# Database & SQL

*Induction Training*

## *Do The Math*

**Chicago, IL**
**Bangalore, India**
**www.mu-sigma.com**

January 16, 2013

**So… what is SQL and why do we need it?**

# SQL stands for:
# Scarcely Qualified Language

Just kidding!

**So… what is SQL and why do we need it?**

# SQL stands for:
# Structured Query Language

That's all you will learn for now!

Kidding!!... But we will hold off SQL for sometime

**Let us spend two minutes on this…**

# What are your expectations from this training?

**Before we dive into SQL, let's talk about DATA…**



Talk about me?

Sorry… when we say DATA,
we mean **information**

**Let's try and understand where and how data is created**

A Retail Supermarket

A Customer





What kind of information is captured?

# A retail store captures a lot of information regarding a customer's shopping behavior

**Who is coming to my store?**

**What products are being bought?**

**What frequently are my products purchased?**

# Let's list out some metrics or data a retail store might collect

**1 Customer Information**
- Customer Name
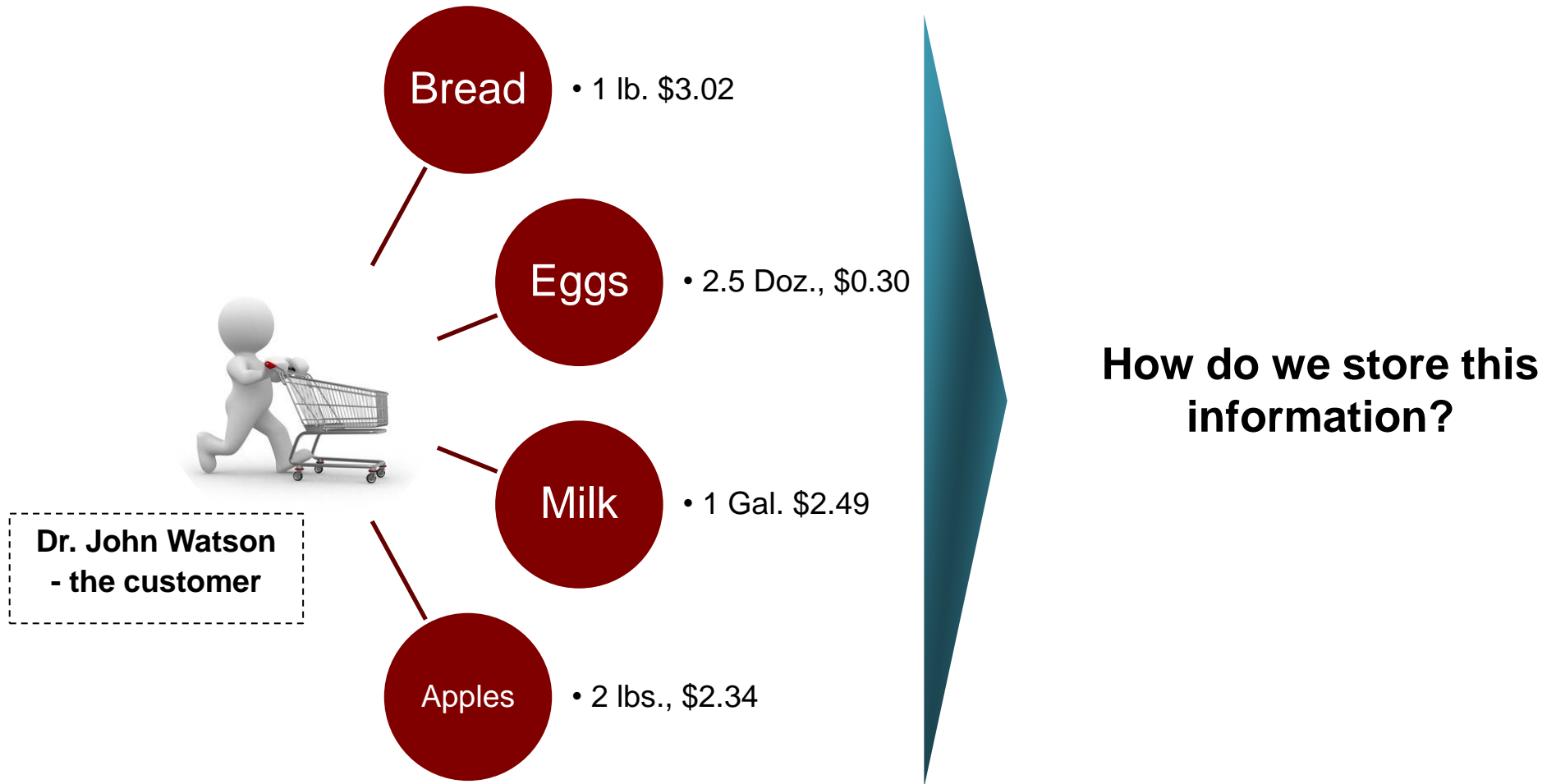- Customer Address
- Age
- Telephone

**2 Product Information**
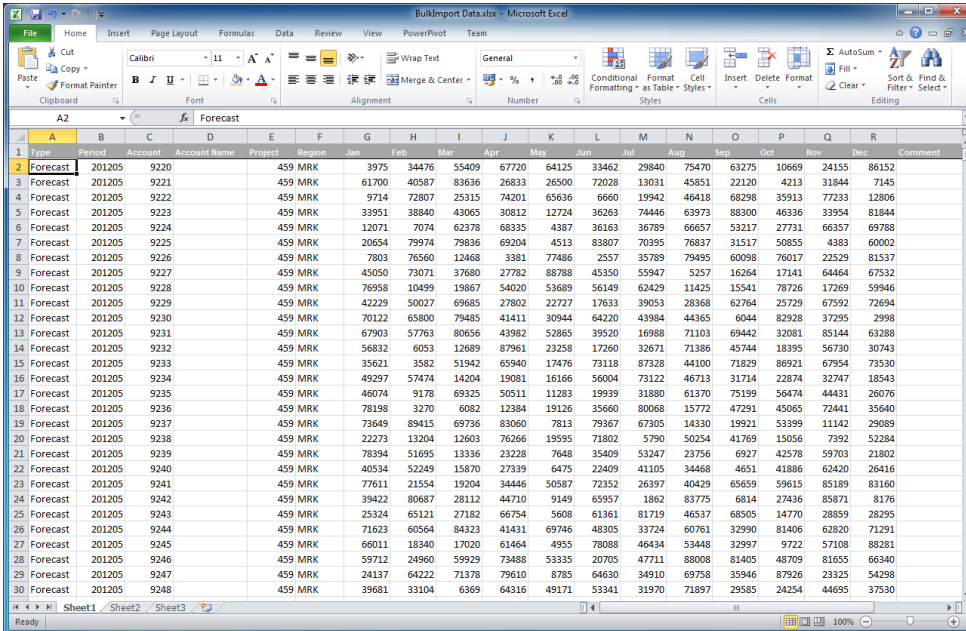- Product Name
- Product Class
- Price
- Discount

**2 Transaction Information**
- Date
- Customer Name
- Product Name
- Cost

# Now, consider purchases of a single customer

**Bread**
- 1 lb. $3.02

**Eggs**
- 2.5 Doz., $0.30

**Milk**
- 1 Gal. $2.49

**Apples**
- 2 lbs., $2.34

Dr. John Watson - the customer

**How do we store this information?**

# What if we use Excel to store this information?



**Excel Spreadsheet**

▶ What if 100 people want to modify the file?

▶ What if the number of transactions cross a million?

▶ How do you manage a million pieces of information?

▶ **An Excel document is clearly not the right choice for storing information in a transactional environment**

▶ **What we need is a file system which can**

