
Data Warehousing and Data Mining

Assignment Project Exam Help

— L2: <https://eduassistpro.github.io/> — OLAP —

Add WeChat edu_assist_pro

Part I

- Why and What are Data Warehouses?
 - Transaction Processing vs. Analytical Processing
 - Databases v

<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
Data is meaningless analysis!

Example in a finance department

- Daily transaction tasks
 - E.g., account receivable, account payable, payroll, etc.

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Columns:

Description

G/L Account

Branch

cost center

G/L account name

Tax code

Total

...

Example/2

- Weekly...monthly...yearly analytical tasks
 - E.g., Finance reports

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Why OLAP Servers?

- Different workload:
 - OLTP (on-line transaction processing)
 - Major task of traditional relational DBMS
 - Day-to-day operations: purchasing, inventory, banking, manufacturing, payroll, registration, accounting, etc.
 - OLAP (on-line <https://eduassistpro.github.io/>)
 - Major task of
 - Data analysis and decision making [Add WeChat edu_assist_pro](#)
- Queries hard/infeasible for OLTP, e.g.,
 - Which **week** we have the largest sales?
 - Does the sales of **dairy products** increase over time?
 - Generate a **spread sheet** of total sales by state and by year.
- Difficult to represent these queries by using SQL ← Why?

OLTP vs. OLAP

	OLTP	OLAP
users	clerk, IT professional	knowledge worker
function	day to day operations	decision support
DB design	application-oriented	subject-oriented
data	current, up-to-date detailed, normalized isolated	historical, aggregated, multidimensional integrated, consolidated
usage	repetitive	
access	read/write index/hash on prim. key	
unit of work	short, simple transaction	complex query
# records accessed	tens	millions
#users	thousands	hundreds
DB size	100MB-GB	100GB-TB
metric	transaction throughput	query throughput, response

Data Analysis Problems

- The same data found in many different systems
 - Example: customer data across different departments
 - The same data differently
- Heterogeneous
 - Relational DBMS, Online Transaction Processing (OLTP)
 - Unstructured data in files (e.g., MS Excel) and documents (e.g., MS Word)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Data Analysis Problems (Cont'd)

- Data is suited for operational systems
 - Accounting, billing, etc.
 - Do not support analysis across business functions
- Data quality issues
 - Missing data, imprecise data, different use of systems
- Data are “volatile”
 - Data deleted in operational systems (6months)
 - Data change over time – no historical information

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Solution: Data Warehouse

- Defined in many different ways, but not rigorously.
 - A decision support database that is maintained **separately** from the organization's operational database
 - Support **inform**g a solid platform of consolidated, h
- "A data warehouse is a subject-oriented, **nonvolatile**, and **integrated** collection of data in support of management's decision-making process."—W. H. Inmon
- Data warehousing:
 - The process of constructing and using data warehouses

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat: edu_assist_pro

Data Warehouse—Subject-Oriented

- Organized around major subjects, such as **customer, product, sales**.
- Focusing on the modeling and analysis of data for decision makers, not on **transaction processing**.
- Provide **a simple and concise** and particular subject issues by **excluding data that are not useful in the decision support process**.

Data Warehouse—Integrated

- Constructed by integrating multiple, heterogeneous data sources
 - relational databases, flat files, on-line transaction records
- Data cleaning techniques are applied.
 - Ensure consistency in naming conventions, encoding structures, attribute measures, etc. among different data sources
 - E.g., Hotel price: currency, tax, breakfast covered, etc.
 - When data is moved to the warehouse, it is converted.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Data Warehouse—Time Variant

- The time horizon for the data warehouse is significantly longer than that of operational systems.
 - Operational database current value data.
 - Data warehouse information from a historical period (years)
- Every key structure in the data warehouse
 - Contains an element of time, explicitly or implicitly
 - But the key of operational data may or may not contain “time element”.

