

**School of Computing**  
**ST1501 Data Engineering CA2**  
**AY2022/2023 Semester 1**

---

**A. Instructions and Guidelines**

1. You are to work in a group of 3 (recommended) or 4 members. The database to use is **MS SQL Server**. Penalty will be accorded if other databases are used.
2. This is a **group** assignment (55%) **with individual** component (45%), which requires students to design and set up a data warehouse. The assignment will require students to demonstrate competency in designing data warehouse and integrate data from various sources.
3. The deadline of this assignment is on **5/8/2022 5 pm**.
4. Submissions should be made via the ST1501 CA2 Assignment Submission link by the stated deadline.
5. Deliverable include:
  - An zip file for the group task with the file-naming convention:  
“**ST1501-YourClass-GroupName.zip**”  
Example **ST1501-DITFT2B01-G1.zip**
  - An zip file for the individual task with the file-naming convention:  
“**ST1501-YourClass-YourStudentID-YourName.zip**”  
Example **ST1501-DITFT2B01-9999999-Benson Ang.zip**
6. The group task should consist of the following.
  - A submission template for group consisting:
    - (a) A Data Warehouse Schema
    - (b) An explanation for each query
  - All the SQL scripts to setup the OLTP database
  - All the SQL scripts to setup the data warehouse
  - All the SQL scripts for queries that can be supported by the data warehouse and their corresponding query results

Include all the above in the zip file “**ST1501-YourClass-GroupName.zip**”.
7. The individual task should consist of the following.
  - A submission template for individual consisting:
    - (a) List the SQL statements to create the documents for each collection
    - (b) All the commands to execute the queries from the MongoDB and the results.

Include all the above in the zip file "**ST1501-YourClass-YourStudentID-YourName.zip**".

8. As part of the assignment requirements, the team will require to present the assignment. The presentation should explain the data warehouse design and demonstrate the process of setup the OLTP database and the data warehouse. Each team member is required to demonstrate his ability to explain the DW and the questions fielded by the tutor during the presentation.
9. Student who is absent from the presentation will be **given zero mark** for the assignment.
10. This assignment will account for **40%** of the module grade.
11. No marks will be awarded, if the work is copied or you have allowed others to copy your work.
12. 50% of the marks will be deducted for assignments that are received within ONE (1) calendar day after the submission deadline. No marks will be given thereafter. Exceptions to this policy will be given to students with valid LOA on medical or compassionate grounds. Students in such cases will need to inform the module tutor as soon as reasonably possible. Students are not to assume on their own that their deadline has been extended.
13. Warning: Plagiarism means passing off as one's own the ideas, works, writings, etc., which belong to another person. In accordance with this definition, you are committing plagiarism if you copy the work of another person and turning it in as your own, even if you would have the permission of that person.  
Plagiarism is a serious offence, and if you are found to have committed, aided, and/or abetted the offence of plagiarism, disciplinary action will be taken against you. If you are guilty of plagiarism, you may fail all modules in the semester, or even be liable for expulsion.

## B. The Business Scenario



Trek Bicycle Store Inc. is a retailer of bicycles specializing in various types of bicycles. Currently, they are using an OLTP system in their day-to-day retail business operations. Trek Bicycle Store has a number of store outlets to serve the large community across various states in America. Each store has its unique name, however the outlets are owned by Trek Bicycle Store Inc.

