

# Assignment No. 1

## Title:

For an organization of your choice, choose a set of business processes. Design star / snow flake schemas for analyzing these processes. Create a fact constellation schema by combining them. Extract data from different data sources, apply suitable transformations and load into destination tables using an ETL tool.

## Problem Definition:

Design a basic ETL model using Rapid Miner Application.

## Prerequisite:

- ☐ Basic concepts of ETL.
- ☐ Knowledge about Rapid miner tool.

## Software Requirements:

- ☐ Rapid Miner

## Learning Objectives:

Understand the implementation of the various ETL model using Rapid Miner tool.

## Outcomes:

After completion of this assignment students can develop and analyze the ETL model and will understand the working.

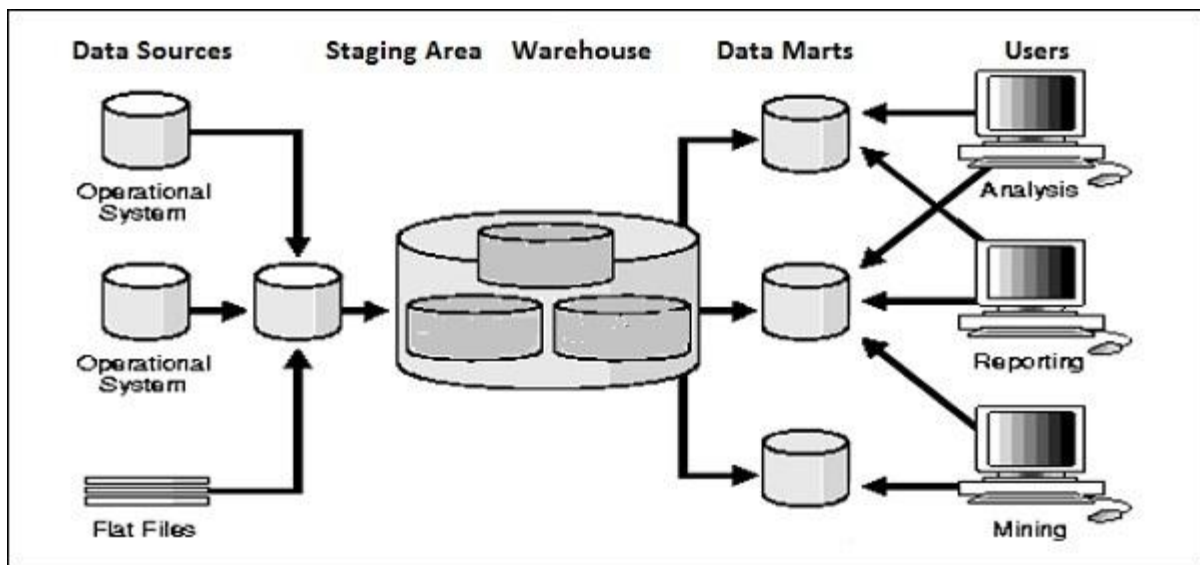
## Theory Concepts:

### What does ETL mean?

ETL stands for Extract, Transform and Load. An ETL tool extracts the data from different RDBMS source systems, transforms the data like applying calculations, concatenate, etc. and then load the data to Data Warehouse system. The data is loaded in the DW system in the form of dimension and fact tables.

### Extraction

- A staging area is required during ETL load. There are various reasons why staging area is required.
- The source systems are only available for specific period of time to extract data. This period of time is less than the total data-load time. Therefore, staging area allows you to extract the data from the source system and keeps it in the staging area before the time slot ends.
- Staging area is required when you want to get the data from multiple data sources together or if you want to join two or more systems together. For example, you will not be able to perform a SQL query joining two tables from two physically different databases.
- Data extractions' time slot for different systems vary as per the time zone and operational hours.
- Data extracted from source systems can be used in multiple data warehouse system, Operation Data stores, etc.
- ETL allows you to perform complex transformations and requires extra area to store the data.



## Transform

In data transformation, you apply a set of functions on extracted data to load it into the target system. Data, which does not require any transformation is known as direct move or pass through data.

You can apply different transformations on extracted data from the source system. For example, you can perform customized calculations. If you want sum-of-sales revenue and this is not in database, you can apply the **SUM** formula during transformation and load the data.