Data Warehouse—Non-Volatile

1. A **physically separate store** of data transformed from the operational environment.
2. Operational **updates of data do not help** in the data warehouse environment.
 - Does not require **transacting, recovery, and concurrency control**
 - Requires only two operations in data accessing:
 - *initial loading of data* and *access of data*.

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Data Warehouse Architecture

- Extract data from operational data sources
 - clean, transform
- Bulk load/refresh
 - warehouse is of
- OLAP-server provide multidimensional view
- Multidimensional-olap
 - (Essbase, oracle express)
- Relational-olap
 - (Redbrick, Informix, Sybase, SQL server)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Data Warehouse Architecture

All subjects,
integrated

Advanced
analysis

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Function-oriented
systems

Subject-oriented
systems

Why Separate Data Warehouse?

- High performance for both systems
 - DBMS— tuned for OLTP: access methods, indexing, concurrency control, recovery
 - Warehouse—t P queries, multidimensional
- Different functions and different data
 - missing data: Decision support rical data which operational DBs do not typically maintain
 - data consolidation: DS requires consolidation (aggregation, summarization) of data from heterogeneous sources
 - data quality: different sources typically use inconsistent data representations, codes and formats which have to be reconciled

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Comparisons

	Databases	Data Warehouses
Purpose	Many purposes; Flexible and general	One purpose: Data analysis
Conceptual Model	ER	multidimensional
Logical Model	(Normalized) Relational	(Normalized) Star schema / cube/cuboids
Physical Model	Relational Tables	ROLAP: Relational tables MOLAP: Multidimensional arrays
Query Language	SQL (hard for analytical queries)	MDX (easier for analytical queries)
Query Processing	B+-tree/hash indexes, Multiple join optimization, Materialized views	Bitmap/Join indexes, Star join, Materialized data cube

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Comparisons/2

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Comparisons/2

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

- The Multidimensional Model

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

The Multidimensional Model

- A data warehouse is based on a multidimensional data model which views data in the form of a **data cube**, which is a multidimensional generalization of 2D spread sheet.
- Key concepts:
 - **Facts**: the s
 - Typically tr <https://eduassistpro.github.io/> er types includes snapshots, etc.
 - Measures: numbers that can d
 - Dimensions: context of the measure
 - Hierarchies:
 - Provide contexts of different granularities (aka. grains)
- Goals for dimensional modeling:
 - Surround facts with as much relevant context (dimensions) as possible ← Why?

Supermarket Example

- Subject: analyze total sales and profits
- Fact: Each Sales **Transaction**
 - Measure: Dollars_Sold, Amount_Sold, Cost
 - Calculated M
- Dimensions: <https://eduassistpro.github.io/>
 - Store [Add WeChat edu_assist_pro](#)
 - Product
 - Time

Visualizing the Cubes

- A valid **instance** of the model is a data cube

total Sales		product			
		p1	p2	p3	p4
city	NY	\$454	-	-	-
	LA	\$468	\$800	-	-
	SD	\$296	-	\$240	-
	SF	\$652	-	\$540	\$745

city	product			
	p1	p2	p3	p4
total Sales	454	-	-	925
total Sales	468	800	-	-
total Sales	296	-	240	-
total Sales	652	-	540	745

Concepts: cell, fact (=non-empty cell), measure, dimensions

Q: How to generalize it to 3D?