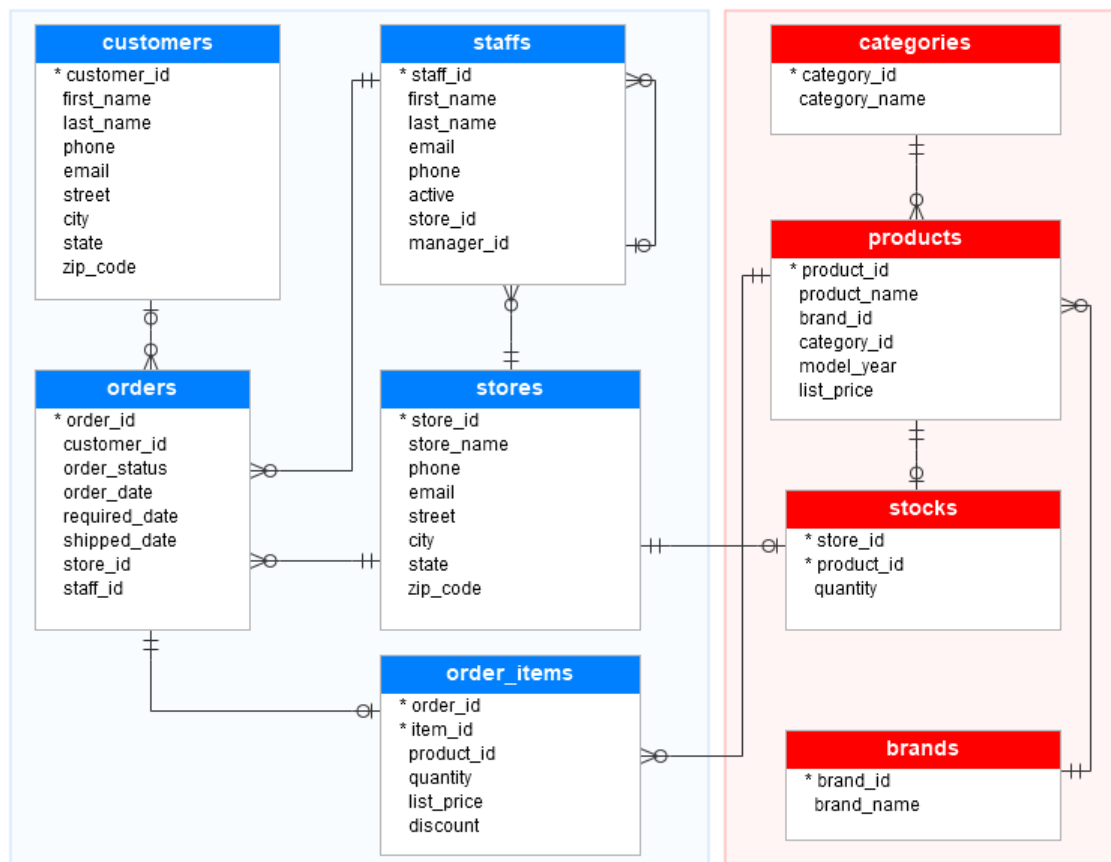
Recently, the owner of the company Fabiola Jackson wanted to incorporate a Data Warehouse solution in this existing IT infrastructure after listening to a talk on business analytics.

The owner wants to have business queries on production as well as sales to be answered from Data Warehouse. He also wishes to keep as much details of time attributes so that he knows the peak and low seasons. Inventory control is also top on his mind so that he knows when to replenish bicycles from suppliers should the stock level runs low or out of stock. Fabiola currently has 3 managers reporting to him across the 3 stores.

**OLTP Database**

Sales

Production



As you can see from the diagram, the database has two schemas sales and production, and these schemas have nine tables. A schema is term to describe a container to hold the various tables serving different functions the retailing business.

**Database Tables**

<b>Table stores</b>	This table includes the store's information. Each store has a store name, contact information such as phone and email, and an address including street, city, state, and zip code.
<b>Table staffs</b>	<p>This table stores the essential information of staffs including first name, last name. It also contains the communication information such as email and phone. A staff works at a store specified by the value in the store_id column.</p> <p>A store can have one or more staffs. A staff reports to a store manager specified by the value in the manager_id column. If the value in the manager_id is null, then the staff is the top manager. If a staff no longer works for any stores, the value in the active column is set to zero.</p>
<b>Table categories</b>	This table stores the bike's categories such as children bicycles, comfort bicycles, and electric bikes.