For example, if you have the first name and the last name in a table in different columns, you can use concatenate before loading.

## Load

During Load phase, data is loaded into the end-target system and it can be a flat file or a Data Warehouse system.

## Data Warehousing Schemas

1. Star Schema
2. Snowflake Schema
3. Fact Constellation

### Star Schema

For example, as you can see in the above-given image that fact table is at the center which contains keys to every dimension table like Deal\_ID, Model ID, Date\_ID, Product\_ID, Branch\_ID & other attributes like Units sold and revenue.

### Characteristics of Star Schema:

- Every dimension in a star schema is represented with the only one-dimension table.
- The dimension table should contain the set of attributes.
- The dimension table is joined to the fact table using a foreign key
- The dimension table are not joined to each other
- Fact table would contain key and measure
- The Star schema is easy to understand and provides optimal disk usage.
- The dimension tables are not normalized. For instance, in the above figure, Country\_ID does not have Country lookup table as an OLTP design would have.
- The schema is widely supported by BI Tools

## Snowflake Schema

A Snowflake Schema is an extension of a Star Schema, and it adds additional dimensions. It is called snowflake because its diagram resembles a Snowflake.

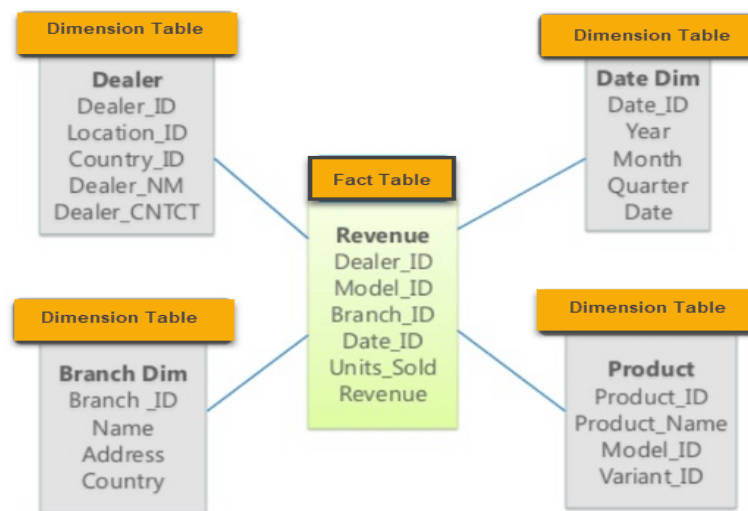
The dimension tables are normalized which splits data into additional tables. In the following example, Country is further normalized into an individual table.

### Characteristics of Snowflake Schema:

- The main benefit of the snowflake schema it uses smaller disk space.
- Easier to implement a dimension is added to the Schema
- Due to multiple tables query performance is reduced
- The primary challenge that you will face while using the snowflake Schema is that you need to perform more maintenance efforts because of the more lookup tables.

Star Schema	Snow Flake Schema
Hierarchies for the dimensions are stored in the dimensional table.	Hierarchies are divided into separate tables.
It contains a fact table surrounded by dimension tables.	One fact table surrounded by dimension table which are in turn surrounded by dimension table
In a star schema, only single join creates the relationship between the fact table and any dimension tables.	A snowflake schema requires many joins to fetch the data.
Simple DB Design.	Very Complex DB Design.
Denormalized Data structure and query also run faster.	Normalized Data Structure.

High level of Data redundancy	Very low-level data redundancy
Single Dimension table contains aggregated data.	Data Split into different Dimension Tables.
Cube processing is faster.	Cube processing might be slow because of the complex join
Offers higher performing queries using Star Join Query Optimization. Tables may be connected with multiple dimensions.	The Snow Flake Schema is represented by centralized fact table which unlikely connected with multiple dimensions.



**Star Schema**

### Tool for ETL: **RAPID MINER**

Rapid Miner is a world-leading open-source system for data mining. It is available as a stand-alone application for data analysis and as a data mining engine for the integration into own products. **Rapid Miner is now Rapid Miner Studio** and Rapid Analytics is now called Rapid Miner Server.