– **Hold multiple tables of information**

– **Do not run out of space**

– **Allow multiple people to access and modify the information simultaneously**

# A database allows multiple users to interact with the same piece of information simultaneously

InterBase        Firebird        Panorama

SQLBase        **SQL Server**        Sybase

# RDBMS

Access

Filemaker

## Oracle   DB2        SQLite

FoxPro

PostGreSQL

MaxDB        # MySQL        Teradata

▸ More on **RDBMS** in the later slides …

# How should we store every transaction in the database?

| Date | First Name | Last Name | Bread | Eggs | Milk | Apple |
|------|------------|-----------|-------|------|------|-------|
| 3-Mar-12 | John | Watson | $3.02 | $0.30 | $2.49 | $2.34 |

| Advantages | Disadvantages |
|------------|---------------|
| ▶ I can find everything about a purchase by looking at single row | ▶ Not feasible as the number of columns is undefined (if a customer buys 20 items) |

# How should we store every transaction in the database?

**Store every __item__ in a transaction as a single row**

| Date | First Name | Last Name | Product | Quantity | Price |
|------|-----------|-----------|---------|----------|-------|
| 3-Mar-12 | John | Watson | Bread | 1 lb. | $3.02 |
| 3-Mar-12 | John | Watson | Eggs | 2.4 doz. | $0.30 |
| 3-Mar-12 | John | Watson | Milk | 1 gal. | $2.49 |
| 3-Mar-12 | John | Watson | Apple | 2 lbs. | $2.34 |

| Advantages | Disadvantages |
|------------|---------------|
| ▸ I don't have to worry about how many items a customer is buying | ▸ There is duplication of information across the rows |

# That brings us to the subject of a relational database system



E-R diagram

▸ A database is just a location to store and retrieve data

▸ A relational database is one which treats all of its data as a collection of relations

▸ Each table in a relational database has unique key (also referred to as the primary key)

