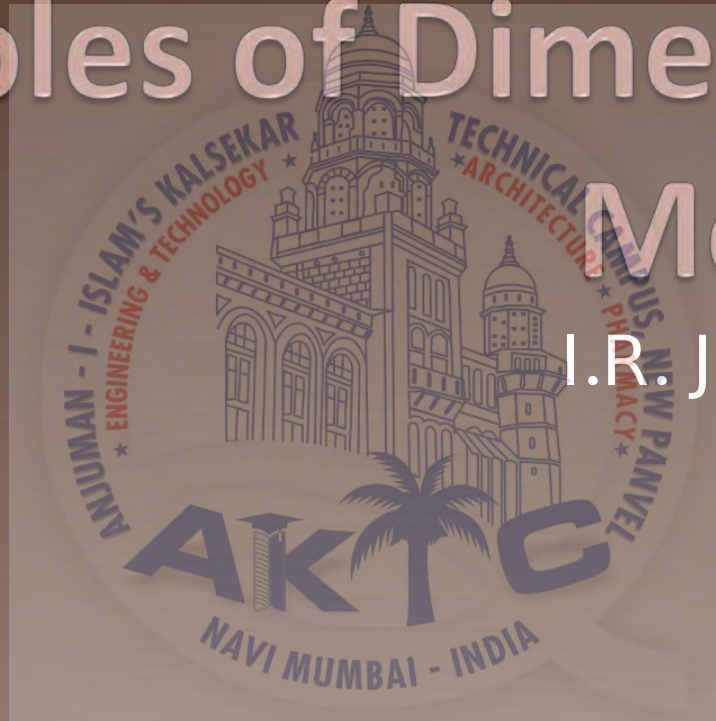# Sub: DWM        Sem : 8

Course Owner : Prof I.R. jamkhandikar

Academi Year : 2017/  18
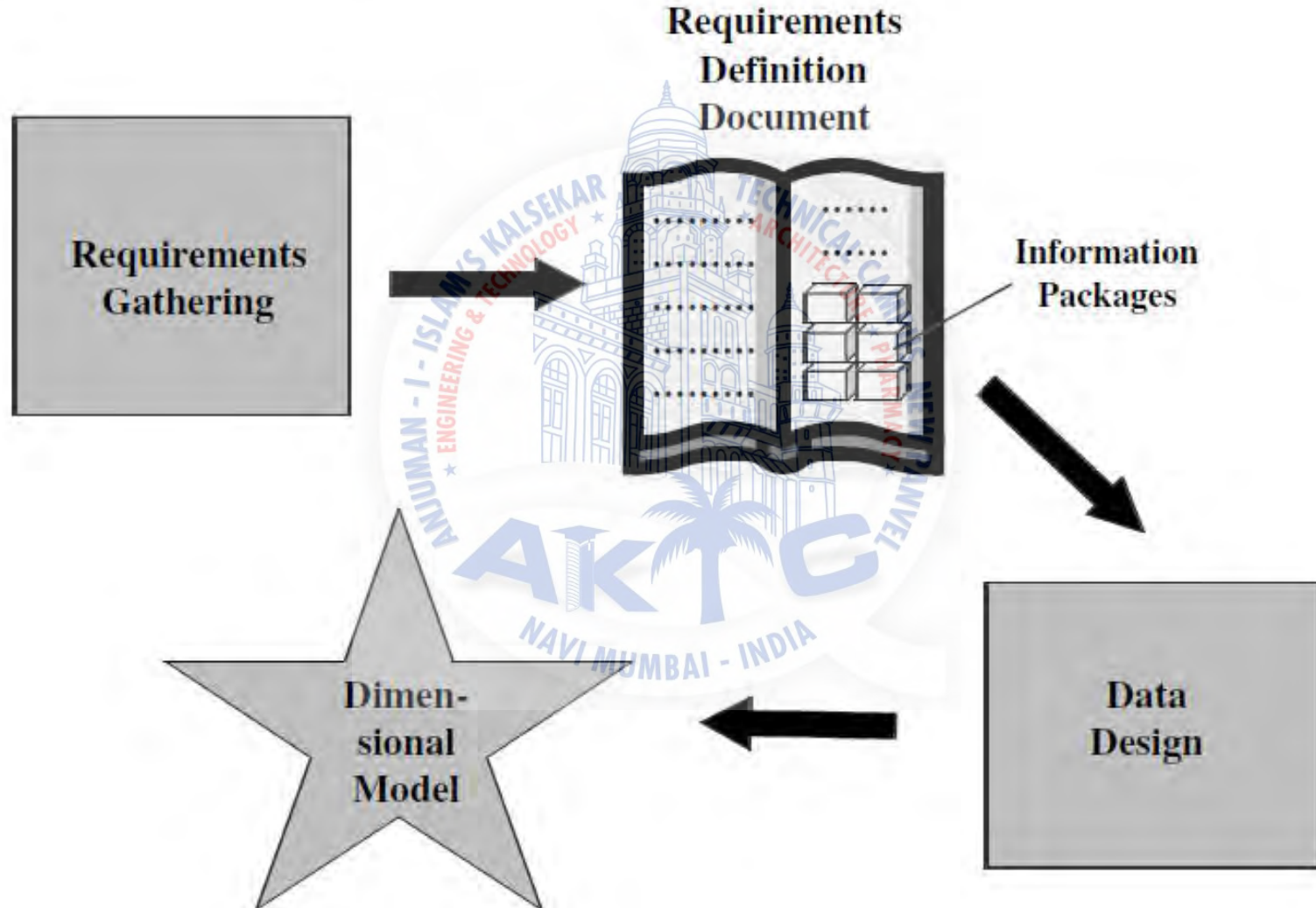
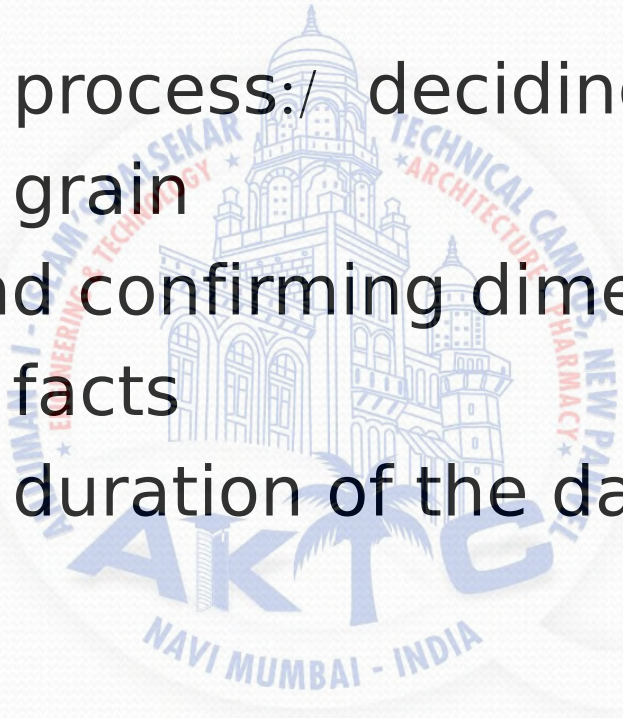# Principles of Dimensional Modeling

I.R. Jamkhandikar

# Objectives

- ಉ Understand how requirements definition determines data design
- ಉ Introduction of dimensional modeling /contrast with E/ R modeling
- ಉ Basics of star schema
- ಉ Contents of fact/dimension tables
- ಉ Advantages of star schema for DW

# Requirements to Design

# Design decisions to be taken

- ಖ Choosing the process:/ deciding subjects
- ಖ Choosing the grain
- ಖ Identifying and confirming dimensions
- ಖ Choosing the facts
- ಖ Choosing the duration of the database

