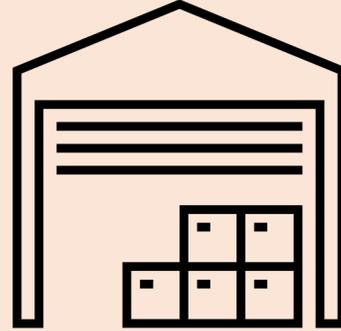*To cite this presentation: Pendyala, V.S. (2023) "Slowly changing dimensions and fast changing facts - the story of the traditional Datawarehouse". PPT Presentation. IEEE Computer Society, Santa Clara Valley Chapter Webinar.*

# Slowly changing dimensions and fast changing facts - the story of the traditional Datawarehouse

Vishnu S. Pendyala, PhD

*Video Recording:*
*https://www.youtube.com/watch?v=BcttdNrbBhk*

# How do you generate actionable insights into an organization's performance?

## The beginnings

- Inmon, W. H. (1992). Building the Data Warehouse. Wiley.
- Codd, E. F., Codd, S. B., & Salley, C. T. (1993). Providing OLAP (on-line analytical processing) to user-analysts. *An IT Mandate. White Paper. Arbor Software Corporation*, *4*.
- Kimball, R. (1996). The Data Warehouse Toolkit: The Definitive Guide to Dimensional Modeling (1st ed.). Wiley.
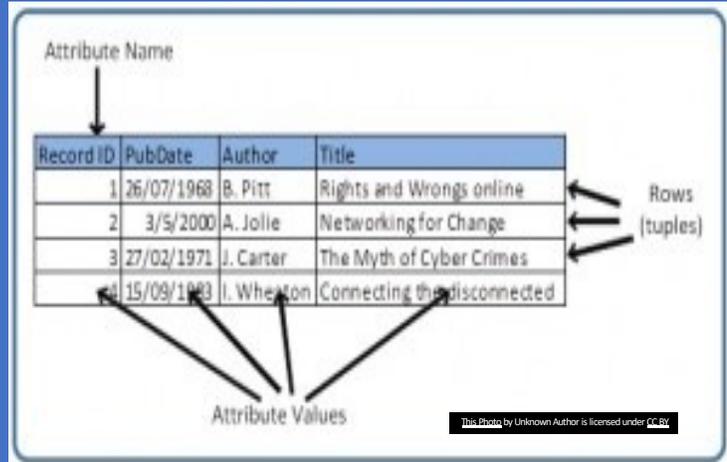
*Example of a business insight:*

Company ABC made $234,567
from the sale of
13" MacBookPro to customer XYZ

How do we generate this?
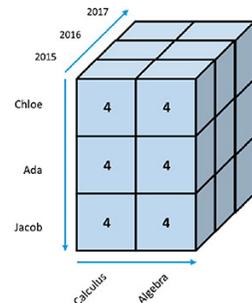
## Take 1: Generate insights from existing databases

**What is the problem?**

## Take 2: Generate the aggregates during the night when systems are not busy and store them in efficient data structures (data cube)

Dimensions of the cube: Subject, Student, Year
Numbers in the subcubes: GPA

*Operations: Slice, Dice,*
*Roll-up, pivot, etc*

Business insight:
Company ABC made $234,567
from the sale of
13" MacBookPro to customer XYZ

*Fact*
*(fast changing, hopefully)*

*Dimension*
*(slowly changing, hopefully)*

*Dimension*
*(slowly changing, hopefully)*

# There's still a problem!

Insights require data from many sources, some of which may not even be databases (files, spreadsheets, etc) => need to Extract, Transform, and Load

Need to preserve historical data with a specific need to find temporal patterns

Performance does not scale as analytical processing needs grow; transaction processing cannot take a hit

Normalization can come in the way of efficient analytical query processing => need for a different way to model data

=> Cannot mix OLTP (operational) and OLAP (informational) systems!

