



# ROSETTA SQL



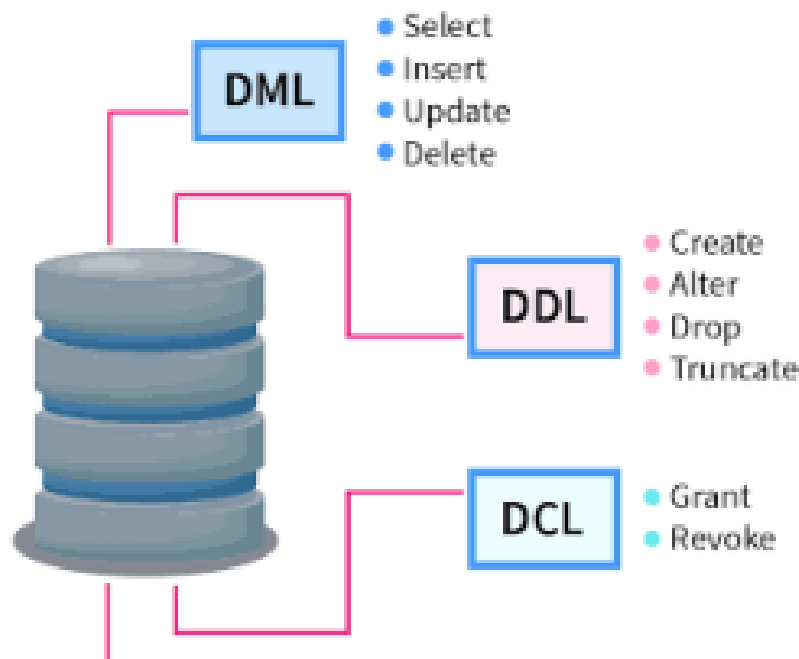
A close-up of Will Smith looking slightly to the left with a thoughtful expression.

**Can an AI  
write efficient  
SQL queries?**

A close-up of a white, bald AI robot with large eyes, looking directly forward.

**Can you?**





## DML - Data Manipulation Language:

Command	Description
SELECT	Retrieves certain records from one or more tables
INSERT	Creates a record
UPDATE	Modifies records
DELETE	Deletes records

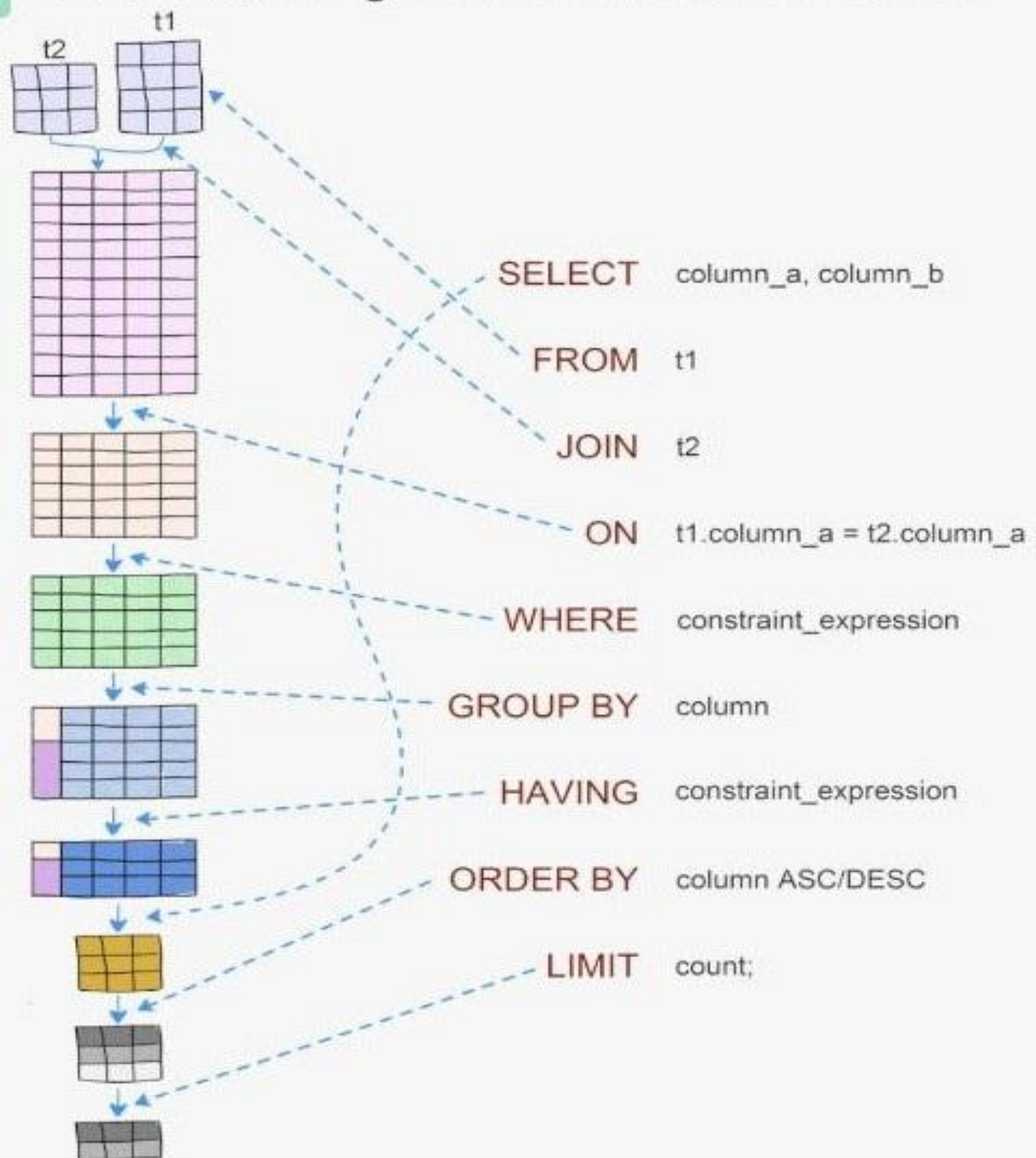
## DDL - Data Definition Language:

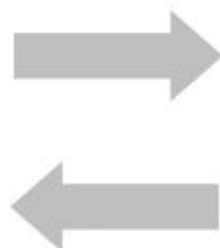
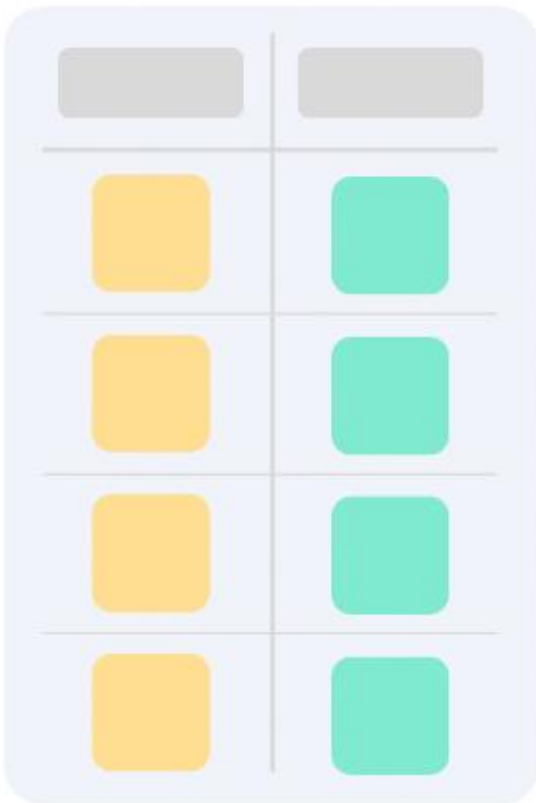
Command	Description
CREATE	Creates a new table, a view of a table, or other object in database
ALTER	Modifies an existing database object, such as a table.
DROP	Deletes an entire table, a view of a table or other object in the database.

## DCL - Data Control Language:

Command	Description
GRANT	Gives a privilege to user
REVOKE	Takes back privileges granted from user

# SQL Query Execution Order





SQLite

```
1 #Check table in SQLite
2 pd.read_sql_query ('select * from employees', con)
```

	index	empid	firstname	lastname
0	0	12985	Michael	Scott
1	1	12986	Dwight	Schrute
2	2	12987	Jim	Halpert
3	3	12988	Pam	Beesly
4	4	12989	Andy	Bernard
5	5	12990	Kevin	Malone
6	6	12991	Toby	Flenderson
7	7	12992	Angela	Martin