# DImensional modeling basics

| Time | Product | Payment Method | Customer Demographics | Dealer | |
|---|---|---|---|---|---|
| Year | Model Name | Finance Type | Age | Dealer Name | |
| Quarter | Model Year | Term (Months) | Gender | City | |
| Month | Package Styling | Interest Rate | Income Range | State | |
| Date | Product Line | Agent | Marital Status | Single Brand Flag | |
| Day of Week | Product Category | | House-hold Size | Date First Operation | |
| Day of Month | Exterior Color | | Vehicles Owned | | |
| Season | Interior Color | | Home Value | | |
| Holiday Flag | First Year | | Own or Rent | | |

**Facts**: Actual Sale Price, MSRP Sale Price, Options Price, Full Price, Dealer Add-ons, Dealer Credits, Dealer Invoice, Down Payment, Proceeds, Finance

# Formation of the automaker sales fact table

## Dimensions

**Automaker Sales**

**Fact Table**

Actual Sale Price
MSRP Sale Price
Options Price
Full Price
Dealer Add-ons
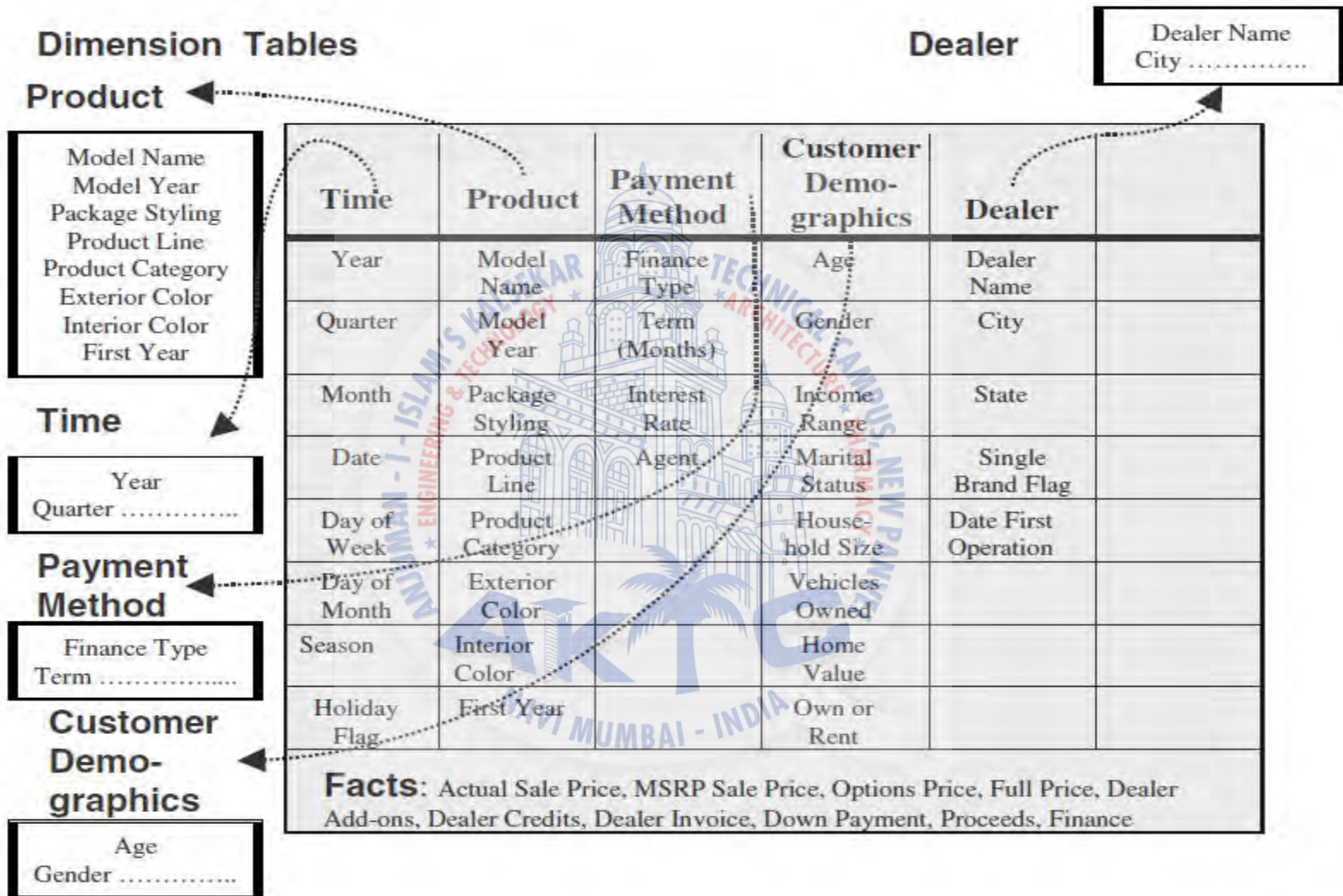Dealer Credits
Dealer Invoice
Down Payment
Proceeds
Finance

| Time | Product | Payment Method | Customer Demo-graphics | Dealer | |
|---|---|---|---|---|---|
| Year | Model Name | Finance Type | Age | Dealer Name | |
| Quarter | Model Year | Term (Months) | Gender | City | |
| Month | Package Styling | Interest Rate | Income Range | State | |
| Date | Product Line | Agent | Marital Status | Single Brand Flag | |
| Day of Week | Product Category | | House-hold Size | Date First Operation | |
| Day of Month | Exterior Color | | Vehicles Owned | | |
| Season | Interior Color | | Home Value | | |
| Holiday Flag | First Year | | Own or Rent | | |

**Facts**: Actual Sale Price, MSRP Sale Price, Options Price, Full Price, Dealer Add-ons, Dealer Credits, Dealer Invoice, Down Payment, Proceeds, Finance
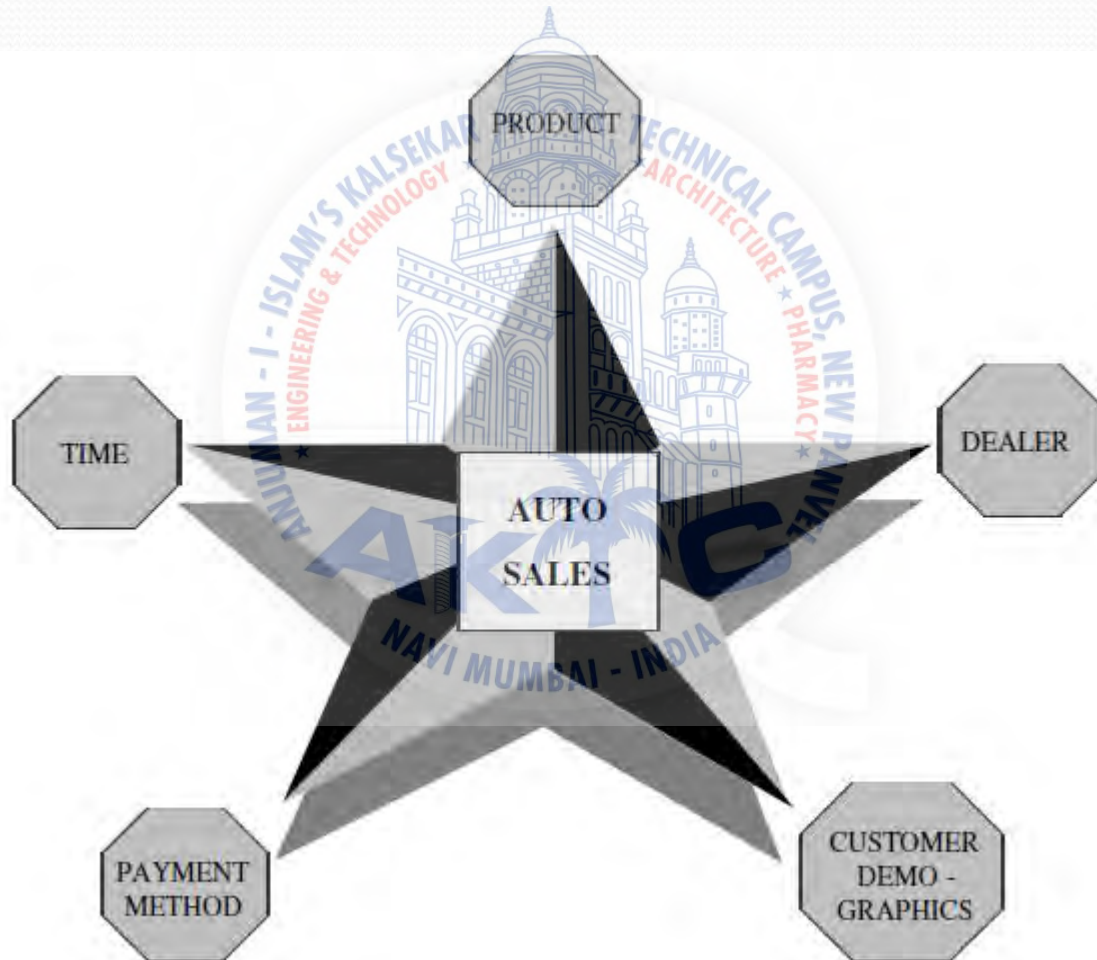
# Formation of the automaker dimension tables