▸ This primary key can also be present in another table as a foreign key and in process creates a relation between two tables

# To create a relational database, we need to understand the granularity or level of each table

| StudentID | FirstName | LastName | CourseID |
|-----------|-----------|----------|----------|
| S002001 | Jack | Black | C004 |
| S002002 | Lucy | Liu | C001 |
| S002003 | Dustin | Hoffman | C002 |
| S002004 | Seth | Rogen | C004 |

**Foreign Key**

**Relationship**

**Primary Key**

**A primary key uniquely identifies each record in a table**

| CourseID | CourseName |
|----------|------------|
| C001 | History |
| C002 | Computing |
| C003 | Biology |
| C004 | English |

# A simple relational database



**DimProduct**
- ProductKey
- ProductLabel
- ProductName
- ProductDescription
- ProductSubcategoryKey
- Manufacturer
- BrandName
- ClassID
- ClassName
- StyleID
- StyleName
- ColorID
- ColorName
- Size
- SizeRange
- SizeUnitMeasureID
- Weight
- WeightUnitMeasureID
- UnitOfMeasureID
- UnitOfMeasureName
- StockTypeID
- StockTypeName
- UnitCost
- UnitPrice
- AvailableForSaleDate
- StopSaleDate
- Status
- ImageURL
- ProductURL
- ETLLoadID
- LoadDate
- UpdateDate

**FactOnlineSales**
- OnlineSalesKey
- DateKey
- StoreKey
- ProductKey
- PromotionKey
- CurrencyKey
- CustomerKey
- SalesOrderNumber
- SalesOrderLineNumber
- SalesQuantity
- SalesAmount
- ReturnQuantity
- ReturnAmount
- DiscountQuantity
- DiscountAmount
- TotalCost
- UnitCost
- UnitPrice
- ETLLoadID
- LoadDate
- UpdateDate

**DimCustomer**
- CustomerKey
- GeographyKey
- CustomerLabel
- Title
- FirstName
- MiddleName
- LastName
- NameStyle
- BirthDate
- MaritalStatus
- Suffix
- Gender
- EmailAddress
- YearlyIncome
- TotalChildren
- NumberChildrenAtHome
- Education
- Occupation
- HouseOwnerFlag
- NumberCarsOwned
- AddressLine1
- AddressLine2
- Phone
- DateFirstPurchase
- CustomerType
- CompanyName
- ETLLoadID
- LoadDate
- UpdateDate

Product table has **ProductKey** as the primary key

Customer table has **CustomerKey** as the primary key

FactOnlineSales table has both **CustomerKey** and **ProductKey** as foreign keys