<b>Table brands</b>	This table stores the brand's information of bikes, for example, Electra, Haro, and Heller
<b>Table products</b>	<p>This table stores the product's information such as name, brand, category, model year, and list price.</p> <p>Each product belongs to a brand specified by the brand_id column. Hence, a brand may have zero or many products.</p> <p>Each product also belongs to a category specified by the category_id column. Also, each category may have zero or many products.</p>
<b>Table customers</b>	This table stores customer's information including first name, last name, phone, email, street, city, state and zip code.
<b>Table orders</b>	<p>This table stores the sales order's header information including customer, order status, order date, required date, shipped date.</p> <p>Order status: 1 = Pending; 2 = Processing; 3 = Rejected; 4 = Completed</p> <p>It also stores the information on where the sales transaction created (store) and who created it (staff). Each sales order has a row in the sales_orders table. A sales order has one or many line items stored in the sales_order_items table. The owner would wish to track the status of the orders in order to assess their efficiency in delivery with respect to the shipping date, required date and order date.</p> <p>As the status is fast changing attribute and transient, we do not track its history and let the most recent value overrides the previous one in DW.</p>
<b>Table order_items</b>	<p>This table stores the line items of a sales order. Each line item belongs to a sales order specified by the order_id column.</p> <p>A sales order line item includes product, order quantity, list price and discount.</p>
<b>Table stocks</b>	This table stores the inventory information i.e. the quantity of a particular product in a specific store. The quantity in stock is read on the date of the latest sales order for that particular product regardless of stores.

## C. Your Task

As a data specialist team, your team is asked to design and setup 2 databases: the OLTP database, the data warehouse using MS SQL Server.

### Group Task

- (a) Design a star or snowflake schema for a data warehouse that will help answer various business questions regarding sales based on the scenario described in Section B. The team **must seek your module tutor's feedback** on the design before proceeding to setup the data warehouse.
- (b) Setup the OLTP database (create the tables and insert the data) in the **Microsoft SQL Server** based on the OLTP database design given in Section B.
  - The database should be named as **BikeSalesXXXX** where XXXX is to be replaced by the team name.
  - All SQL scripts for loading data to the OLTP database are to be submitted.
  - The data to be loaded to the OLTP are given as follows:
    - The data from the tables in the OLTP database have come from various sources. Some data are in CSV format, TEXT format and JSON format.
    - The **Customers data**, **Orders data**, **Order items data** and **Stocks** have been prepared in **CSV format**.
    - The **Products data** are in **JSON format**.
    - The **stores, staff, brands and categories data** are given in **textual format**.
- (c) Create the data warehouse after your tutor has given feedback to your data warehouse design.

Create the database in the **Microsoft SQL Server** based on your data warehouse design.

- The database should be named as **BikeSalesDWXXXX** where XXXX is to be replaced by the team name.
  - Implement surrogate keys for all dimensions including the Time dimension and fact tables.
  - The SQL script to 'CREATE Table' is to be submitted
- (d) Load data from the source system (the OLTP database).
    - All the SQL scripts for loading the data to the warehouse are to be submitted.
  - (e) Besides the sales, the owner wants to have knowledge of the total stock quantity of bikes in the DW. As it is a dynamic value and changes from day to day depending on their orders. As an example, for product\_id CHB21, Electra Cruiser 1 (24-Inch) – 2016, we have the following distribution across the 3 stores:

store_id	product_id	quantity
ST1	CHB21	24
ST2	CHB21	16
ST3	CHB21	8
Total quantity		48

The combined stock quantity across all stores is 48 based on the most recent sale of this product on 2018-01-09. The owner wants to keep track of such stock levels. However, there are products that have never been sold in any orders. For such product stocks, the date for stock take is assumed to be the current date of loading data in the DW.

Some products are not even found in the stocks and this have to reflected as 0 for total quantity and date for stock take is assumed to be the current date of loading data in the DW.

The owner does not wish to track the history of inventory changes but he needs to know the current stock of his products at point in time. The DW must be able to reflect the stock inventory.

**You can ignore Slow and Fast Changing Dimensions in your DW.**

- (f) Provide meaningful queries that can be supported by your data warehouse. The SQL script, the corresponding business question and its query result are to be submitted.