Solution

STAR SCHEMA

Time dimension Table

**TIME**

| PK | Time_ID |
| --- | --- |
| | Day |
| | Month |
| | Quarter |
| | Year |

Item Dimension Table

**ITEM**

| PK | Item_Key |
| --- | --- |
| | Item_Name |
| | Brand |
| | Sold_By |
| | Category |

Sales Fact Table

| FK | ITEM_KEY |
| --- | --- |
| FK | TIME_ID |
| FK | BRANCH_ID |
| FK | LOCATION_ID |
| | QTY_SOLD |
| | AMT_SOLD |
| | AVG_SALES |

Location dimension Table

**LOCATION**

| PK | Location_ID |
| --- | --- |
| | Name |
| | State |
| | Pincode |

Branch dimension Table

**BRANCH**

| PK | Branch_ID |
| --- | --- |
| | Branch_NM |
| | Owner |

**An Example**

# Fact Table: Sales

| sale_id | product_id | time_id | location_id | branch_id | quantity_sold | amount |
|---------|------------|---------|-------------|-----------|---------------|--------|
| 1 | 101 | 1 | 1 | 1 | 5 | 150 |
| 2 | 102 | 2 | 2 | 2 | 3 | 90 |
| 3 | 101 | 3 | 1 | 1 | 2 | 60 |
| 4 | 103 | 4 | 3 | 3 | 7 | 210 |
| 5 | 102 | 5 | 2 | 2 | 4 | 120 |

Facts are fast changing
=> Fact tables are very long
=> Typically, billions of rows
=> Fact tables need to be lean (for storage considerations)
=> Fact tables contain only numbers

Dimensions are descriptive
=> They are wide
=> Mostly contain strings and some numbers
=> Less number of rows compared to fact table
=> Short and stout!

## Dimension Table: Time

| time_id | date | day_name |
|---|---|---|
| 1 | 2023-05-10 | Monday |
| 2 | 2023-05-11 | Tuesday |
| 3 | 2023-05-12 | Wednesday |
| 4 | 2023-05-13 | Thursday |
| 5 | 2023-05-14 | Friday |

## Dimension Table: Location

| location_id | location_name |
|---|---|
| 1 | City A |
| 2 | City B |

## Dimension Table: Product

| product_id | product_name | category_id |
|------------|--------------|-------------|
| 101 | Laptop | 1 |
| 102 | Smartphone | 2 |
| 103 | Tablet | 1 |

## Dimension Table: Branch

| branch_id | branch_name |
|-----------|-------------|
| 1 | Store A |
| 2 | Store B |
| 3 | Store C |

## What are the total sales amount for each product in City A during the month of May 2023?

```
SELECT p.product_name, SUM(s.amount) AS
total_sales_amount
FROM Sales s
JOIN Product p ON s.product_id = p.product_id
JOIN Location l ON s.location_id = l.location_id
JOIN Time t ON s.time_id = t.time_id
WHERE l.location_name = 'City A'
  AND t.date BETWEEN '2023-05-01' AND '2023-05-31'
GROUP BY p.product_name;
```

## Snowflake schema

DIMENSION

DIMENSION

FACT TABLE

DIMENSION

DIMENSION

Dimensional Modeling

## Key considerations

| | | |
|---|---|---|
| | What are the dimensions and how many of them? | I/O bound queries => minimize storage |

| | |
|---|---|
| | What numeric quantities will be stored in the Fact table? |

| | |
|---|---|
| | At what granularity will be facts be captured? | Tradeoff: finer grain => better detail for drilling down, but more rows |

# An example

**Suppose there are 60 dimensional attributes and 3 facts to be captured in the fact table**

**How do we design the star schema? How many dimension tables? Choices:**

A. One dimension with all 60 attributes
B. Each attribute gets its own dimension table
C. Four dimensions

# A: One dimension with all 60 attributes

| FK | Fact 1 | Fact 2 | Fact 3 |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

| PK | Attr 1 | Attr 2 | … | Attr 60 |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# B: Each attribute gets its own dimension table

| FK 1 | FK … | FK 60 | Fact 1 | Fact 2 | Fact 3 |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