**Dimension Tables**                                    **Dealer**

**Product**

| Model Name, Model Year, Package Styling, Product Line, Product Category, Exterior Color, Interior Color, First Year |

**Time**

| Year, Quarter ............... |

**Payment Method**

| Finance Type, Term ............... |

**Customer Demo-graphics**

| Age, Gender ............... |

| Dealer Name, City ............... |

| | Time | Product | Payment Method | Customer Demo-graphics | Dealer | |
|---|---|---|---|---|---|---|
| | Year | Model Name | Finance Type | Age | Dealer Name | |
| | Quarter | Model Year | Term (Months) | Gender | City | |
| | Month | Package Styling | Interest Rate | Income Range | State | |
| | Date | Product Line | Agent | Marital Status | Single Brand Flag | |
| | Day of Week | Product Category | | House-hold Size | Date First Operation | |
| | Day of Month | Exterior Color | | Vehicles Owned | | |
| | Season | Interior Color | | Home Value | | |
| | Holiday Flag | First Year | | Own or Rent | | |

**Facts**: Actual Sale Price, MSRP Sale Price, Options Price, Full Price, Dealer Add-ons, Dealer Credits, Dealer Invoice, Down Payment, Proceeds, Finance

How much sales proceeds did the jeep tata mahindra, 2009 model with vxi  options, generate in january 2000 at spectra  auto dealership for buyers who owned their homes, financed by icici prudential financing?

# Tips for combining data into dimensional model

- Provide best data access
- Model should be query/ centric
- Model should be optimized for queries and analyses
- Model should reveal the interactions between the dimension and fact tables
- There should be drilling down or rolling up along dimension hierarchies

# STAR SCHEMA for automaker sales

# ER Model v/s Dimension Model

- ಬ ER diagram is a complex diagram, used to represent multiple processes. A single ER diagram can be broken down into several DM diagrams.
- ಬ In DM, we prefer keeping the tables de/ normalized, whereas in a ER diagram, our main aim is to remove redundancy
- ಬ ER model is designed to express microscopic relationships between elements. DM captures the business measures
- ಬ DM is designed to answer queries on business process, whereas the ER model is designed to record the business processes via their transactions.

# Entity-Relationship vs. Dimensional Models

## E- R DIAGRAM

- ౪ One table per entity
- ౪ Minimize data redundancy
- ౪ Optimize update
- ౪ The Transaction Processing Model

## DIMENSIONAL MODEL

- ౪ One fact table for data organization
- ౪ Maximize understandability
- ౪ Optimized for retrieval
- ౪ The data warehousing model

# Star Schema-example of order analysis

Query result

# Understanding query from the star schema

# Understanding drill down analysis from the star schema

# Dimension table

ꙮ Contain information about a particular dimension.

  ꙮ Dimension table key

  ꙮ Table is wide

  ꙮ Textual attributes

  ꙮ Attributes not directly related

  ꙮ Not normalized

  ꙮ Drilling down, rolling up

  ꙮ Multiple hierarchies

  ꙮ Fewer number of records

# Facts

- ஃ Numeric measurements (values) that represent a specific business aspect or activity
- ஃ Stored in a fact table at the center of the star scheme
- ஃ Contains facts that are linked through their dimensions
- ஃ Can be computed or derived at run time
- ஃ Updated periodically with data from operational databases

# Fact table

ಎ Contains primary information of the warehouse

- ಎ Concatenated key
- ಎ Data grain
- ಎ Fully additive measures
- ಎ Semi/ additive measures(derived attributes)
- ಎ Table deep, not wide
- ಎ Sparse data
- ಎ Degenerate dimensions(attributes which are neither fact or a dimension)

# Star schema for a retail chain

**Time Dimension Table**
- Time key
- Year
- Quarter
- Month
- Week
- Date

**Sales Fact Table**
- Time key
- Product key
- Customer key
- Store key
- Mode key
- Actual sales
- Forecast sales
- Price
- Discount

**Customer Dimension Table**
- Customer key
- Name
- Age
- Income
- Gender
- Marital status

**Store Dimension Table**
- Store key
- Name
- City
- State
- Op from year
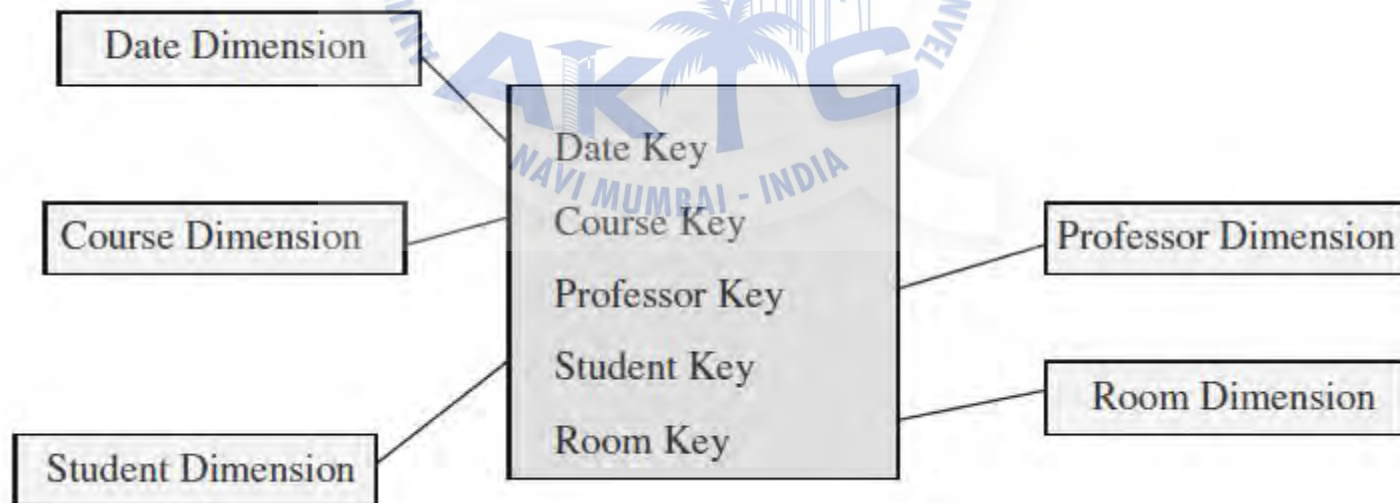
**Payment Mode Dimension Table**
- Mode key
- Payment mode