# Select

Selecting data with Pandas vs SQL



<code>df</code>	<b>SELECT * FROM tab;</b>	query data from all columns
<code>df[['col_1', 'col_2']]</code>	<b>SELECT col_1, col_2 FROM tab;</b>	select subset of columns and get all rows
<code>df.assign(new_col = df['col_1'] / df['col_2'])</code>	<b>SELECT *, col_1/col_2 as new_col FROM tab;</b>	Aliases - add a calculated column based on other columns
<code>df[['col 1']]</code>	<b>SELECT `col 1` FROM tab;</b>	Column name with space - `col 1`
<code>df.sort_values(by='col_1', ascending=True)</code>	<b>SELECT * FROM tab ORDER BY col_1 ASC;</b>	Sort values by col_1 in ascending or descending (DESC) order

# Where

Filtering in SQL vs Pandas



<code>df[df['col_1'] == '11']</code>	<b>SELECT * FROM tab</b> <b>WHERE col_1 = '11';</b>	Filtering on single condition
<code>df[(df['col_1'] == 11) &amp; (df['col_2'] &gt; 5)]</code>	<b>SELECT * FROM tab</b> <b>WHERE col_1 = 11 AND col_2 &gt; 5;</b>	Filtering on multiple conditions
<code>df[df['col_2'].isna()]</code>	<b>SELECT *</b> <b>FROM tab</b> <b>WHERE col_2 IS NULL;</b>	NULL checking is done using the isna()
<code>df[df['col_1'].notna()]</code>	<b>SELECT *</b> <b>FROM tab</b> <b>WHERE col_1 IS NOT NULL;</b>	NOT NULL checking is done using the notna()
<code>df[df.col_1 &gt; df.col_2]</code>	<b>SELECT * FROM tab</b> <b>WHERE col_1 &gt; col_2;</b>	Where clause with 2 SQL columns
<code>df.query('col_1 == col_2')</code>	<b>SELECT * FROM tab</b> <b>WHERE col_1 = col_2;</b>	Filter with Pandas Query
<code>df.query('col_1 == `col 2`')</code>	<b>SELECT * FROM tab</b> <b>WHERE col_1 = `col 2`;</b>	Column name with space - `col 2` in where clause



## Like, and, or

Operators(Text, Logical) in Pandas vs SQL

<b>contains</b>	→	<b>LIKE</b>
<b>startswith</b>	→	
<b>endswith</b>	→	
<b>&amp;</b>	→	<b>AND</b>
<b> </b>	→	<b>OR</b>

<pre>df[df['col_1'].str.contains('i', na=False)]</pre>	<b>WHERE col_1 LIKE '%i%'</b>	Finds values which contain `i`. Column need to be string - .astype(str)
<pre>df[df['col_1'].str.contains('sh rd', regex=True, na=True)]</pre>	<b>WHERE col_1 LIKE '%sh%'</b> <b>OR col_1 LIKE '%rd%'</b>	Finds any values which contain `sh` or `rd`
<pre>df[df['col_1'].str.startswith('h', na=False)]</pre>	<b>WHERE col_1 LIKE 'h%'</b>	Finds any values that start with `h`
<pre>df[df['col_1'].astype(str).str.endswith('k', na=False)]</pre>	<b>WHERE col_1 LIKE '%k'</b>	Finds any values that ends with `k`
<pre>(df['col_1'] == '11') &amp; (df['col_2'] &gt; 5)</pre>	<b>WHERE col_1 = '11' AND col_2 &gt; 5;</b>	AND = SQL - `and`, Pandas - `&`
<pre>(df['col_1'] == '11')   (df['col_2'] &gt; 5)</pre>	<b>WHERE col_1 = '11' OR col_2 &gt; 5;</b>	OR = SQL - `or`, Pandas - ` `
<pre>df[df['col_1'].isin([1,2,3])]</pre>	<b>SELECT * FROM tab</b> <b>WHERE col_1 in (1,2,3);</b>	IN operator - find values from list of values
<pre>df[df['col_1'].between(1, 5)]</pre>	<b>SELECT * FROM tab</b> <b>WHERE col_1 BETWEEN 1 AND 5;</b>	BETWEEN operator - find values in a range