# A database management system is not just a place to dump data

Creating and changing the logical structure of a database

Querying and making changes to the information

Menus, data entry screens, reports and application software

Who can use and see what information; methods for backup and recovery

Queries to see the effect of structural changes

Data Definition → DBMS ENGINE

Data Manipulation → DBMS ENGINE

Application Generation → DBMS ENGINE

Database Administration → DBMS ENGINE

DBMS ENGINE ⇄ Relational Database

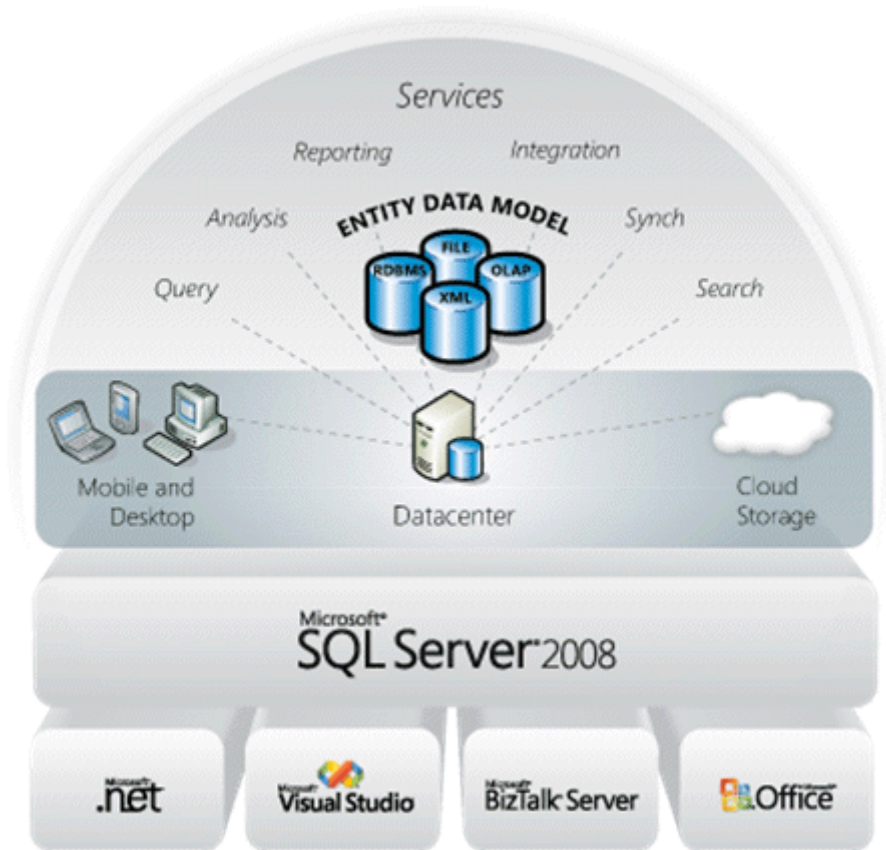- **Security**
- **Integrity**
- **Performance**
- **Reliability**
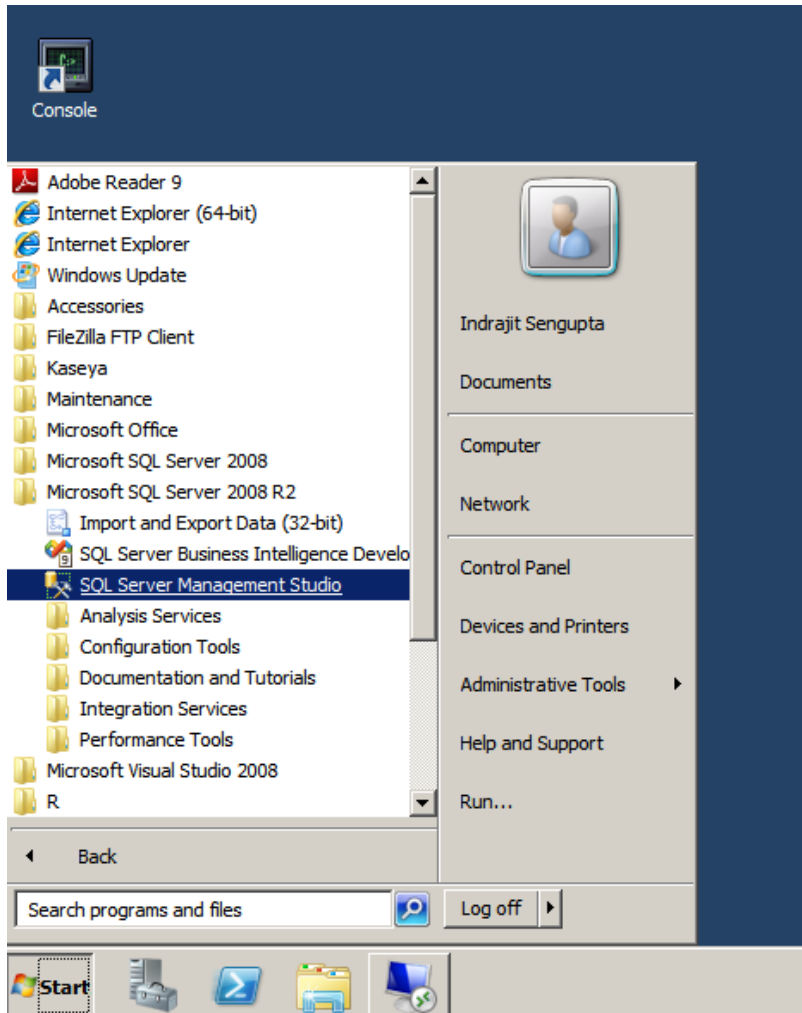
# Coming back to what is SQL and why do we need it?

▸ SQL is the standard language for communicating with relational database management systems

▸ SQL statements are used to perform tasks such as update data on a database, or retrieve data from a database

▸ Although most database systems use SQL, most of them also have their own additional proprietary extensions that are usually only used on their system

▸ However the standard SQL commands can be used to accomplish almost everything one needs

# For our learning purpose we will be using the Microsoft SQL Server 2008 Express edition



▶ SQL Server 2008 R2 Express edition is a free version of SQL Server 2008 R2

▶ Almost everything that can be done in the commercial version is supported in the express edition