3D Cube and Hierarchies

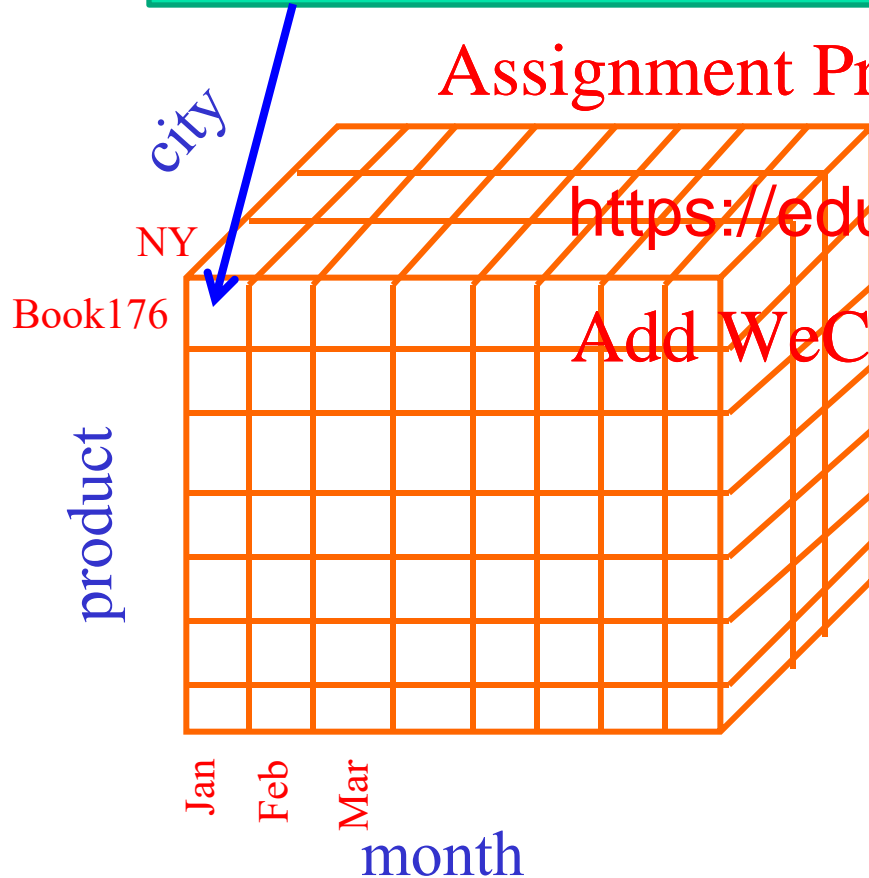
Concepts: hierarchy (a tree of dimension values), level

Sales of **book176** in **NY** in **Jan** can be found in this cell

Assignment Project Exam Help

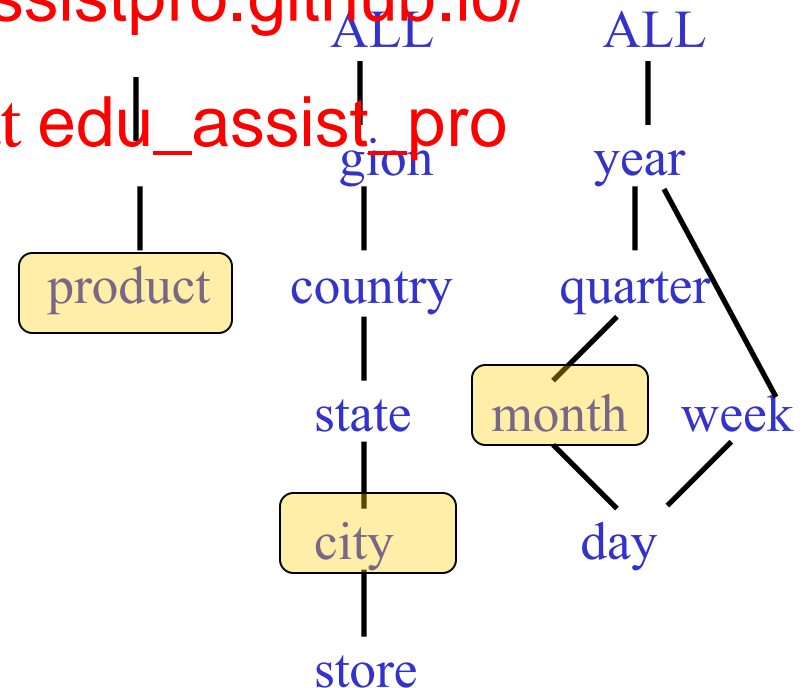
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



