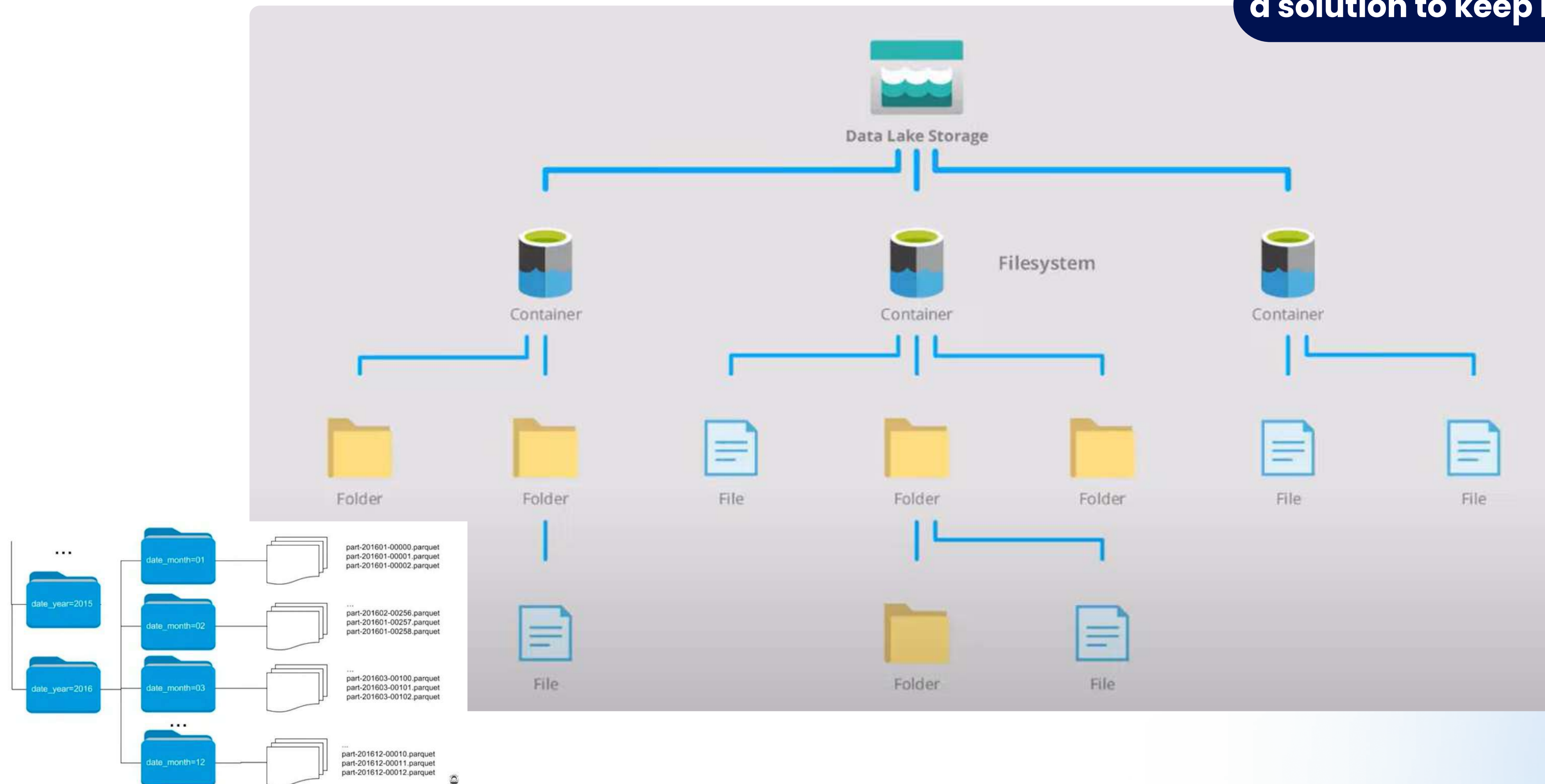


Data Warehouse



Data Lakes

a solution to keep raw data



Kenapa Data Lakes

Data Lakes



Raw

Data Lakes contain unstructured, semi structured and structured data with minimal processing. It can be used to contain unconventional data such as log and sensor data

Large

Data Lakes contain vast amounts of data in the order of petabytes. Since the data can be in any form or size, large amounts of unstructured data can be stored indefinitely and can be transformed when in use only

Undefined

Data in data lakes can be used for a wide variety of applications, such as Machine Learning, Streaming analytics, and AI

Data Warehouse



Refined

Data Warehouses contain highly structured data that is cleaned, pre-processed and refined. This data is stored for very specific use cases such as BI.