**Product Dimension Table**
- Product key
- Name
- Brand
- Category
- Colour
- Price

# Star Schema characteristics

- Star schema is a relational model with one/ to/ many relationship between the fact table and the dimension tables.
- De/ normalized relational model
- Easy to understand. Reflects how users think. This makes it easy for them to query and analyse the data.
- Optimizes navigation.
- Enhances query extraction.
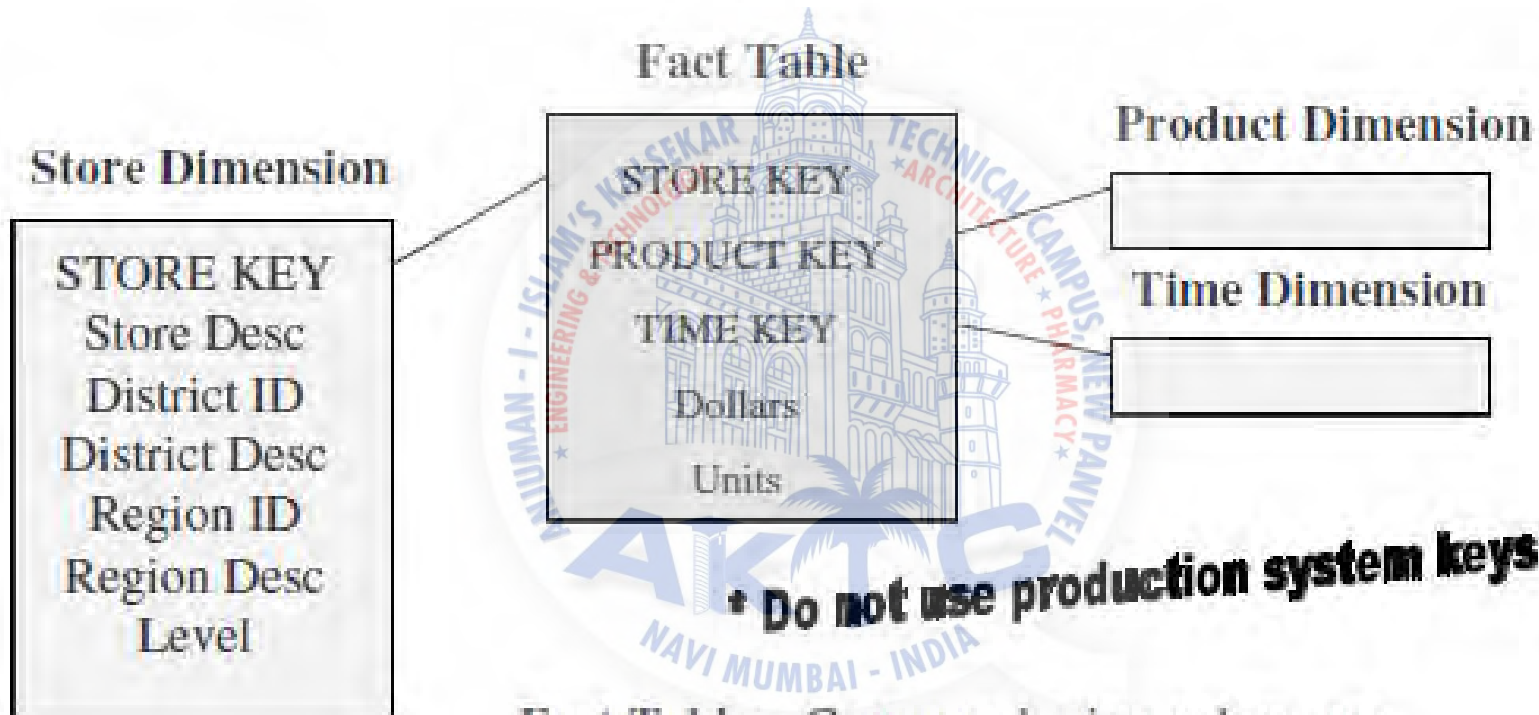- Ability to drill down or roll up.

# Factless fact table

ಐ A fact table is said to be empty if it has no measures to be displayed. Fact table represents **events**

ಐ Contains no data, only keys.

# Data Granularity

ಉ When fact table at the lowest grain, the users can as well drill down to the lowest grain of details

ಉ But when data is kept till the lowest level of data, we have to compromise on the storage and maintenance of DW

ಉ Advantages

  ಪ Easier to extract from operational data and load into DW

  ಪ Can be feed directly to the DM application

# Star Scheme Keys

**Fact Table**

**Store Dimension**

**Product Dimension**

| STORE KEY |
|-----------|
| Store Desc |
| District ID |
| District Desc |
| Region ID |
| Region Desc |
| Level |

| STORE KEY |
|-----------|
| PRODUCT KEY |
| TIME KEY |
| Dollars |
| Units |

**Time Dimension**

* Do not use production system keys *

**Fact Table:**   Compound primary key, one segment for each dimension

**Dimension Table:**   Generated primary key

# Star schema keys contd...

- ಉ Primary keys: should not be same as production system
- ಉ Surrogate keys: System generated sequence numbers having no built/ in meanings
- ಉ Foreign keys: primary key of each dimension table must be a foreign key in the fact table.

# Primary key for Fact table

- ಜ A single compound primary key whose length is the total length of the keys of the individual dimension tables

- ಜ **Concatenated primary key that is the concatenation of all the primary keys of the dimension tables,**

- ಜ A generated primary key independent of the keys of the dimension tables.

# Advantages of the star schema

- ಉ Easy for users to understand
- ಉ Optimizes navigation
- ಉ Most suitable for query processing

# Starjoin and Starindex

ಊ Star join:/ high/ speed, single pass parallelizable, multitable join.

  ಊ Boots query performance

ಊ Star index:/ specialized index to accelerate join performance

  ಊ Speed up joins between the dimension tables and fact tables

# Summing up

ಹ Derived from the information packages in the requirements definition.

ಹ The STAR schema used for data design is a relational model consisting of fact and dimension tables.

ಹ The fact table contains the business metrics or measurements; the dimensional tables contain the business dimensions. Hierarchies within each dimension table are used for drilling down to lower levels of data.

ಹ STAR schema advantages are: easy for users to understand, optimizes navigation, most suitable for query processing, and enables specific

# Objectives

ಉ Slowly changing dimensions

ಉ Large dimensions

ಉ Snowflake schema

ಉ Aggregate tables

ಉ Family of starts and their applications

# Updating the Dimension table

ಞ Dimension tables are non⁄ volatile and mostly read⁄ only.

ಞ More rows are added to the Dimension tables over time.

ಞ Changes to certain attributes of a row become eminent at times.

ಞ There are many types of changes that affect the dimension tables.

# Slowly changing dimensions

ಲ Most dimensions are generally constant over time

ಲ Many dimensions change slowly

ಲ Though the key does not change other description and attributes change slowly over time

ಲ Dimension table attributes are not overwritten

ಲ The ways changes are made in dimension tables depend on the types of changes and what information must be preserved.

# Type 1: Correction of errors

- ಖ Usually relate to correction of errors in the source systems.
- ಖ E.g., spelling error in customer names; change of names of customers;
- ಖ There is no need to preserve the old values here.
- ಖ The old value in the source system needs to be discarded.
- ಖ The changes made need not be preserved or noted.