DIMENSIONS

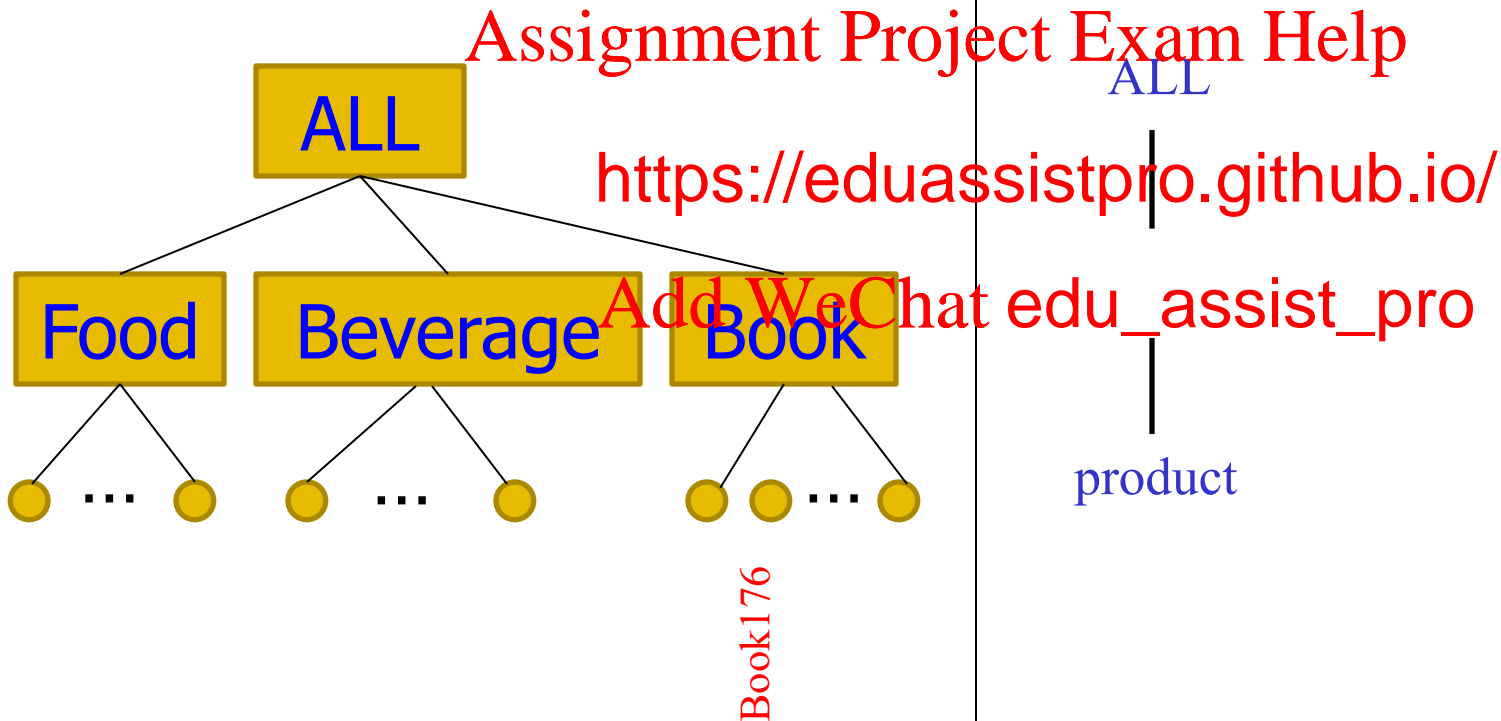
LOCATION TIME



Hierarchies

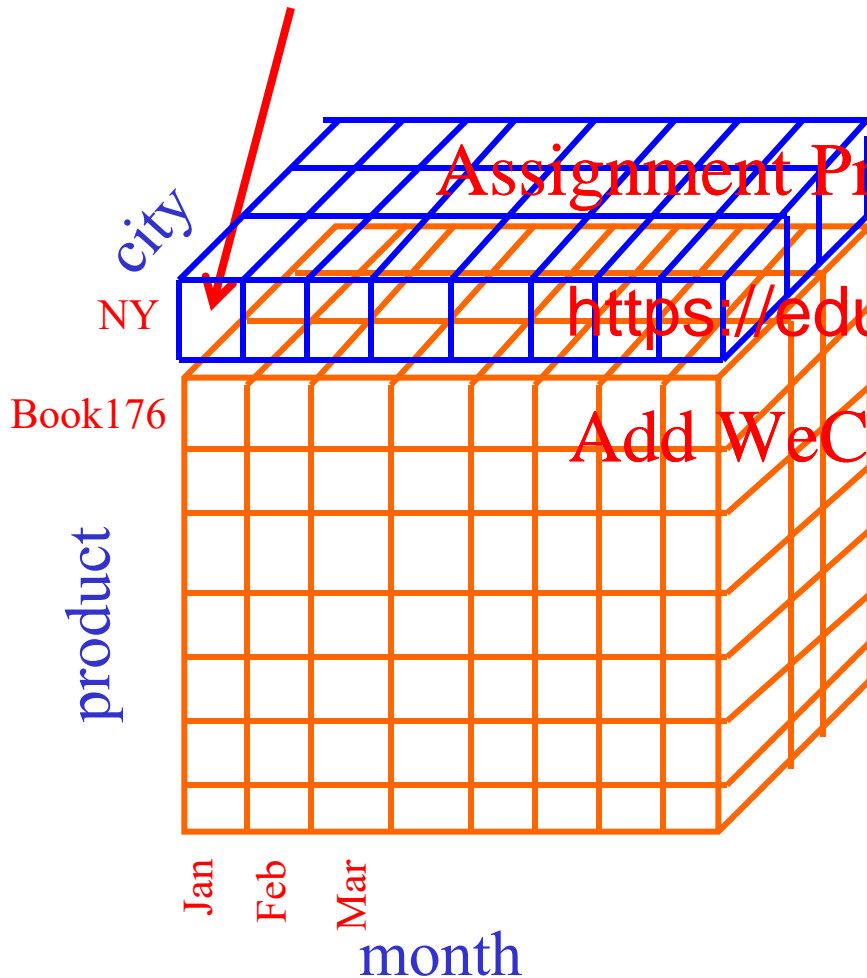
Concepts: hierarchy (a tree of dimension values), level