Smaller

Data Warehouses contain less data in the order of terabytes. In order to maintain data cleanliness and health of the warehouse, Data must be processed before ingestion and periodic purging of data is necessary

Relational

Data Warehouses contain historic and relational data, such as transaction systems, operations etc



Kenapa Data Lakes

Kepemilikan dapat dibagi secara terbatas

Simple Data Management with more Context

Fondasi bagi AI dan ML

More Security & Governance

Bahas sedikit tentang DBMS

SQL

S = Structured

Q = Query

L = Language

NoSQL

N = Not

O = Only

S = Structured

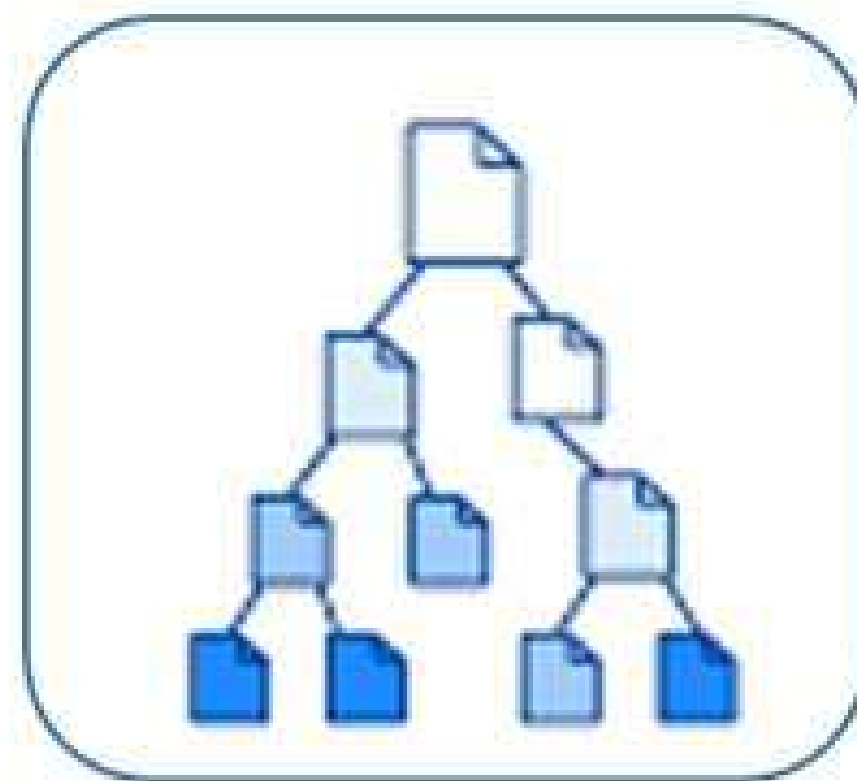
Q = Query

L = Language

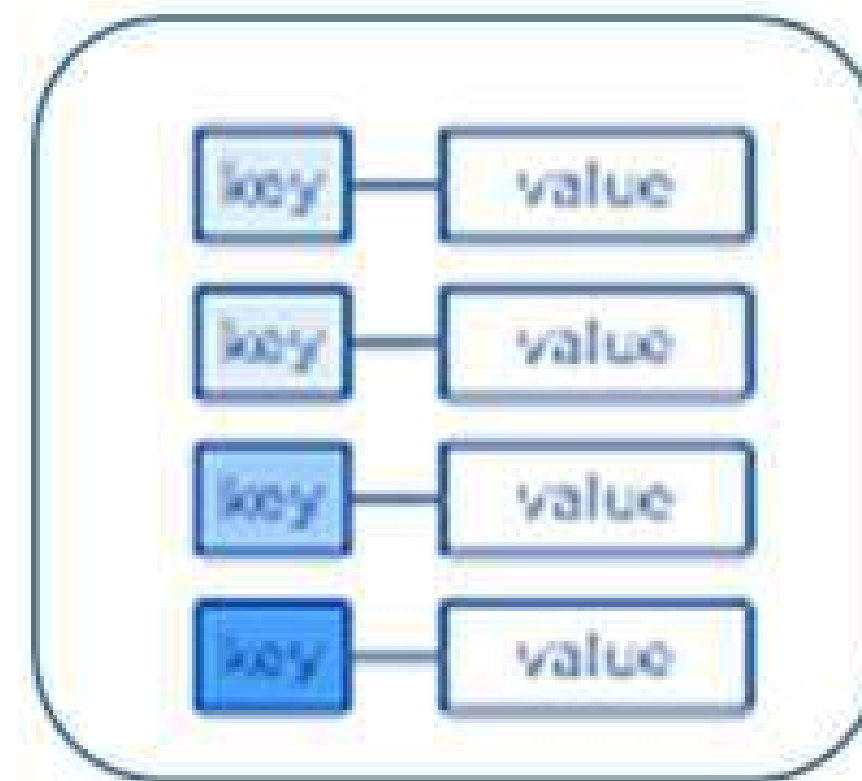
Bahas sedikit tentang DBMS

NoSQL

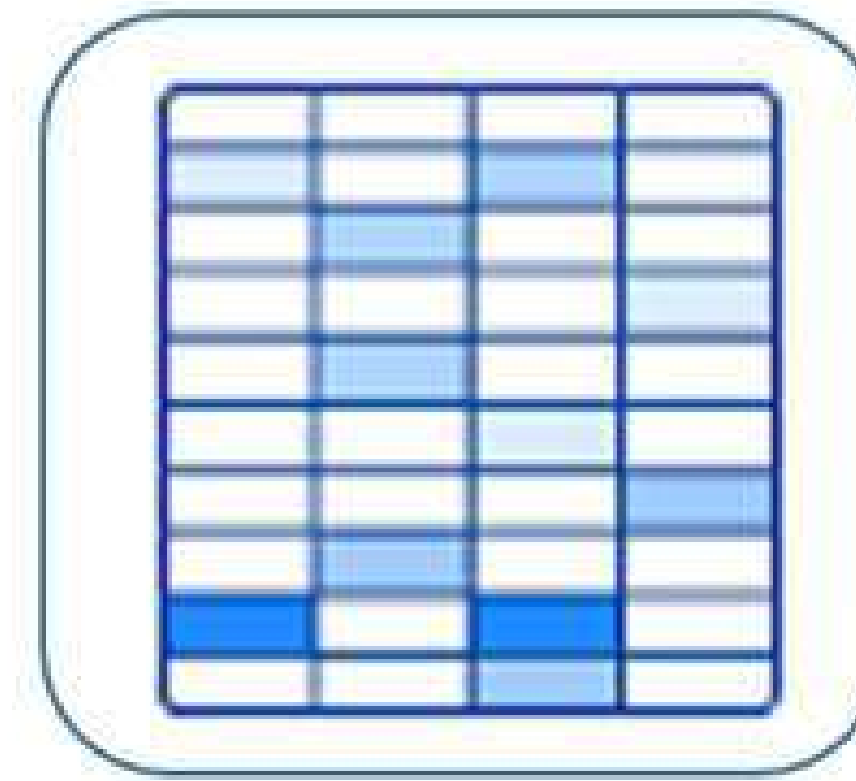
-> sistem manajemen database yang tidak memiliki relasi.