# Type 1.. continued

ಜ Overwrite attribute value in the dimension table row with new value

ಜ No other changes are made to the dimension table row.

ಜ The key is not disturbed

Easiest type of change to implement



**KEY RESTRUCTURING**

33154112 ← K12356

**INCREMENTAL LOAD -- TYPE 1 CHANGE**

Customer Code: K12356

Customer Name: Kristin Samuelson

**BEFORE**

| Customer Key: | 33154112 |
| Customer Name: | Kristin Daniels |
| Customer Code: | K12356 |
| Marital Status: | Single |
| Address: | 733 Jackie Lane, Baldwin Harbor |
| State: | NY |
| Zip: | 11510 |

**AFTER**

33154112

Kristin Samuelson

K12356

Single

733 Jackie Lane, Baldwin Harbor

NY

11510

# Type 2 : preservation of history

- True changes in the source systems.
- E.g., change of marital status; change of address
- There is a need to preserve history
- This type of changes partition the warehouse
- Every change for the same attribute must be preserved.
- Applying these changes:
  - Add a new dimension table row with new value of the changed attribute
  - No changes are made to the existing row.
  - New rows are inserted with a new surrogate key.

# Type 2.. continued

# Type 3: tentative soft revision

- ಕ Tentative changes in the source system
- ಕ E.g., if an employee will get posted for a short period to a different location
- ಕ Need to keep track of history with old and new values
- ಕ Used to compare performances across the transition
- ಕ Applying these changes
  - ಕ An "old" field is added in the dimension table
  - ಕ Push existing value of attribute from "current" to "old"
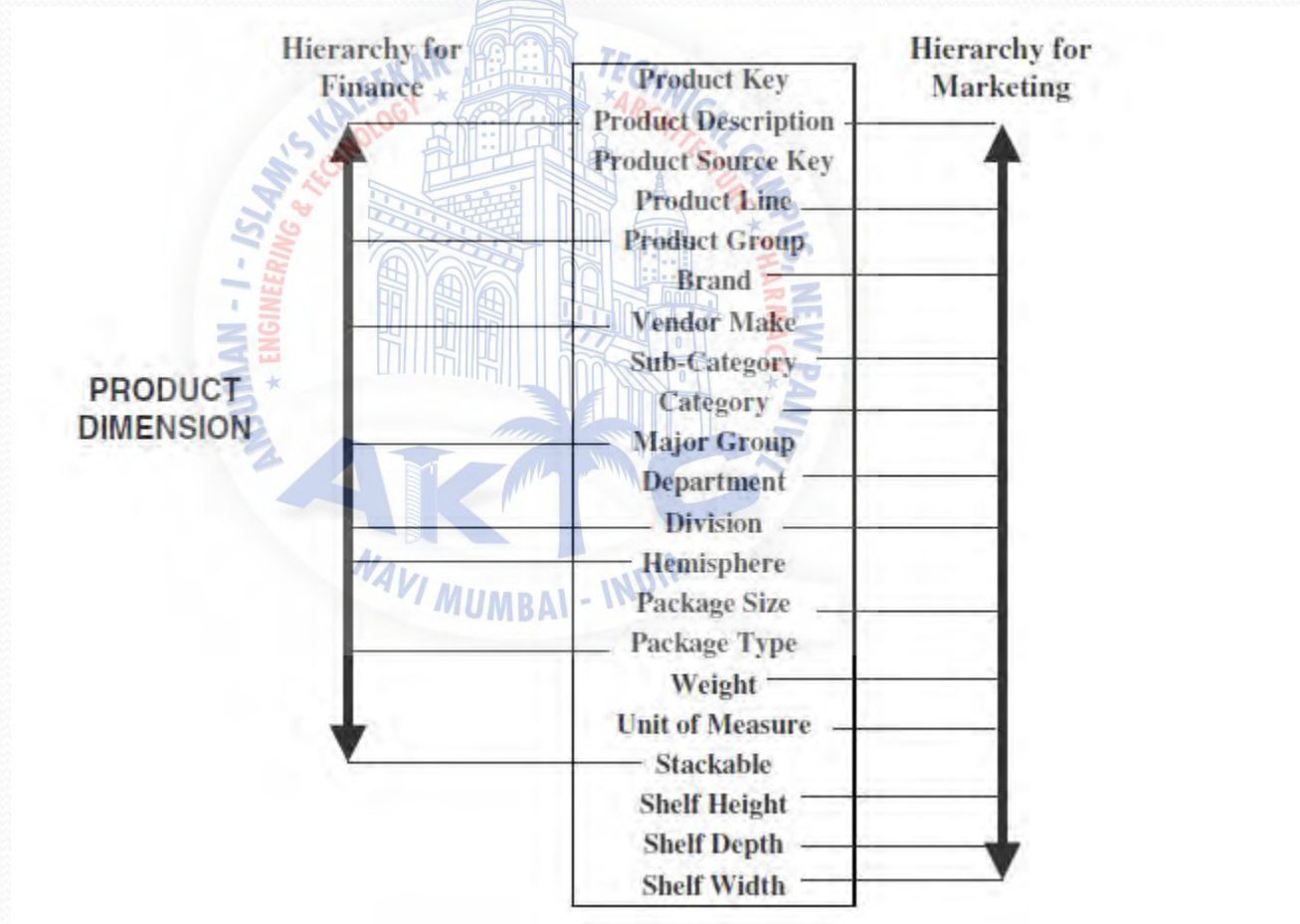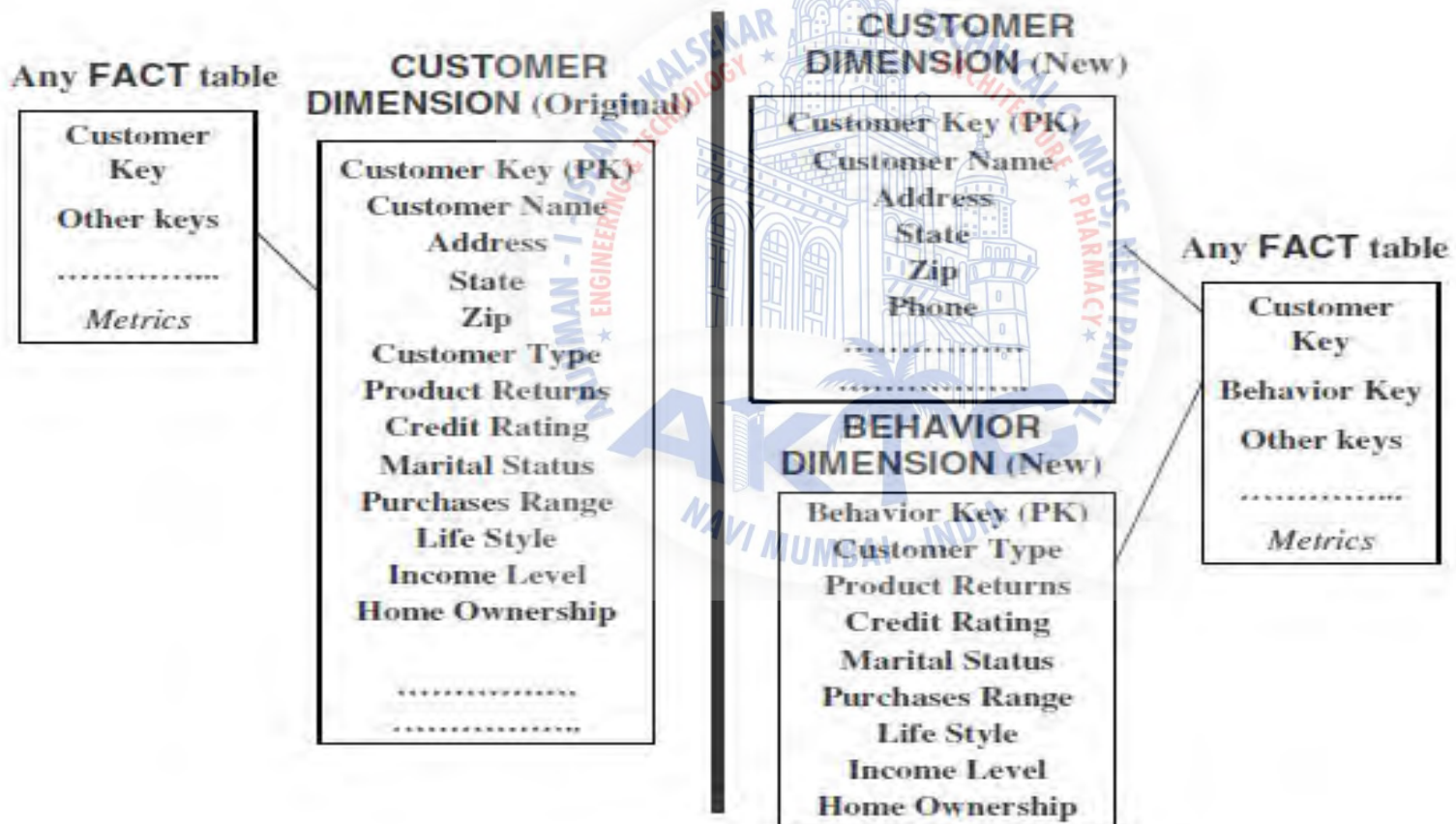
# Type 3.. continued