- Query 1 - Sales/profits/discounts/revenue
- Query 2 - Sales/Staff/stores
- Query 3 - Sales/Seasons of Sales/time
- Query 4 - Sales/Orders/Customers
- Query 5 - Sales/Products/brands/categories/inventory

The query should provide insightful findings to the owner. Recommended to incorporate aggregate functions, row wise functions, group by, sub-query etc.

Remember, there is no right and wrong answers, only insightful or trivial observations. That will depend on how much you want to explore and your examination of the data. Stating the obvious will not be eligible for high grade. Your tutor can tell whether enough thoughts have been put into the queries presented.

### **Individual Task**

In the following, you shall experience how to store and manipulate data in NoSQL databases such as MongoDB. You are to show all commands and results in using MongoDB command-line tool (mongo.exe) for each item listed below:

## (a) Mongo database creation

Taking reference from OLTP, create a Mongo database called Bikes with 3 collections:

*Hint: use SQL with JSON format*

- UnSold – create a collection to store those products which have never been sold
- ZeroStock – create a collection to store products which have zero stock
- Stock – create a collection to store products which have stock
  - For UnSold and ZeroStock collections, the documents must have the fields product\_id, product\_name, brand\_name, category\_name, model\_year and list\_price
  - For the Stock collection, the documents must have fields product\_id, store\_id and quantity
  - Find the count of documents in each collection

## (b) Queries

Use MongoDB commands to demonstrate the following:

- Query 1 - Implement a query to find out those bike names which are less than \$2000 and belonging to Road Bikes and have zero stock
- Query 2 - Implement a query to find out the unduplicated category names and unduplicated brand names separately for these unsold products

*Hint: Use MongoDB Aggregate pipeline & \$Lookup*

- Query 3 - Implement a query to find the unduplicated sets of category names and brand names in combination
- Query 4 - Implement a query to find the information of those products which are unsold and have stock
- Query 5 – Making use of the product\_id from Query 4, implement a query to find their total stock available arranged in descending order
- Query 6 – Modify Query 4 and implement a query to find the information of those products which are sold but do not have stock

**-- End of Assignment --**



## D. Submission Template

### Group Task

Data Warehouse Schema

--

Queries explanation

Query	Explanation	SQL script
1		
2		
3		
4		
5		

**Individual Task**

Item	SQL statement for JSON format
UnSold – create a collection to store those products which have never been sold	
ZeroStock – create a collection to store products which have zero stock	
Stock – create a collection to store products which have stock	

Item	MongoDB command	Results
Count of Unsold Collection		
Count of ZeroStock Collection		
Count of Stock Collection		
Query 1		
Query 2		
Query 3		
Query 4		
Query 5		
Query 6		

## The Assessment Criteria

Component	Tasks	Weightage	
Group Work 55%	<b>Data Warehouse Design (Database Diagram)</b> <ul style="list-style-type: none"><li>The design supports the described business scenario.</li><li>The chosen table names, field names and attributes are descriptive.</li><li>The explanation of the design is clear and concise.</li></ul>	5%	
	<b>OLTP Database Setup</b> <ul style="list-style-type: none"><li>The Create Table SQL script implements the database design, including the primary key and foreign key definition.</li><li>The Insert Records SQL script to load the given data.</li><li>The SQL Script to load the Product and Customer records</li><li>Explanation and the demonstration of OLTP setup is without error during the presentation.</li></ul>	5%	
	<b>The database warehouse setup in MS SQL</b> <ul style="list-style-type: none"><li>The Create Table SQL script implements the star or snowflake design, including the primary key, surrogate key and foreign key definition.</li><li>The SQL Script to load the data into the data warehouse</li></ul>	15%	
	<b>The queries in MS SQL</b> <ul style="list-style-type: none"><li>Five insightful query that covers each of the category (Sales, Staff/Office, Seasons of Sales, Orders/Customer, Products/brands/categories/inventory)</li></ul>	6%	30%
		6%	
6%			
6%			
6%			
Individual Work 45%	<b>Setup MongoDB</b> <ul style="list-style-type: none"><li>The commands to create the database</li><li>The commands to create the collections</li><li>The command to execute the viewing of the collection documents</li></ul>	10%	
	<b>The queries in the MongoDB</b> <ul style="list-style-type: none"><li>The queries as implemented</li><li>The commands to execute</li><li>The results from MongoDB</li></ul>	3%	25%
		3%	
		3%	
		6%	
		7%	
	3%		
<b>Question and answer during the presentation</b>		10%	
<b>Total</b>		<b>100%</b>	