Which design is better? Why?



The (city, moth) Cuboid

Sales of **ALL_PROD** in **NY** in Jan



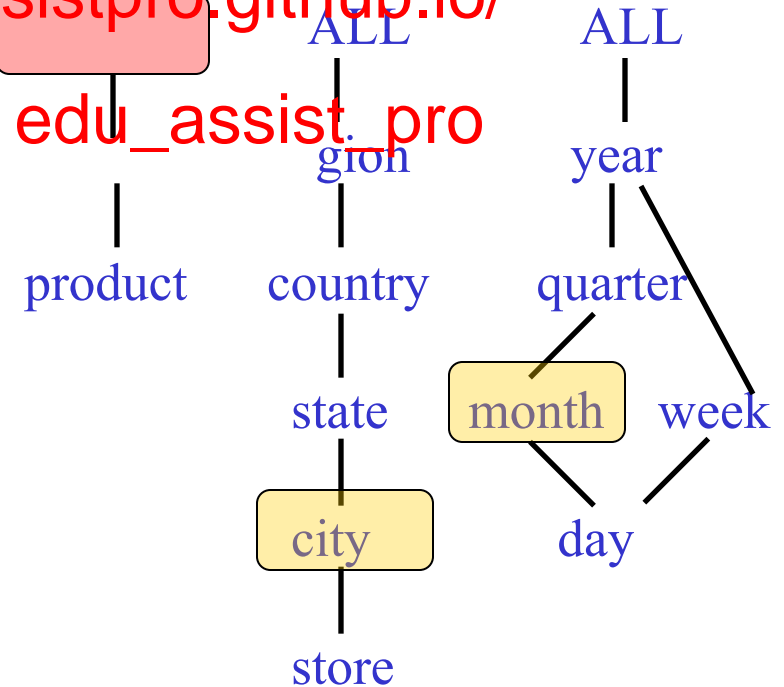
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