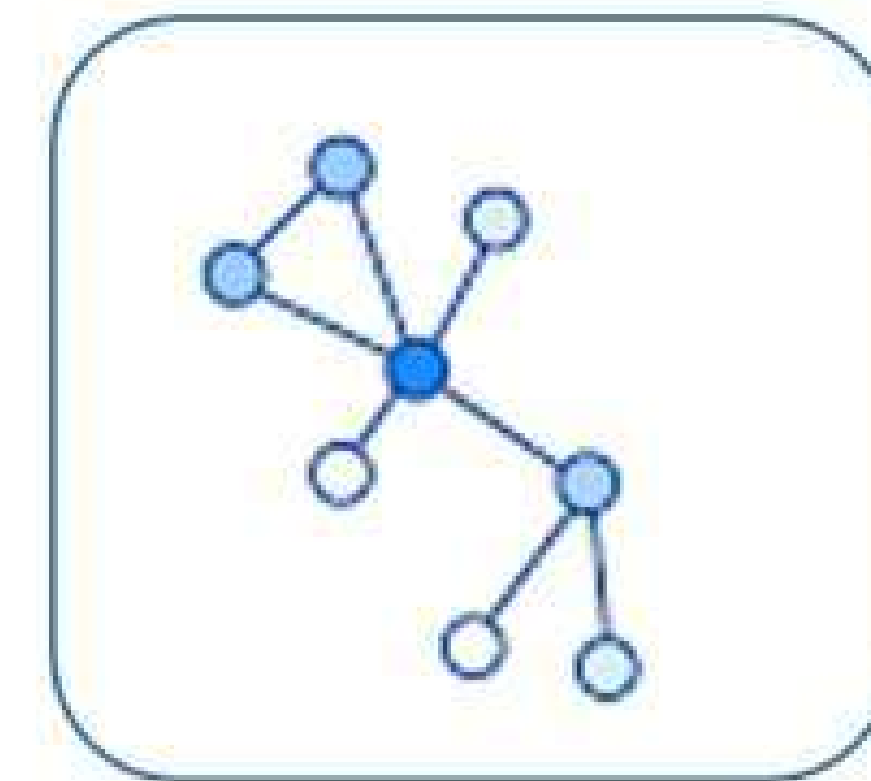
Document
Store



Key-Value
Store



Wide-Column
Store



Graph
Store

Bahas sedikit tentang DBMS

SQL

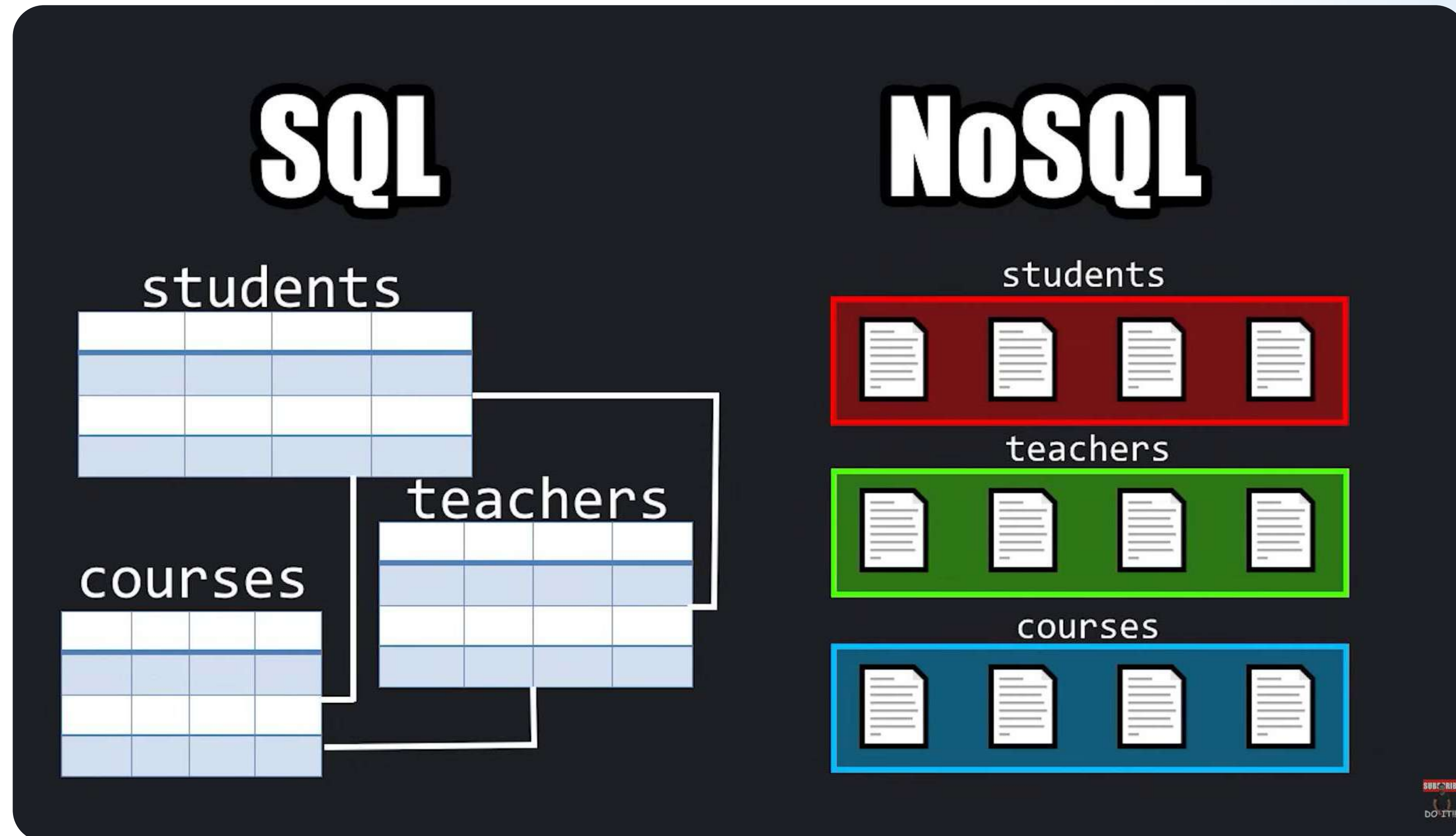
name	age	gpa	fullTime
Spongebob	32	3.2	false
Patrick	38	1.5	false
Sandy	27	4.0	true