# Large dimensions

- ಉ Very deep(large number of rows)
- ಉ Very wide(large number of attributes)
- ಉ Have multiple hierarchies
- ಉ Rapidly changing dimensions
- ಉ Junk dimensions
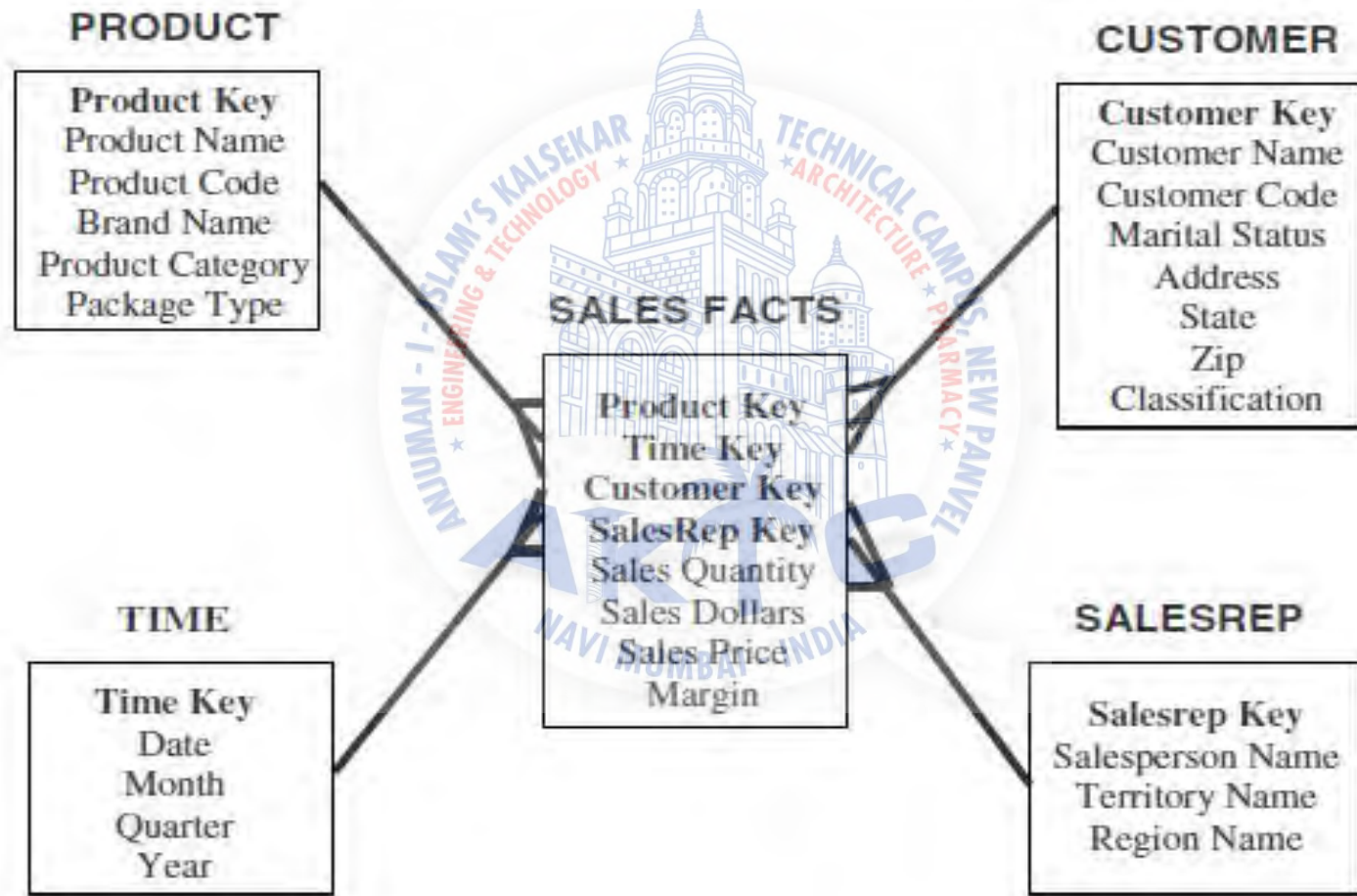
# Multiple hierarchies
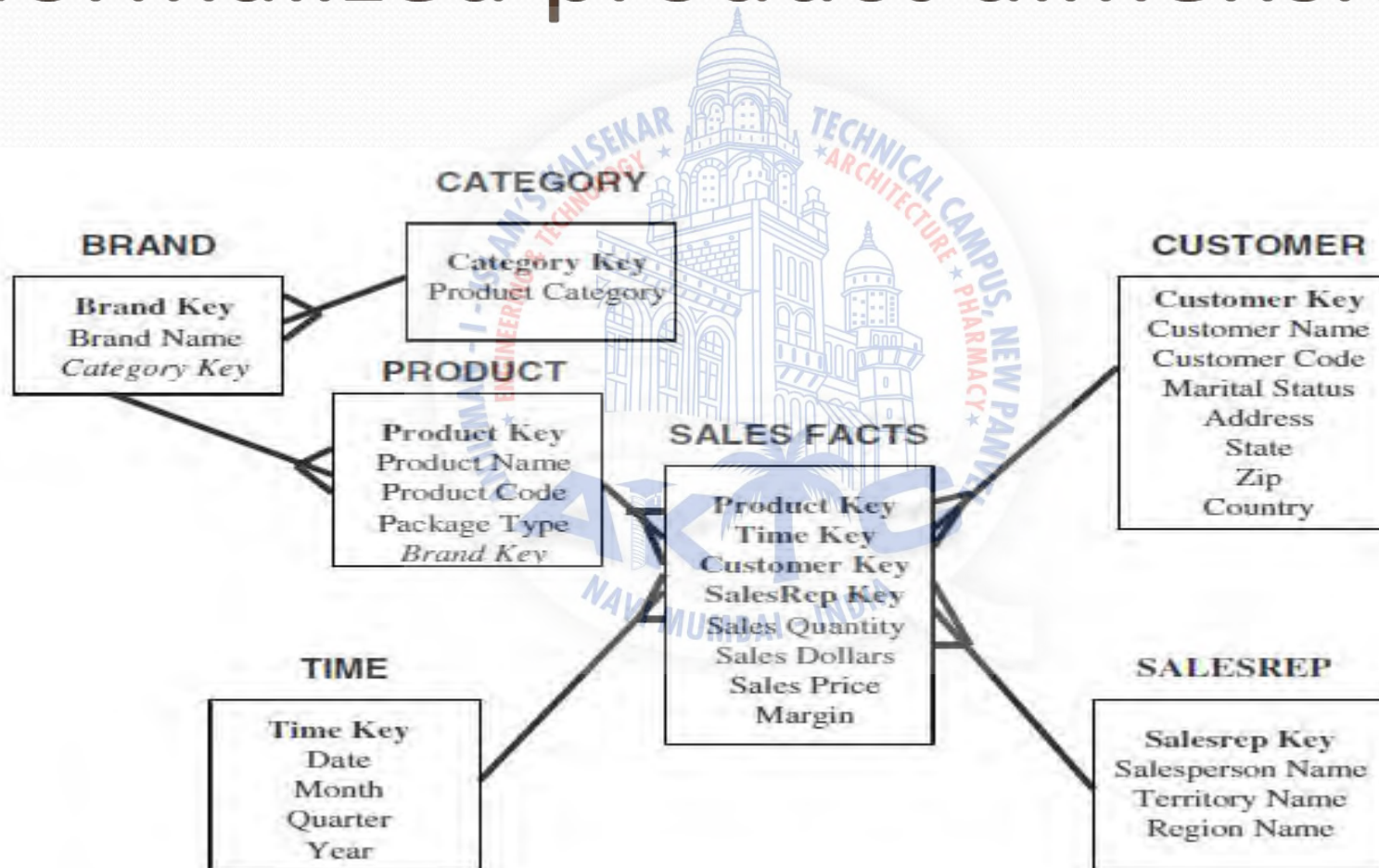
# Rapidly changing dimensions

