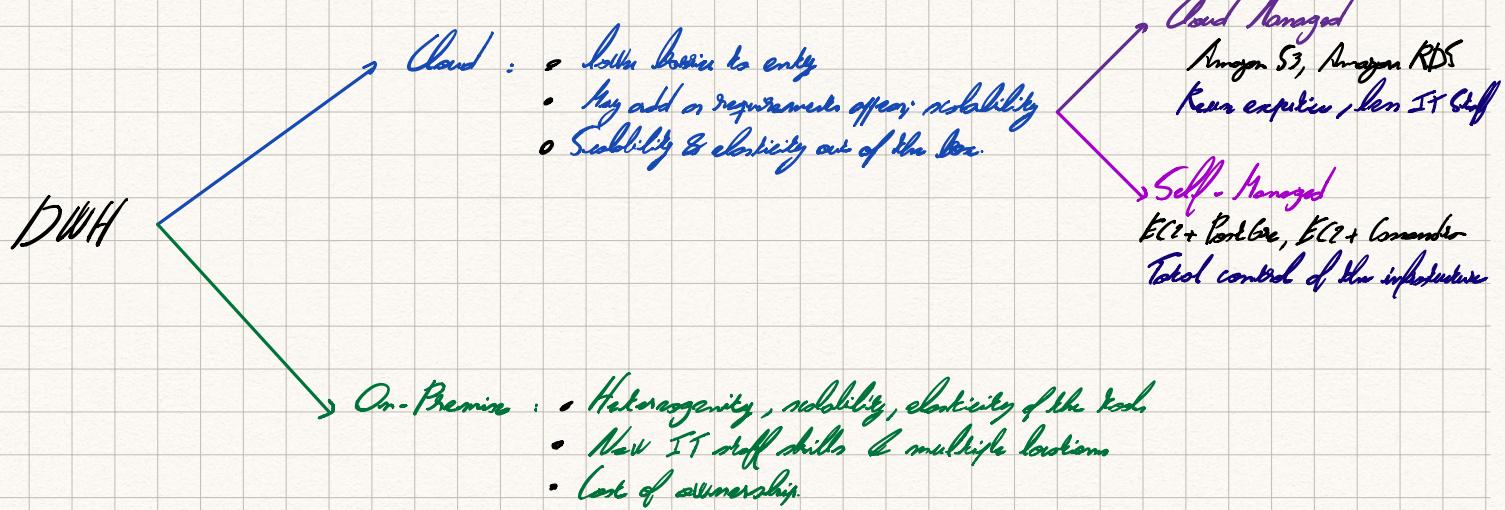


## # Choice for Implementing a Data Warehouse:



## # Amazon Relational:

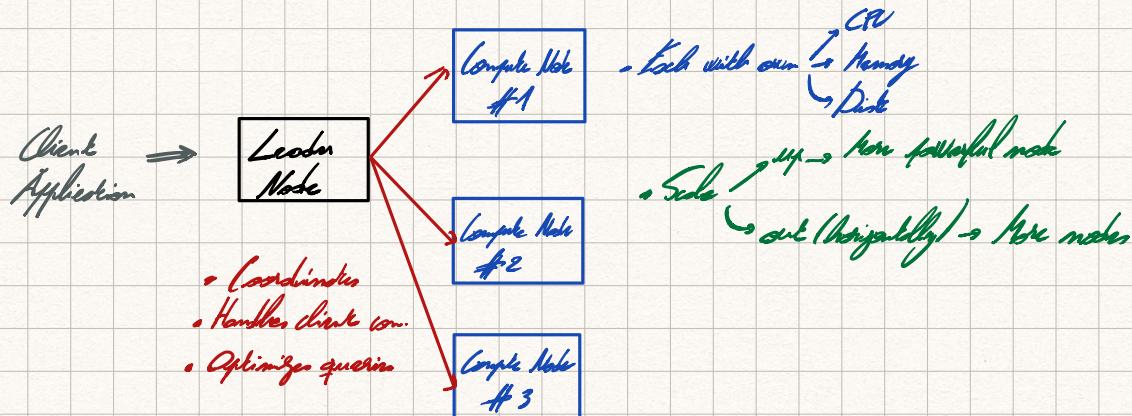
### # Relational Database:

- Every query is always executed on a single CPU of a single machine
- Acceptable for OLTP, mostly systems and low data reliability.

### # Massive Parallel Processing (MPP) Database

- Parallelize the execution of one query on multiple CPU/machines.
- Table partition (distribution across node 'clusters') becomes critical
- Partitions are processed in parallel.

## # Architecture:



## # Real-life ETL Context:



## # Optimizing Table Design:

- When a table is partitioned up into many pieces and distributed across disk in different machines, this is done blindly.
- If one has an idea about the frequent access pattern of a table, one can choose a more clever strategy
  - Two possible strategies:
    - Distribution style
    - Sorting style

## # Distribution Style:

Types:  
EVEN distribution  
ALL distribution  
RND distribution  
KEY distribution

→ sliced → dimensionally distributed.

### # EVEN distribution:

- Round robin over all slices to achieve load balancing
- Good if tables won't be joined → High cost of shuffling operations

### # ALL distribution:

- Small tables could be replicated on all slices to speed up joins
- Most frequently for dimension tables
- AKA "Broadcasting" → Prevents "massive shuffling" of dimension tables while joining.

## # FVO distribution: • Leave decision to Redshift

- Table Size ↗
- "Small enough" tables are distributed with an "All Strategy"
  - Large tables are distributed with "ELEN Strategy"

## # KBY distribution: • Rows with similar values are placed in the same slice.



This can lead to a skewed distribution of some slices of the disk key or more frequent than others

Very useful when a dimension table is too big to be distributed with All strategy.  
In this case, we distribute both the fact table and the dimension table using the same disk key

- If 2 tables are distributed on the joining keys, Redshift collects the data from both tables on the same slice.

## # Sorting key:

- One can define columns as sort key
- Upon loading, rows are sorted before distribution to slices
- Minimize the query time since each node already has contiguous ranges of rows based on the sorting key
- Useful for columns that are used frequently in sorting like the fact dimension and its corresponding foreign key in the fact table.