▶ Database size is restricted to 10 GB in size

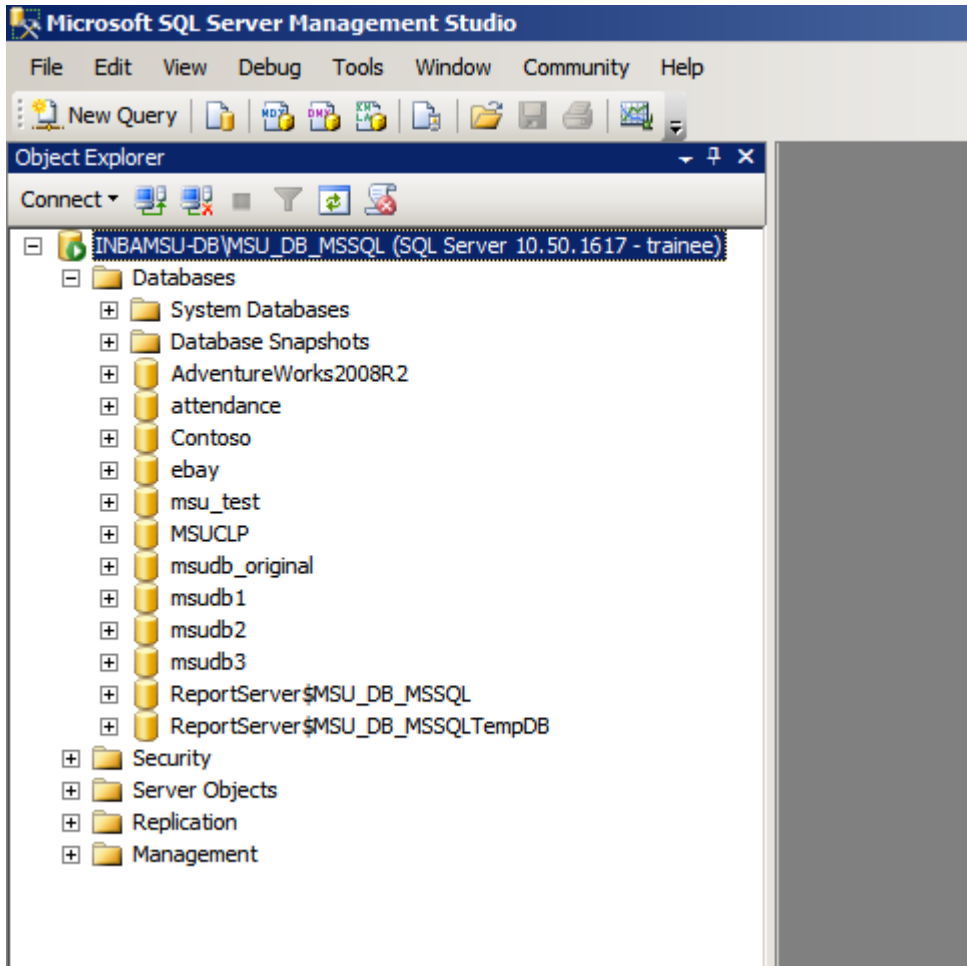▶ Uses only a single processor

▶ Uses only 1 GB of internal memory

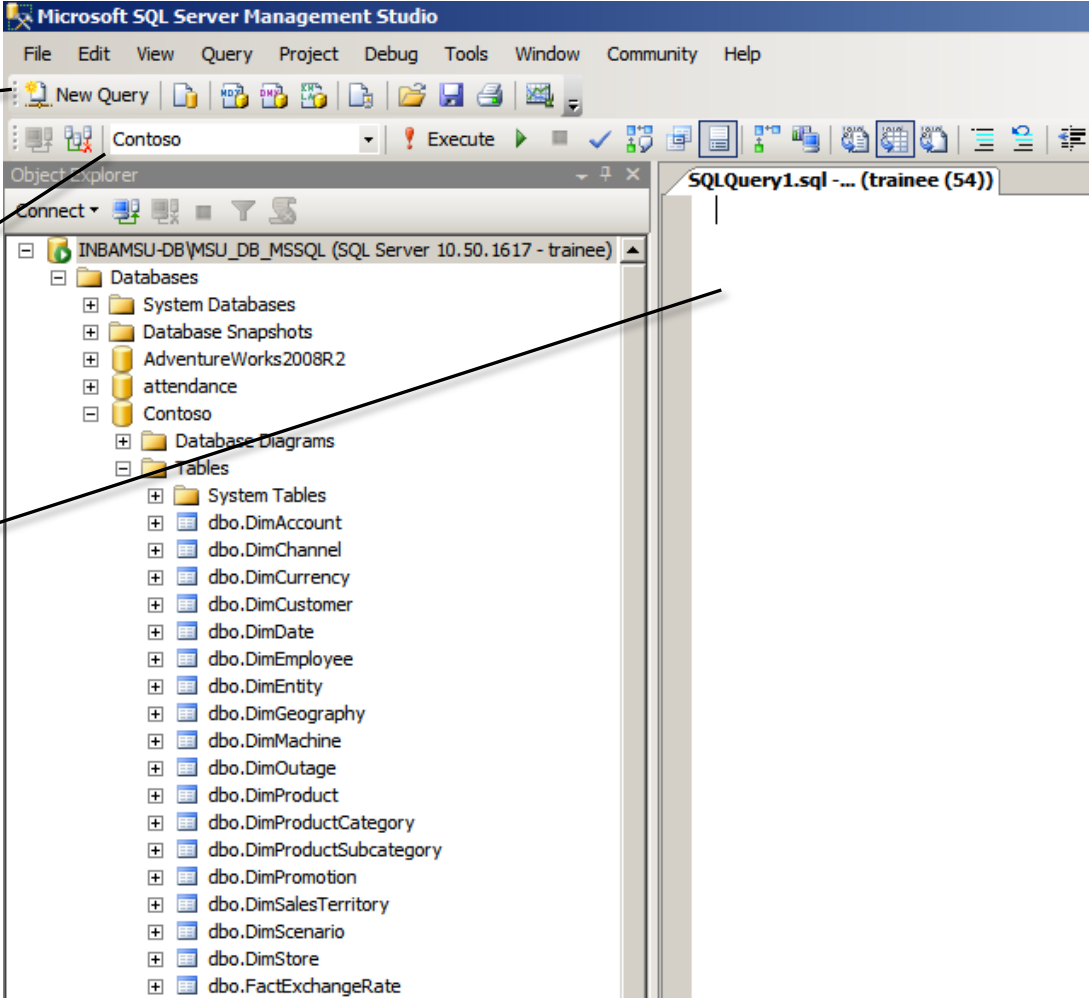# Start up SQL Server Management Studio from your Start menu