# Snowflake schema

- ಉ A variation of the star schema, in which all or some of the dimension tables may be normalized.

- ಉ Eliminates redundancy

- ಉ Generally used when a dimension table is wide.

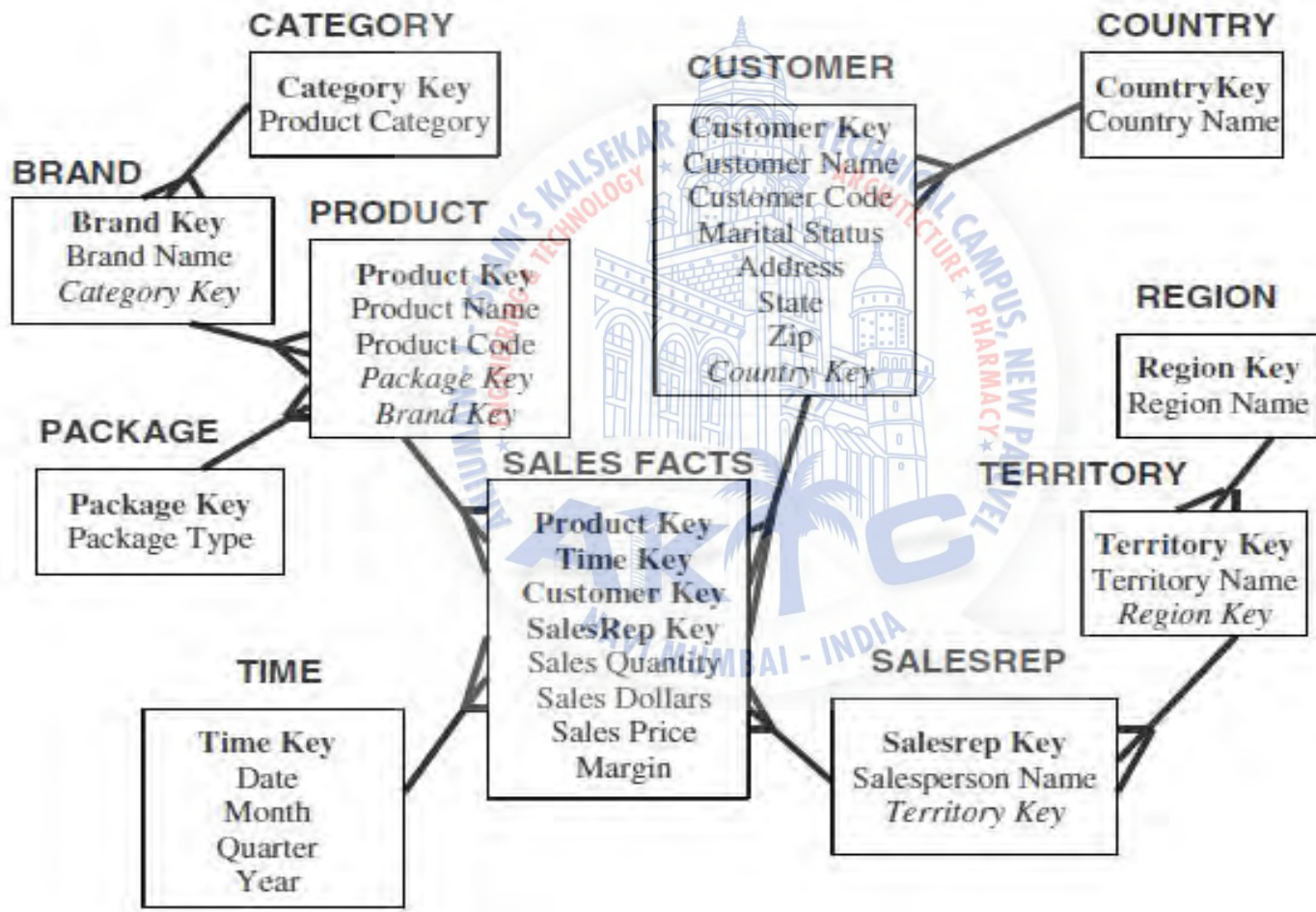- ಉ Saves space

- ಉ Complex querying is required.