NoSQL

```
{  
  name: 'Spongebob',  
  age: 30,  
  gpa: 3.2,  
  fullTime: false,  
},  
{  
  name: 'Patrick',  
  age: 38,  
  gpa: 1.5,  
  fullTime: false,  
},  
{  
  name: 'Sandy',  
  age: 27,  
  gpa: 4,  
  fullTime: true,  
}
```



Bahas sedikit tentang DBMS



each row = each document

Pros & Cons

SQL

- + Easy querying on relationships
- + More organized & structured (minim error)
- + Atomic
- Prepared structure
- Hard to Scale (horizontal scaling)

NoSQL

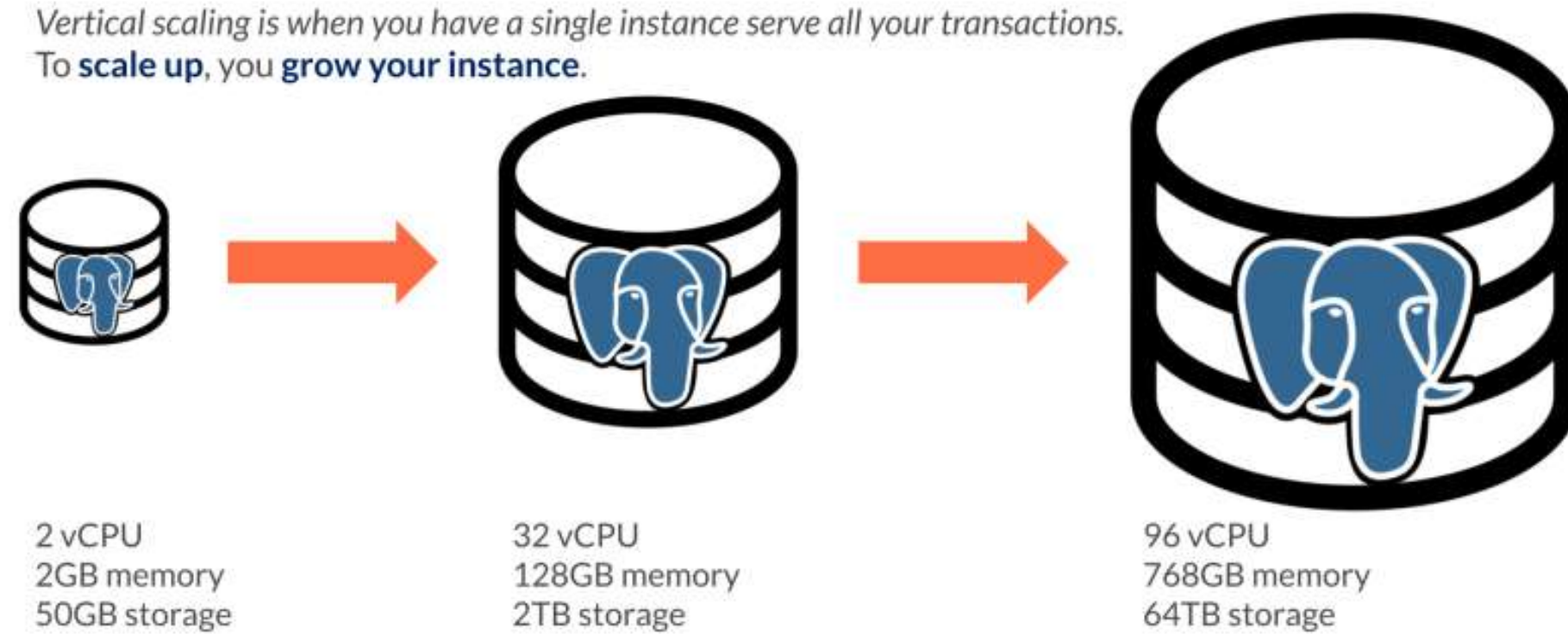
- + Flexible
- + Enabled Sharding
- Ease loss of consistency