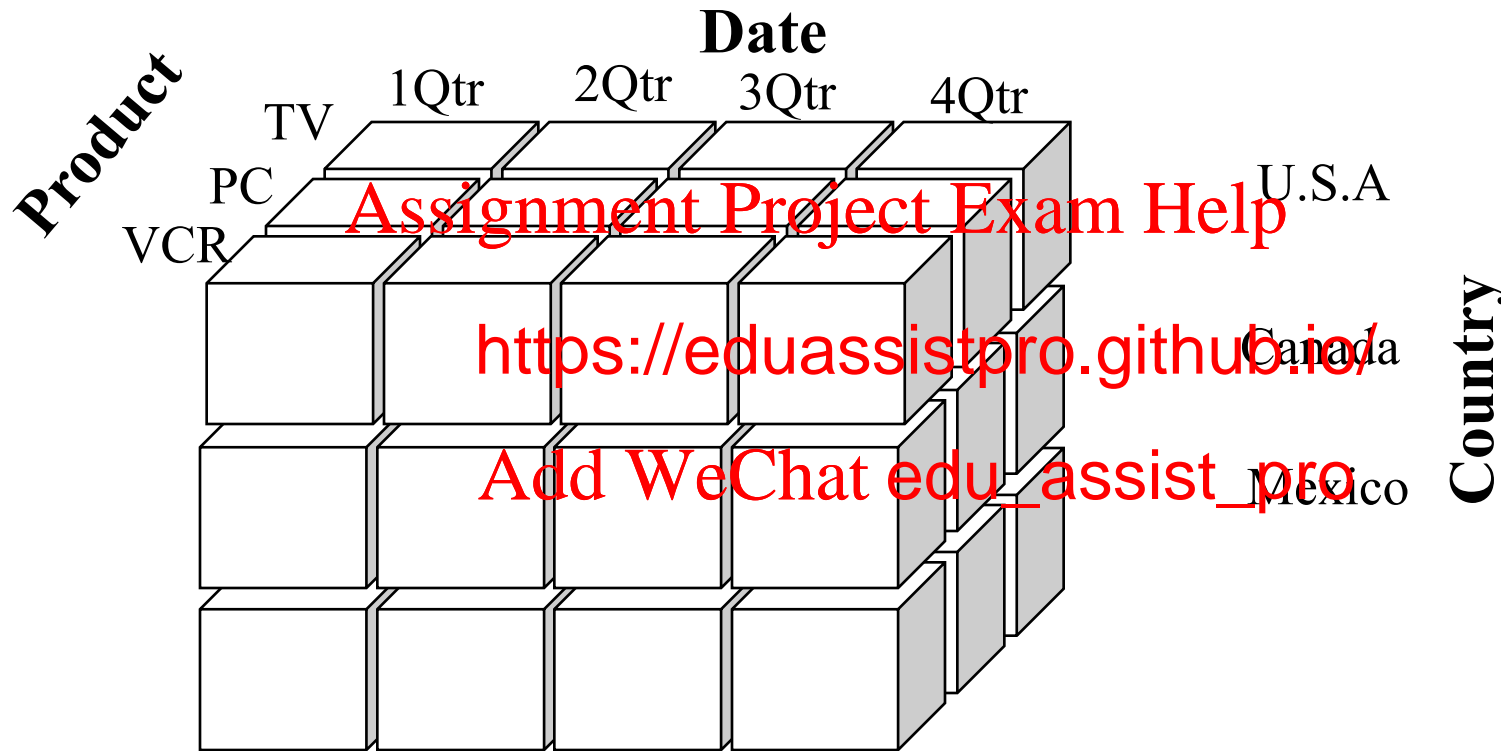
DIMENSIONS

LOCATION TIME