| PK | Attr 1 |
|---|---|
| | |
| | |
| | |

| PK | Attr 2 |
|---|---|
| | |
| | |

| PK | Attr … |
|---|---|
| | |
| | |
| | |

| PK | Attr 1 |
|---|---|
| | |
| | |
| | |

| PK | Attr … |
|---|---|
| | |
| | |
| | |

| K | Attr … |
|---|---|
| | |

| PK | Attr 1 |
|---|---|
| | |
| | |
| | |

| PK | Attr 60 |
|---|---|
| | |
| | |

# C: Four dimensions

| F K 1 | F K ... | F K 4 | Fact 1 | Fact 2 | Fa ct 3 |
|---|---|---|---|---|---|
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |
|  |  |  |  |  |  |

| P K | Date | Attr | ... | Attr |
|---|---|---|---|---|

| P K | Pro duct | Attr | ... | Attr |
|---|---|---|---|---|

| P K | Stor e | Attr | ... | Attr |
|---|---|---|---|---|

| P K | Pro mo | Attr | ... | Attr |
|---|---|---|---|---|

**Which one fares better?**

- Option A is the worst – dimension can have as many rows as the fact table!
- Option B – dimensions occupy less storage, but the fact tables will be very wide and occupy huge storage, still less than in Option A
- Option C is the best – do the math!

# Slowly Changing Dimensions

SCD Type 1, 2, and 3
and also
Type 0, 4, 5, …

## How do we handle rarely occurring changes to dimensions?

- Do nothing – ignore the changes (type 0)

- Overwrite the existing values – erase history (type 1)!

- Find ways to preserve the history
  - Type 2, 3, 4, …

# SCD Type 1: In-Place Update

**Initial Customer Dimension**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | New York | Gold |
| 102 | Jane Doe | Los Angeles | Silver |

**Requested Customer update**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | Boston | Platinum |

**Updated Customer Dimension**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | Boston | Platinum |
| 102 | Jane Doe | Los Angeles | Silver |

# SCD Type 2: Historical Tracking

**Initial Customer Dimension**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | New York | Gold |
| 102 | Jane Doe | Los Angeles | Silver |

**Requested Customer update**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | New York | Platinum |

**Updated Customer Dimension**

| customer_id | customer_name | city | loyalty_status | effective_date | end_date |
|---|---|---|---|---|---|
| 101 | John Smith | New York | Gold | 2023-01-01 | 2023-08-10 |
| 101 | John Smith | New York | Platinum | 2023-08-11 | (null) |
| 102 | Jane Doe | Los Angeles | Silver | (null) | (null) |

# SCD Type 3: Alternate Reality

**Initial Customer Dimension**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | New York | Gold |
| 102 | Jane Doe | Los Angeles | Silver |

**Requested Customer update**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | New York | Platinum |

**Updated Customer Dimension**

| customer_id | customer_name | city | loyalty_status | prev_loyalty_status |
|---|---|---|---|---|
| 101 | John Smith | New York | Platinum | Gold |
| 102 | Jane Doe | Los Angeles | Silver | (null) |

# SCD Type 4: Monster / Mini Dimension

**Initial Customer Dimension**

| customer_id | customer_name | city | loyalty_status |
|---|---|---|---|
| 101 | John Smith | New York | Gold |
| 102 | Jane Doe | Los Angeles | Silver |

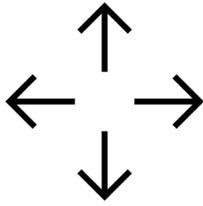**Loyalty_status is changing rapidly => create a new mini dimension for it**

| customer_id | customer_name | city |
|---|---|---|
| 101 | John Smith | New York |
| 102 | Jane Doe | Los Angeles |

| customer_id | loyalty_status | start_date | end_date |
|---|---|---|---|
| 101 | Gold | 2023-01-01 | 2023-08-10 |

**Updated Mini Dimension**

| customer_id | loyalty_status | start_date | end_date |
|---|---|---|---|
| 101 | Gold | 2023-01-01 | 2023-08-10 |
| 101 | Platinum | 2023-08-11 | (null) |

# More SCD types

- Hybrid approaches:
  - Type 5: Add mini-dimension (Type 4) and Type 1
  - Type 6: Type 1 + Type 2 + Type 3
  –

SAN JOSE STATE UNIVERSITY

*https://www.sjsu.edu/people/vishnu.pendyala/*
*@vishnupendyala*