▸ Enter the server name: **INBAMSU-DB\MSU_DB_MSSQL**

▸ Enter user id: **trainee**

▸ Enter password: **password@1234**

▸ Remember to choose **SQL Server Authentication**

# The Object Explorer window will list out all the databases



▸ Once you login to the server, the object explorer window on the left will list out all the objects

▸ Expand the **Database** object to display all the databases currently installed

▸ We will focus on the **Contoso** database for all our practice sessions

# Open a query window by clicking on the New Query button



Open query window by clicking here

Select **Contoso** database from the drop down

Enter all your queries here

# Your first SQL query



Enter the query and press **F5** or **Execute** to submit the query

Query **output**

# The SELECT query extracts rows from a table

```
SELECT *

FROM dbo.DimProductCategory;
```

| | Results | Messages | | |
|---|---|---|---|---|
| | ProductCategoryKey | ProductCategoryLabel | ProductCategoryName | ProductCategoryDescription |
| 1 | 1 | 01 | Audio | Audio |
| 2 | 2 | 02 | TV and Video | TV and Video |
| 3 | 3 | 03 | Computers | Computers |
| 4 | 4 | 04 | Cameras and camcorders | Cameras and camcorders |
| 5 | 5 | 05 | Cell phones | Cell phones |
| 6 | 6 | 06 | Music, Movies and Audio Books | Music, Movies and Audio Books |
| 7 | 7 | 07 | Games and Toys | Games and Toys |
| 8 | 8 | 08 | Home Appliances | Home Appliances |

▸ The **SELECT** query is used to select rows from a table in a database

▸ The asterisk (*) indicates all rows and all columns

▸ The **FROM** statement indicates the table you want to access

# You can specify the column names to focus on only a few fields

```sql
SELECT
ProductCategoryLabel,
ProductCategoryName
FROM dbo.DimProductCategory;
```

| | ProductCategoryLabel | ProductCategoryName |
|---|---|---|
| 1 | 01 | Audio |
| 2 | 02 | TV and Video |
| 3 | 03 | Computers |
| 4 | 04 | Cameras and camcorders |
| 5 | 05 | Cell phones |
| 6 | 06 | Music, Movies and Audio Books |
| 7 | 07 | Games and Toys |
| 8 | 08 | Home Appliances |

▸ You can specify the columns you want to display by listing out the names separated by a comma

▸ SQL is not case sensitive, hence lower case syntax is also correct

# Use the WHERE statement to apply conditions on rows

```
SELECT
        FirstName,
        LastName,
        YearlyIncome,
        Gender,
        Occupation
FROM dbo.DimCustomer
WHERE YearlyIncome > 100000;
```

| | First Name | Last Name | YearlyIncome | Gender | Occupation |
|---|---|---|---|---|---|
| 1 | Donald | Gonzalez | 160000.00 | M | Management |
| 2 | Damien | Chander | 170000.00 | M | Management |
| 3 | Savannah | Baker | 120000.00 | F | Management |
| 4 | Angela | Butler | 130000.00 | F | Management |
| 5 | Alyssa | Cox | 130000.00 | F | Management |
| 6 | Sarah | Thomas | 110000.00 | F | Management |
| 7 | Nicholas | Robinson | 110000.00 | M | Management |
| 8 | Jose | Flores | 110000.00 | M | Management |
| 9 | Molly | Rodriguez | 120000.00 | F | Management |
| 10 | April | Anand | 160000.00 | F | Management |
| 11 | Devin | Martin | 170000.00 | M | Management |

▸ You can specify multiple conditions by separating each condition by the keywords **AND** or **OR**