In a few words, Rapid Miner Studio is a "downloadable GUI for machine learning, data mining, text mining, predictive analytics and business analytics". It can also be used (for most purposes) in batch mode (command line mode)

Rapid Miner Support to Nominal, Numerical values, Integers, Real numbers, 2-value nominal, multi-value nominal etc.

## Dataset Description (Tables CSV and EXCEL used for ETL) :

### Table EXCEL:-

s_id	s_marks
1	350
2	400
3	450
4	500

### Table CSV:-

s_id	s_name	s_marks
5	Amruta	555
6	Vaishnavi	600
7	shailesh	650
8	Anand	750

## Steps and Operators Description with Screenshots

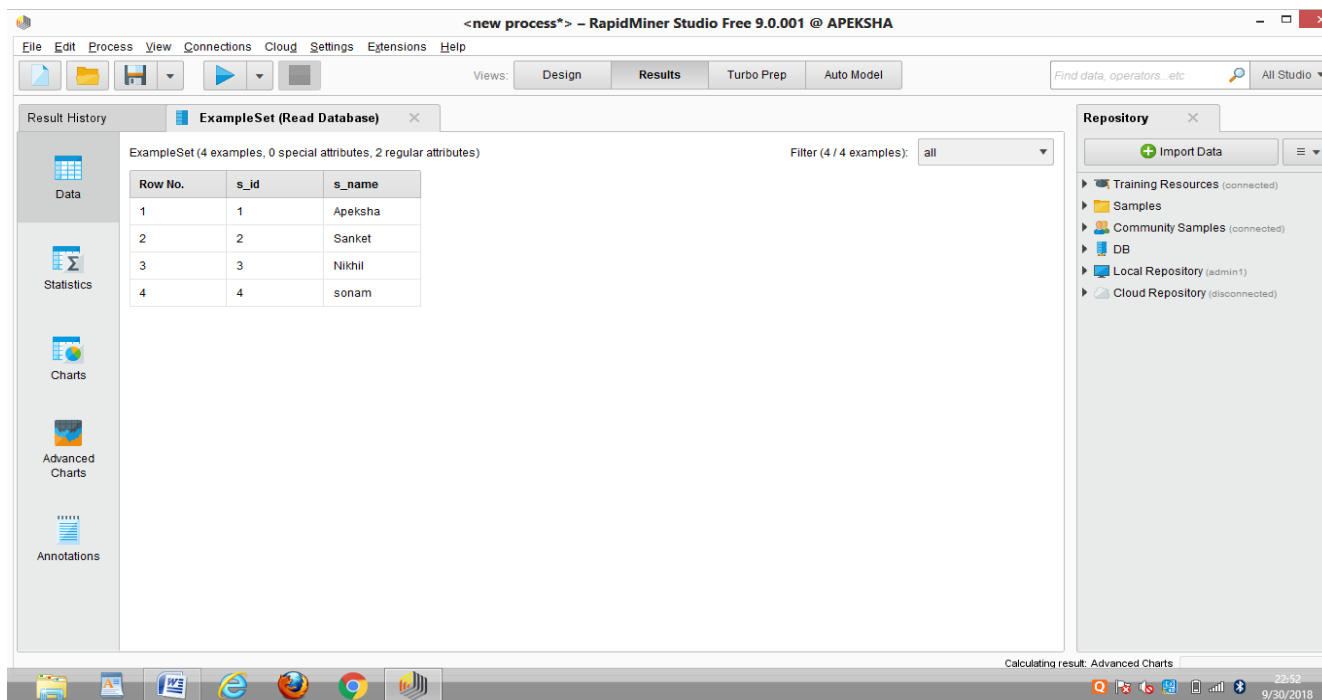
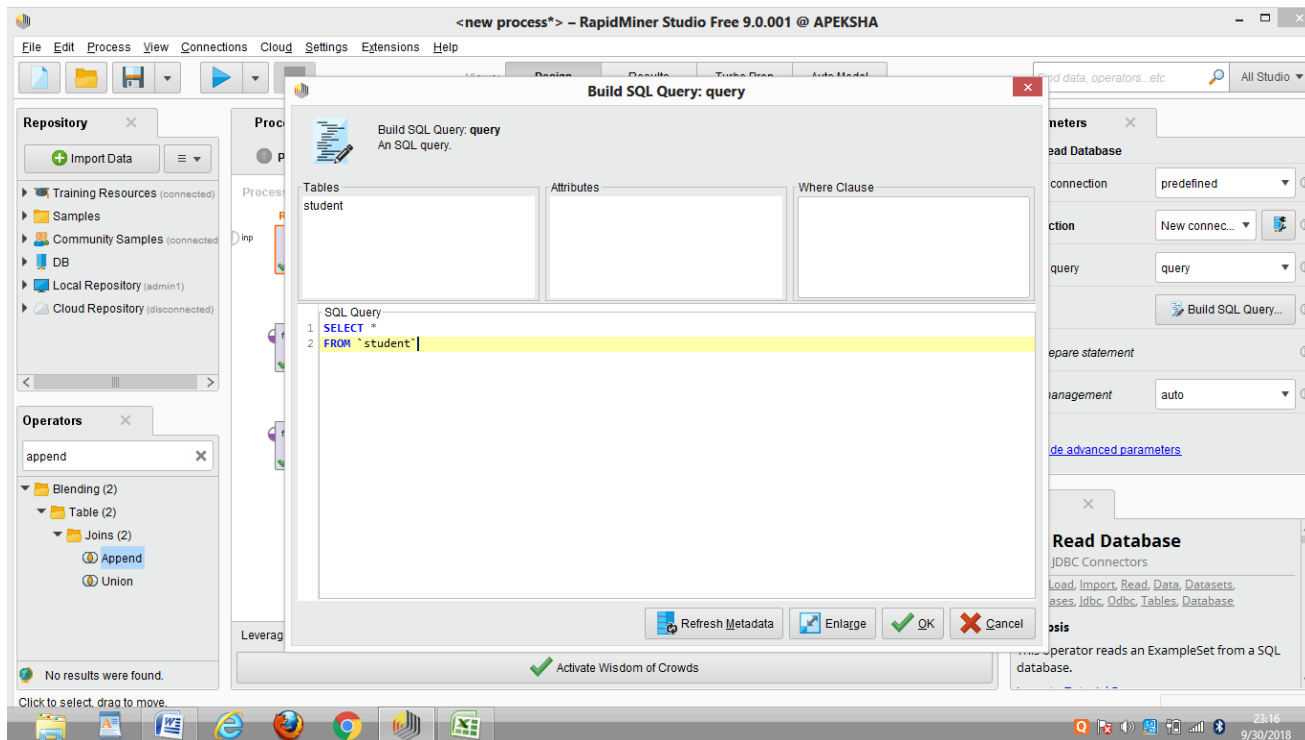
### 1. Read Database:-