each row = each document

Pros & Cons

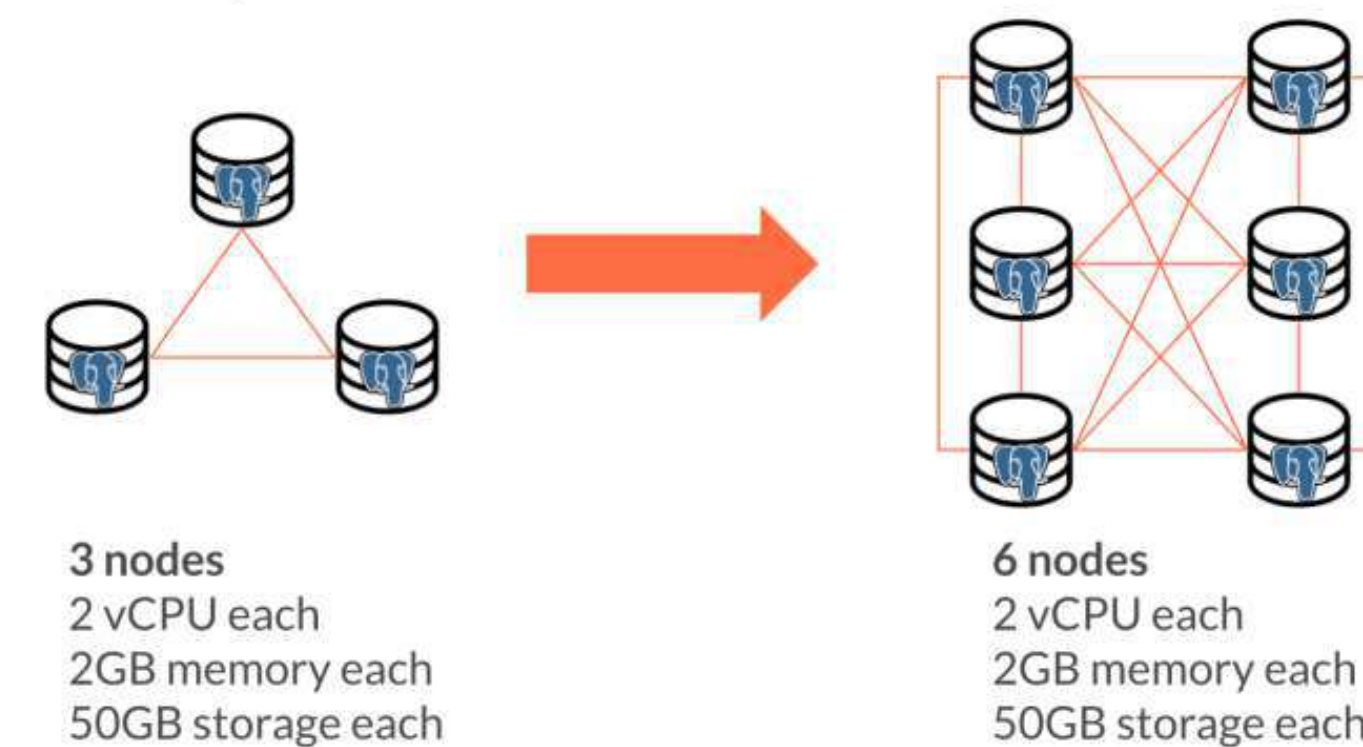
Vertical Scaling

Vertical scaling is when you have a single instance serve all your transactions.
To **scale up**, you **grow your instance**.



Horizontal Scaling

Horizontal scaling is when you have a multiple instances serving your transactions.
To **scale out**, you **add instances** (called nodes).



Pros & Cons

```
{
  title : 'movie',
  actors : [
    {
      _id: 'actor_id1',
      name: 'actor1'
    },
    {
      _id: 'actor_id'2,
      name: 'actor2'
    }
  ],
  plot: '...',
  reviews: [...],
  ...
}

{
  name : 'actor1',
  movies : [
    {
      _id: 'movie_id1',
      name: "movie1"
    },
    {
      _id: 'movie_id2',
      name: "movie2"
    }
  ],
  biography: '...',
  pictures: [...],
  ...
}
```


So which one?