▸ For better readability enclose the conditions with brackets

# The ORDER statement sorts the output in ascending or descending order of a variable

```
SELECT
        FirstName,
        YearlyIncome,
        Occupation
FROM dbo.DimCustomer
WHERE YearlyIncome > 100000
ORDER BY YearlyIncome;
```

| | FirstName | YearlyIncome | Occupation |
|----|-----------|--------------|------------|
| 1 | Sarah | 110000.00 | Management |
| 2 | Nicholas | 110000.00 | Management |
| 3 | Jose | 110000.00 | Management |
| 4 | Larry | 110000.00 | Professional |
| 5 | Tristan | 110000.00 | Professional |
| 6 | Megan | 110000.00 | Management |
| 7 | Arturo | 110000.00 | Management |
| 8 | Theresa | 110000.00 | Management |
| 9 | Cindy | 110000.00 | Management |
| 10 | Sean | 110000.00 | Professional |

▸ The default sort order is ascending – to sort in descending order, use the keyword **DESC** after the variable name

# The DISTINCT keyword identifies unique rows from a table

```
SELECT DISTINCT

        Occupation

FROM dbo.DimCustomer;
```

| Results | Messages |
| --- | --- |
| | Occupation |
| 1 | Professional |
| 2 | NULL |
| 3 | Manual |
| 4 | Clerical |
| 5 | Management |
| 6 | Skilled Manual |

▸ Specifying multiple column names with **DISTINCT** keyword results in selecting unique combinations of all the columns from the table

# The GROUP BY statement is used with summary functions to roll up tables

```sql
SELECT
        Occupation,
        Gender,
        COUNT(CustomerKey) AS
TotalCustomers
FROM dbo.DimCustomer
GROUP BY
        Occupation,
        Gender;
```

| | Occupation | Gender | TotalCustomers |
|---|---|---|---|
| 1 | Management | M | 1592 |
| 2 | Skilled Manual | M | 2293 |
| 3 | Professional | M | 2727 |
| 4 | Manual | F | 1133 |
| 5 | Management | F | 1483 |
| 6 | Clerical | M | 1488 |
| 7 | Clerical | F | 1440 |
| 8 | Manual | M | 1251 |
| 9 | Professional | F | 2793 |
| 10 | Skilled Manual | F | 2284 |
| 11 | NULL | NULL | 385 |

Results | Messages

▸ Various summary functions are available in SQL language like **SUM**, **AVG**, **MIN**, **MAX**, **STDEV**

▸ Additional functions may be available which can be specific to the database you are using

▸ A simple **COUNT(*)** will count the number of rows

# The HAVING statement applies conditions on a summary variable

```sql
SELECT
        Occupation,
        Gender,
        COUNT(CustomerKey) AS
        TotalCustomers
FROM dbo.DimCustomer
GROUP BY
        Occupation, Gender
HAVING COUNT(CustomerKey) >
        2000
```



| Results | Messages | | |
|---|---|---|---|
| | Occupation | Gender | TotalCustomers |
| 1 | Skilled Manual | M | 2293 |
| 2 | Professional | M | 2727 |
| 3 | Professional | F | 2793 |
| 4 | Skilled Manual | F | 2284 |

▸ **HAVING** clause works on aggregation of data rather than the rows of the data prior to the aggregation
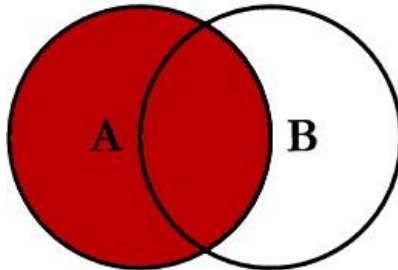
A **SUBQUERY** is a query on the data that is found within another query statement

```
SELECT
        FirstName,
        EmailAddress
FROM dbo.DimCustomer
WHERE GeographyKey IN (
        SELECT GeographyKey
        FROM dbo.DimGeography
        WHERE   CityName =
        'Beverley Hills');
```
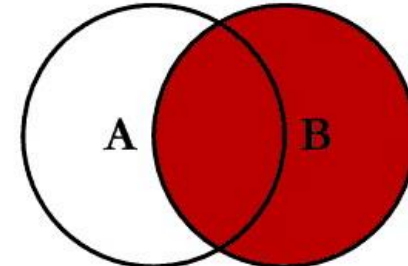
| | FirstName | EmailAddress |
|---|---|---|
| 1 | Luis | luis24@adventure-works.com |
| 2 | Molly | molly18@adventure-works.com |
| 3 | Carson | carson17@adventure-works.com |
| 4 | Lauren | lauren35@adventure-works.com |
| 5 | Rafael | rafael26@adventure-works.com |
| 6 | Nina | nina7@adventure-works.com |
| 7 | Darrell | darrell4@adventure-works.com |
| 8 | Jada | jada15@adventure-works.com |
| 9 | Chad | chad13@adventure-works.com |
| 10 | Logan | logan33@adventure-works.com |