The screenshot displays the RapidMiner Studio Free 9.0.001 interface. The main workspace shows a process diagram with a 'Read Database' operator. The 'Parameters' panel on the right is configured as follows:

- define connection: predefined
- connection: New connect...
- define query: query
- query: Build SQL Query...
- prepare statement: ☐
- data management: auto

The 'Help' panel on the right provides information about the 'Read Database' operator, including its tags (Load, Import, Read, Data, Datasets, Databases, jdbc, Odbc, Tables, Database) and a synopsis: 'This operator reads an ExampleSet from a SQL database.'

The 'Repository' panel on the left shows the 'Data Access (1)' folder expanded, with the 'Read Database' operator listed. The 'Operators' panel on the left shows a search for 'read datab' with no results found.



## 2. Read Excel:-

The screenshot shows the RapidMiner Studio interface. The 'Process' tab is active, displaying a workflow. The 'Read Database' operator is connected to a 'Set Role' operator. The 'Read Excel' operator is also connected to a 'Set Role' operator. The 'Set Role' operator has parameters: 'attribute name' set to 's\_marks', 'target role' set to 'regular', and 'set additional roles' set to 'Edit List (0)...'. The 'Parameters' panel on the right shows these settings. The 'Help' panel on the right provides information about the 'Set Role' operator, including its tags and synopsis.

Process Design:

```
graph LR; Inp(( )) --> ReadDB[Read Database]; ReadDB --> SetRole1[Set Role]; ReadExcel[Read Excel] --> SetRole2[Set Role (2)]; SetRole1 --> Out(( )); SetRole2 --> Out;
```

Parameters for Set Role (2):

- attribute name: s\_marks
- target role: regular
- set additional roles: Edit List (0)...