## Appendix 1

### Table creation

```
CREATE TABLE sales.stores (  
  store_id varchar(5) PRIMARY KEY,  
  store_name VARCHAR (255) NOT NULL,  
  phone VARCHAR (25),  
  email VARCHAR (255),  
  street VARCHAR (255),  
  city VARCHAR (255),  
  state VARCHAR (10),  
  zip_code VARCHAR (5)  
);  
  
CREATE TABLE sales.staffs (  
  staff_id varchar(5) PRIMARY KEY,  
  first_name VARCHAR (50) NOT NULL,  
  last_name VARCHAR (50) NOT NULL,  
  email VARCHAR (255) NOT NULL UNIQUE,  
  phone VARCHAR (25),  
  active int NOT NULL,  
  store_id varchar(5) NOT NULL,  
  manager_id varchar(5),  
  FOREIGN KEY (store_id) REFERENCES sales.stores (store_id)  
  ON DELETE CASCADE  
  ON UPDATE CASCADE,  
  FOREIGN KEY (manager_id) REFERENCES sales.staffs (staff_id)  
  ON DELETE NO ACTION  
  ON UPDATE NO ACTION  
);  
  
CREATE TABLE production.categories (  
  category_id varchar(5) PRIMARY KEY,  
  category_name VARCHAR (255) NOT NULL  
);  
CREATE TABLE production.brands (  
  brand_id varchar(5) PRIMARY KEY,  
  brand_name VARCHAR (255) NOT NULL  
);  
  
CREATE TABLE production.products (  
  product_id varchar(10) PRIMARY KEY,  
  product_name VARCHAR (255) NOT NULL,  
  brand_id varchar(5) NOT NULL,  
  category_id varchar(5) NOT NULL,  
  model_year int NOT NULL,
```

```

list_price DECIMAL (10, 2) NOT NULL,
FOREIGN KEY (category_id) REFERENCES production.categories (category_id)
ON DELETE CASCADE ON UPDATE CASCADE,
FOREIGN KEY (brand_id) REFERENCES sales.brands (brand_id)
ON DELETE CASCADE ON UPDATE CASCADE
);
CREATE TABLE sales.customers (
customer_id varchar(10) PRIMARY KEY,
first_name VARCHAR (255) NOT NULL,
last_name VARCHAR (255) NOT NULL,
phone VARCHAR (25),
email VARCHAR (255) NOT NULL,
street VARCHAR (255),
city VARCHAR (50),
state VARCHAR (25),
zip_code VARCHAR (5)
);
CREATE TABLE sales.orders (
order_id varchar(10) PRIMARY KEY,
customer_id varchar(10),
order_status int NOT NULL,
-- Order status: 1 = Pending; 2 = Processing; 3 = Rejected; 4 = Completed
order_date DATE NOT NULL,
required_date DATE NOT NULL,
shipped_date DATE,
store_id INT NOT NULL,
staff_id INT NOT NULL,
FOREIGN KEY (customer_id) REFERENCES sales.customers (customer_id)
ON DELETE CASCADE ON UPDATE CASCADE,
FOREIGN KEY (store_id) REFERENCES sales.stores (store_id)
ON DELETE CASCADE ON UPDATE CASCADE,
FOREIGN KEY (staff_id) REFERENCES sales.staffs (staff_id)
ON DELETE NO ACTION
ON UPDATE NO ACTION
);
CREATE TABLE sales.order_items(
order_id varchar(10),
item_id INT,
product_id varchar(10) NOT NULL,
quantity INT NOT NULL,
list_price DECIMAL (10, 2) NOT NULL,
discount DECIMAL (4, 2) NOT NULL DEFAULT 0,
PRIMARY KEY (order_id, item_id),
FOREIGN KEY (order_id) REFERENCES sales.orders (order_id)
ON DELETE CASCADE ON UPDATE CASCADE,
FOREIGN KEY (product_id) REFERENCES production.products (product_id)