▸ The sub query can only return one column at a time and the inner query gets executed before the outer one
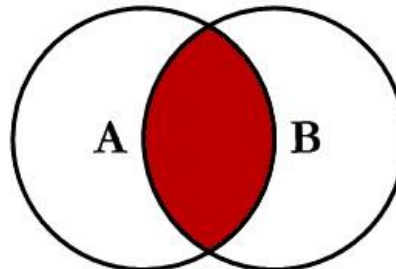
**JOIN keyword is used to query data from two or more tables based on a relationships between the columns in these tables.**
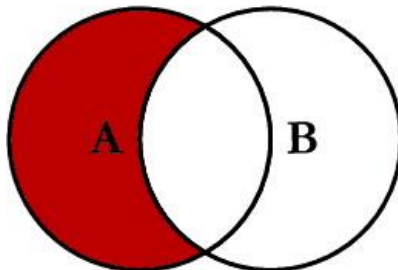


SELECT <select_list>
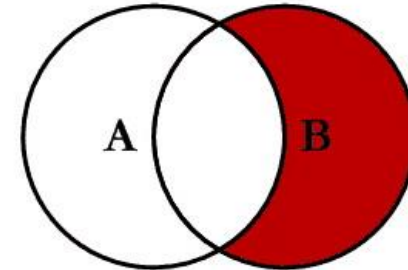FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key

SELECT <select_list>
FROM TableA A
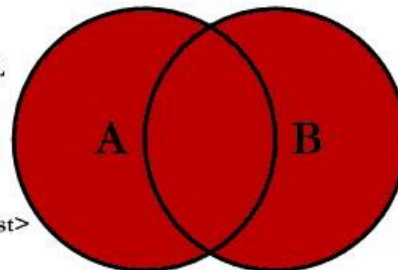INNER JOIN TableB B
ON A.Key = B.Key

SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
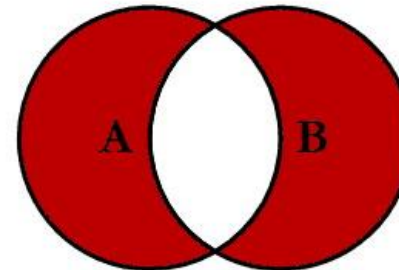
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
WHERE B.Key IS NULL

SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL

SELECT <select_list>
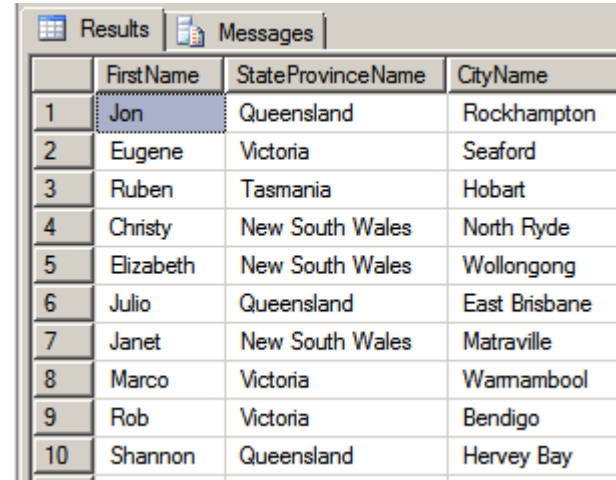FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key

SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL
OR B.Key IS NULL

# An SQL join can merge multiple tables at a time

```
SELECT
        t1.FirstName,
        t2.StateProvinceName,
        t2.CityName
FROM dbo.DimCustomer AS t1
INNER JOIN dbo.DimGeography
AS t2 ON (t1.GeographyKey =
t2.GeographyKey)
```

| Results | Messages | | |
|---|---|---|---|
| | First Name | StateProvinceName | CityName |
| 1 | Jon | Queensland | Rockhampton |
| 2 | Eugene | Victoria | Seaford |
| 3 | Ruben | Tasmania | Hobart |
| 4 | Christy | New South Wales | North Ryde |
| 5 | Elizabeth | New South Wales | Wollongong |
| 6 | Julio | Queensland | East Brisbane |
| 7 | Janet | New South Wales | Matraville |
| 8 | Marco | Victoria | Warrnambool |
| 9 | Rob | Victoria | Bendigo |
| 10 | Shannon | Queensland | Hervey Bay |

▸ The advantage of using table alias (t1, t2 etc.) is that merging column can be of different names in different tables

▸ Joins can happen on multiple columns too

# Further reading and practice

▶ Data warehousing concepts ([Link](#))

▶ Interactive SQL Tutorial ([Link](#))

▶ Beginning SQL Server 2008 for Developers (Chapters 11 and 12) ([Link](#))

▶ SQL Server 2012: T-SQL Fundamentals ([Link](#))

▶ SQL Server 2008: T-SQL Querying ([Link](#))