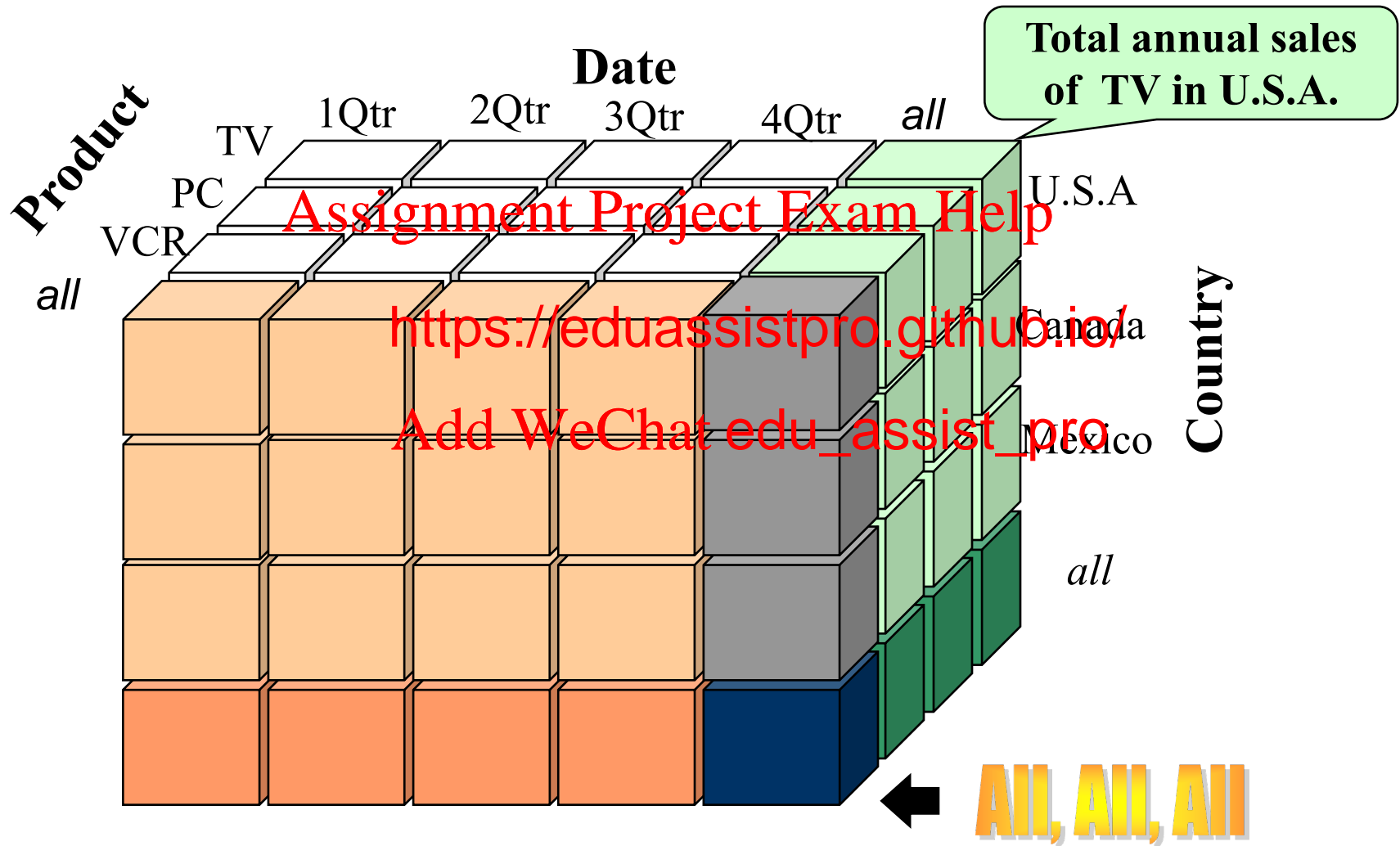
Assume: no other non-ALL levels on all dimensions.