```

```
ON DELETE CASCADE
ON UPDATE CASCADE
);
CREATE TABLE production.stocks (
  store_id varchar(5),
  product_id varchar(10),
  quantity INT,
  PRIMARY KEY (store_id, product_id),
  FOREIGN KEY (store_id) REFERENCES sales.stores (store_id)
  ON DELETE CASCADE ON UPDATE CASCADE,
  FOREIGN KEY (product_id) REFERENCES production.products (product_id)
  ON DELETE CASCADE
  ON UPDATE CASCADE
);
```

## Appendix 2

### Data

#### Staff

staff_id	first_name	last_name	email	phone	active	store_id	manager_id
3031	Fabiola	Jackson	<a href="mailto:fabiola.jackson@bikes.shop">fabiola.jackson@bikes.shop</a>	(831) 555-5554	1	ST1	NULL
30310	Bernardine	Houston	<a href="mailto:bernardine.houston@bikes.shop">bernardine.houston@bikes.shop</a>	(972) 530-5557	1	ST3	3037
3032	Mireya	Copeland	<a href="mailto:mireya.copeland@bikes.shop">mireya.copeland@bikes.shop</a>	(831) 555-5555	1	ST1	3031
3033	Genna	Serrano	<a href="mailto:genna.serrano@bikes.shop">genna.serrano@bikes.shop</a>	(831) 555-5556	1	ST1	3032
3034	Virgie	Wiggins	<a href="mailto:virgie.wiggins@bikes.shop">virgie.wiggins@bikes.shop</a>	(831) 555-5557	1	ST1	3032
3035	Jannette	David	<a href="mailto:jannette.david@bikes.shop">jannette.david@bikes.shop</a>	(516) 379-4444	1	ST2	3031
3036	Marcelene	Boyer	<a href="mailto:marcelene.boyer@bikes.shop">marcelene.boyer@bikes.shop</a>	(516) 379-4445	1	ST2	3035
3037	Venita	Daniel	<a href="mailto:venita.daniel@bikes.shop">venita.daniel@bikes.shop</a>	(516) 379-4446	1	ST2	3035
3038	Kali	Vargas	<a href="mailto:kali.vargas@bikes.shop">kali.vargas@bikes.shop</a>	(972) 530-5555	1	ST3	3031
3039	Layla	Terrell	<a href="mailto:layla.terrell@bikes.shop">layla.terrell@bikes.shop</a>	(972) 530-5556	1	ST3	3037

#### Brand

brand_id	brand_name
BRD1	Electra
BRD2	Haro
BRD3	Heller
BRD4	Pure Cycles
BRD5	Ritchey
BRD6	Strider
BRD7	Sun Bicycles
BRD8	Surly
BRD9	Trek

**Category**

<b>category_id</b>	<b>category_name</b>
CHB1	Children Bicycles
CMB2	Comfort Bicycles
CRB3	Cruisers Bicycles
CYB4	Cyclocross Bicycles
ELB5	Electric Bikes
MOB6	Mountain Bikes
RDB7	Road Bikes

**Store**

<b>store_id</b>	<b>store_name</b>	<b>phone</b>	<b>email</b>	<b>street</b>	<b>city</b>	<b>state</b>	<b>zip_code</b>
ST1	Santa Cruz Bikes	(831) 476-4321	santacruz@bikes.shop	3700 Portola Drive	Santa Cruz	CA	95060
ST2	Baldwin Bikes	(516) 379-8888	baldwin@bikes.shop	4200 Chestnut Lane	Baldwin	NY	11432
ST3	Rowlett Bikes	(972) 530-5555	rowlett@bikes.shop	8000 Fairway Avenue	Rowlett	TX	75088