Help for Set Role:

Tags: Label, Target, Id, Class, Dependent, Independent, Special, Regular, Inputs, Columns, Attributes, Features, Variables, Types, Names & Roles

Synopsis: This Operator is used to change the role of one

The screenshot shows the 'Results' tab in RapidMiner Studio. The 'ExampleSet (Set Role (2))' is selected, showing a table with 4 examples, 1 special attribute, and 1 regular attribute. The table has columns 'Row No.', 's\_id', and 's\_marks'. The 'Repository' panel on the right shows the 'Import Data' button and a list of connected repositories.

Results:

Row No.	s_id	s_marks
1	1	350
2	2	400
3	3	450
4	4	500

### 3. Join Operation:-

The screenshot shows the RapidMiner Studio interface in Design view. The process flow is as follows:

- Read Database** (Input) connects to **Set Role**.
- Read Excel** (Input) connects to **Set Role (2)**.
- Set Role** and **Set Role (2)** both connect to the **Join** operator.

The **Parameters** panel on the right shows the following settings for the **Process** operator:

- logverbosity: init
- logfile: (empty)
- resultfile: (empty)
- random seed: 2001
- send mail: never
- encoding: SYSTEM

The **Help** panel shows the **Process** operator description: "The root operator which is the outer most operator of every process."

The screenshot shows the RapidMiner Studio interface in Results view. The **ExampleSet (Join)** is displayed as a table with 4 rows and 4 columns:

Row No.	s_id	s_name	s_marks
1	1	Apeksha	350
2	2	Sanket	400
3	3	Nikhil	450
4	4	sonam	500

The **Repository** panel on the right shows the same structure as in the Design view.



## 4. Read CSV:-

**Process Design:**

- Read Database** (inp) → **Set Role** (exa) → **Join** (lef, rig, joi)
- Read Excel** (fil) → **Set Role (2)** (exa) → **Join** (lef, rig, joi)
- Read CSV** (fil) → **Join** (lef, rig, joi)

**Parameters for Read CSV:**

- csv file: C:\Users\admin1\...
- column separators: ,
- trim lines: ☐
- use quotes: ☒
- quotes character: "

**Help for Read CSV:**

RapidMiner Studio Core

Tags: Load, Import, Read, Data, Files, Text, Comma, Spreadsheet, Excel, Datasets, Tsv

**Synopsis:** This Operator reads an ExampleSet from the specified CSV file.

**Result History:**

- ExampleSet (Read CSV)**
- ExampleSet (Join)**

**ExampleSet (4 examples, 1 special attribute, 2 regular attributes)**

Row No.	s_id	s_name	s_marks
1	5	Amruta	555
2	6	Vaishnavi	600
3	7	shailesh	650
4	8	Anand	750

Filter (4 / 4 examples): all

## 5. Append Operation:-

The screenshot shows the RapidMiner Studio interface with a process design in the 'Design' view. The process flow is as follows:

- Read Database** (input) connects to **Set Role**.
- Read Excel** (input) connects to **Set Role (2)**.
- Read CSV** (input) connects to the **Join** operator.
- Set Role** and **Set Role (2)** both connect to the **Join** operator.
- The **Join** operator connects to the **Append** operator.
- The **Append** operator has a warning icon, indicating a potential issue with the operation.

The **Parameters** panel on the right shows settings for the **Process** operator:

- logverbosity: init
- logfile: (empty)
- resultfile: (empty)
- random seed: 2001
- send mail: never
- encoding: SYSTEM

The **Help** panel shows the **Process** operator description:

**Synopsis**  
The root operator which is the outer most operator of every process.

**Description**

The screenshot shows the 'Results' view of the RapidMiner Studio interface. The **ExampleSet (Append)** is displayed, showing 8 examples with 3 attributes: s\_id, s\_name, and s\_marks.

Row No.	s_id	s_name	s_marks
1	5	Amruta	555
2	6	Vaishnavi	600
3	7	shailesh	650
4	8	Anand	750
5	1	Apeksha	350
6	2	Sanket	400
7	3	Nikhil	450
8	4	sonam	500

## **Conclusion**

With the help such Tools we can Perform ETL operations on Sample Data sets and can perform analysis on sample data sets.