All the Cuboids

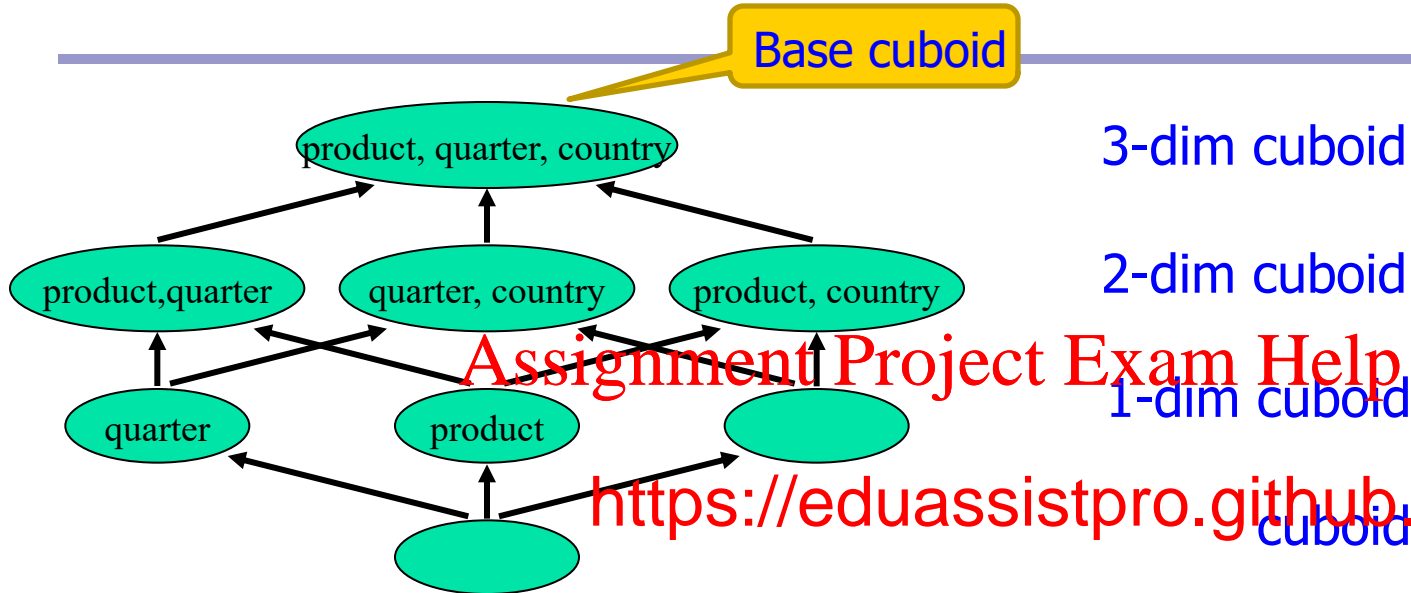


All the Cuboids /2

Assume: no other non-ALL levels on all dimensions.



Lattice of the cuboids



- n-dim cube can be represented as $(\prod_{i=1}^n D_i)$, where D_i is the set of allowed values on the i-th dimen
 - if $D_i = L_i$ (a particular level), then $D_i =$ all descendant dimension values of L_i .
 - ALL can be omitted and hence reduces the effective dimensionality
- A complete cube of d-dimensions consists of $\prod_{i=1}^d (n_i + 1)$ cuboids, where n_i is the number of levels (excluding ALL) on i-th dimension.
 - They collectively form a lattice.

Properties of Operations

- All operations are closed under the multidimensional model
 - i.e., both input and output of an operation is a cube
 - So that they
- <https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Q: What's the analogy in the Relational Model?