# Star schema for sales

# Normalized product dimension

# Sales snowflake schema
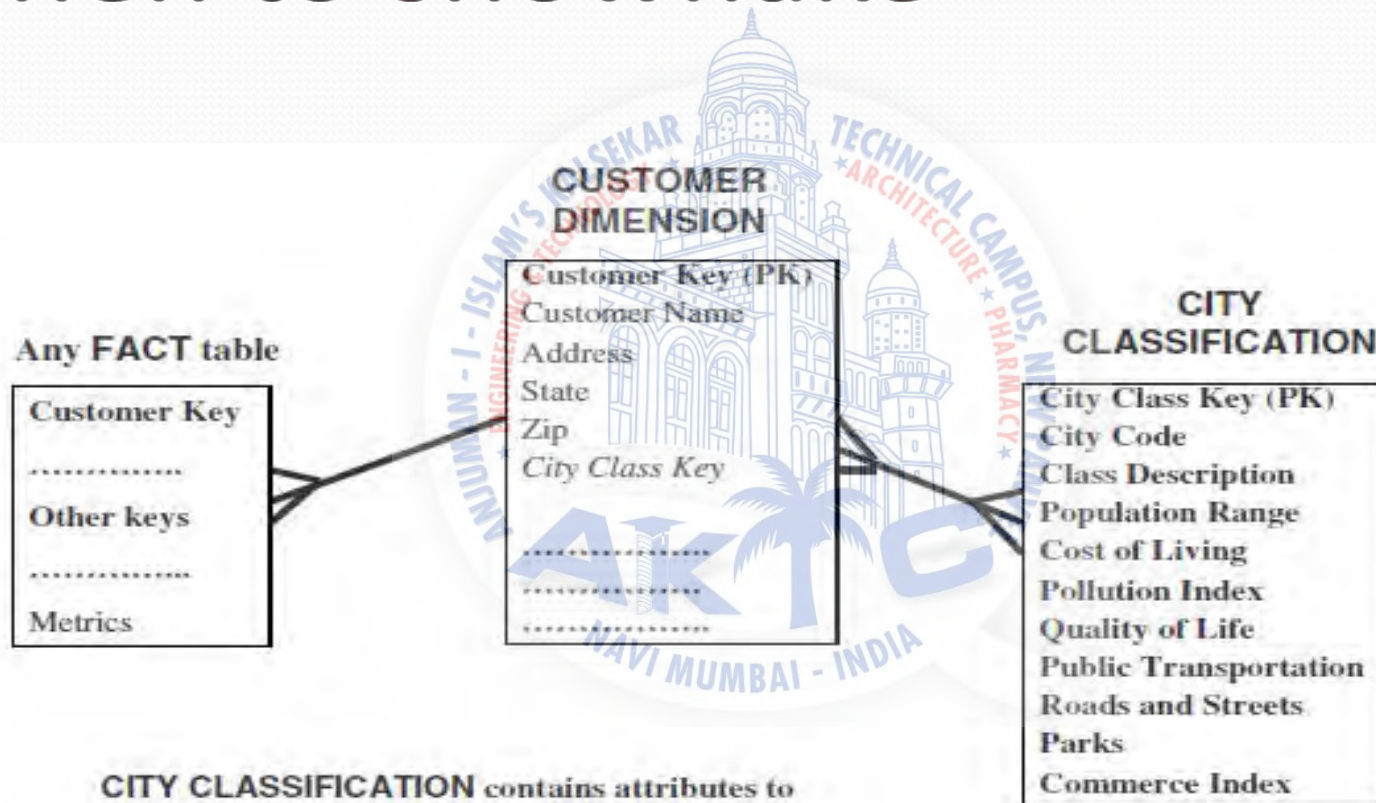
# Advantages and disadvantages

- Advantages
  - Small savings in storage space
  - Normalized structures are easier to update and maintain
- Disadvantages
  - Schema is less intuitive
  - Browsing becomes difficult
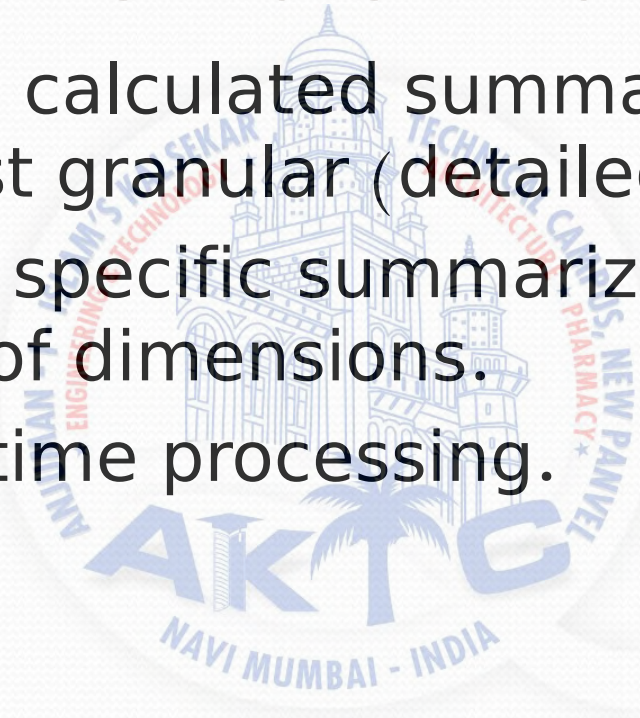  - Degraded query performance because of additional joins

# When to snowflake



**CUSTOMER DIMENSION**
- Customer Key (PK)
- Customer Name
- Address
- State
- Zip
- *City Class Key*

**Any FACT table**
- Customer Key
- ..............
- Other keys
- ..............
- Metrics

**CITY CLASSIFICATION**
- City Class Key (PK)
- City Code
- Class Description
- Population Range
- Cost of Living
- Pollution Index
- Quality of Life
- Public Transportation
- Roads and Streets
- Parks
- Commerce Index

**CITY CLASSIFICATION** contains attributes to classify each city within a limited set of classes. These attributes are separated from the **CUSTOMER DIMENSION** to form a separate sub-dimension as **CITY CLASSIFICATION**.

# Aggregate fact tables

- ಉ Contain pre/ calculated summaries derived from the most granular (detailed) fact table.
- ಉ Created as a specific summarization across any number of dimensions.
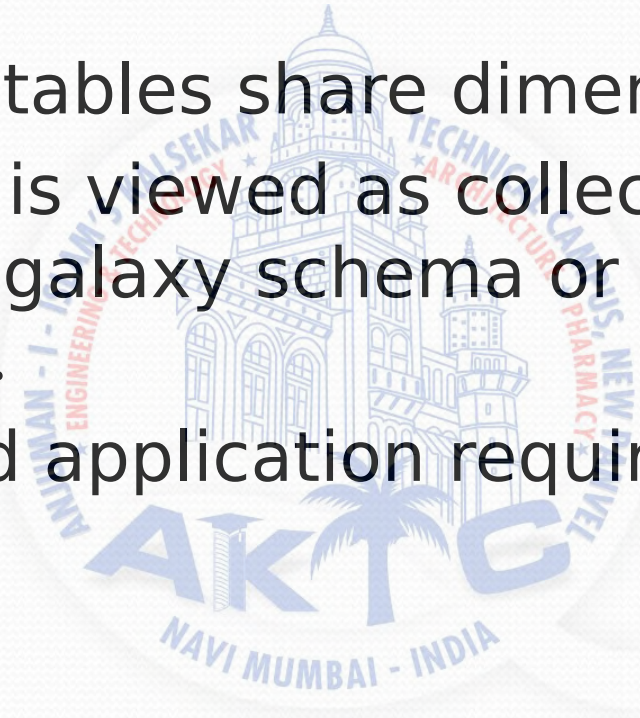- ಉ Reduces runtime processing.

# Why need aggregate fact tables?

- Large size of the fact table
- To speed up query extraction

- Limitations
  - Must be re/ aggregated each time there is a change in the source data
  - Do not support exploratory analysis
  - Limited interactive use.

# Fact Constellation

- Multiple fact tables share dimension tables.
- This schema is viewed as collection of stars hence called galaxy schema or fact constellation.
- Sophisticated application requires such schema.

# Fact Constellation (contd..)

**Sales Fact Table**

| Store Key |
|-----------|
| Product Key |
| Period Key |
| Units |
| Price |

**Product Dimension**

| Product Key |
|-----------|
| Product Desc |

**Shipping Fact Table**

| Shipper Key |
|-----------|
| Store Key |
| Product Key |
| Period Key |
| Units |
| Price |

**Store Dimension**

| Store Key |
|-----------|
| Store Name |
| City |
| State |
| Region |

# Fact Constellation

ಉ Multiple fact tables share dimension tables.

ಉ This schema is viewed as collection of stars hence called galaxy schema or fact constellation.

ಉ Sophisticated application requires such schema.

# Fact Constellation (contd..)

**Sales
Fact Table**

**Shipping
Fact Table**

**Product
Dimension**

| Store Key |
|---|
| Product Key |
| Period Key |
| Units |
| Price |

| Product Key |
|---|
| Product Desc |

| Shipper Key |
|---|
| Store Key |
| Product Key |
| Period Key |
| Units |
| Price |

**Store
Dimension**

| Store Key |
|---|
| Store Name |
| City |
| State |
| Region |