# Group by

Group by operations in SQL vs Pandas



<pre>df.groupby('col_1').size() df.groupby('col_1')['col_2'].count() df.col_1.value_counts()</pre>	<b>SELECT col_1, count(*)</b> <b>FROM tab</b> <b>GROUP BY col_1;</b>	count records in each group( 3 versions in Pandas )
<pre>import numpy as np df.groupby('col_1').agg({'col_2': np.mean, 'col_1': np.size})</pre>	<b>SELECT col_1, AVG(col_2), COUNT(*)</b> <b>FROM tab</b> <b>GROUP BY col_1;</b>	Apply multiple statistical functions
<pre>df.groupby(['col_1', 'col_2']).agg({'col_3': [np.size, np.mean]})</pre>	<b>SELECT col_1, col_2, COUNT(*),</b> <b>AVG(col_3)</b> <b>FROM tab</b> <b>GROUP BY col_1, col_2;</b>	Grouping by multiple columns, multiple functions
<pre># group by g = df.groupby('col_1') # having count(*) &gt; 10 g.filter(lambda x: len(x) &gt; 10) ['col_1']</pre>	<b>SELECT col_1, count(*)</b> <b>FROM tab</b> <b>GROUP BY col_1</b> <b>HAVING count(*) &gt; 10;</b>	HAVING - Group by column and filtering contidion on the groups
<pre>df.groupby('col_1').col_2.nunique()</pre>	<b>SELECT count(distinct col_2)</b> <b>FROM tab</b> <b>GROUP BY col_1;</b>	count(distinct) - count unique elements in group

# Join

Join in SQL and Pandas



```
pd.merge(df1, df2, on='key')
```

```
SELECT *  
FROM t1  
INNER JOIN t2  
ON t1.key = t2.key;
```

Inner join of 2 table/dataframes(1)

```
pd.merge(df1, df2, on='key',  
how='left')
```

```
SELECT *  
FROM t1  
LEFT OUTER JOIN t2  
ON t1.key = t2.key;
```

Left Outer join

```
pd.merge(df1, df2, on='key',  
how='right')
```

```
SELECT *  
FROM t1  
RIGHT OUTER JOIN t2  
ON t1.key = t2.key;
```

Right Join

```
pd.merge(df1, df2, on='key',  
how='outer')
```

```
SELECT *  
FROM t1  
FULL OUTER JOIN t2  
ON t1.key = t2.key;
```

Full Join (not working on MySQL)

```
pd.merge(df1, df2, left_on= ['col_1',  
'col_2'], right_on= ['col_3', 'col_4'],  
how = 'right')
```

```
SELECT *  
FROM t1  
INNER JOIN t2  
ON t1.col_1 = t2.col_3  
AND t1.col_2 = t2.col_4;
```

Join on columns with different names

```
m = pd.merge(df1, df2, how='left', on=  
['col_1', 'col_2'])  
pd.merge(m, df3[['col_1', 'col_2',  
'col_3']], how='left', on=['col_1',  
'col_3'])
```

```
SELECT t1.col_a, t2.col_b, t3.col_c  
FROM t1  
LEFT OUTER JOIN t2  
ON t1.col_1 = t2.col_1  
AND t1.col_2 = t2.col_2  
LEFT OUTER JOIN t3  
ON t1.col_1 = t3.col_1  
AND t1.col_3 = t3.col_3
```

join multiple dataframes on multiple columns

# Union

Union in SQL and Pandas



```
pd.concat([df1, df2])
```

```
SELECT *  
FROM df1  
UNION ALL  
SELECT *  
FROM df2;
```

Union All (columns must have same number of columns)

```
cols= ['col_1', 'col_2']  
pd.concat([df1[cols], df2[cols]])
```

```
SELECT col_1, col_2  
FROM t1  
UNION ALL  
SELECT col_1, col_2  
FROM t2;
```

Union All

```
pd.concat([df1, df2]).drop_duplicates()
```

```
SELECT col_1, col_2  
FROM t1  
UNION  
SELECT col_1, col_2  
FROM t1;
```

Union All ( remove duplicate rows)

# Limit

Limit in SQL and Pandas



```
df.head(10)
```

**SELECT \* FROM tab  
LIMIT 10;**

Get top rows

```
df.tail(10)
```

**SELECT \* FROM tab  
ORDER BY id DESC  
LIMIT 10;**

Get last N rows

```
lim = 2  
offset = 5  
df.sort_values('col_1',  
ascending=False).iloc[offset:lim+offset  
]
```

**SELECT \* FROM tab  
ORDER BY col\_1 DESC  
LIMIT 2 OFFSET 5;**

return only top 2 records, start on record 6  
(OFFSET 5)

# Update

Update in SQL vs Pandas



```
df.loc[df['col_1'] < 2, 'col_1'] *= 2
```

**UPDATE tab  
SET col\_1 = col\_1\*2  
WHERE col\_1 < 2;**

Update all rows for 1 column with condition

```
df1['name'] =  
np.where(df2['id']==1,df2['name'],df1['  
name'])
```

**UPDATE t1,  
(  
SELECT \* FROM t2 WHERE id = 1  
) AS temp  
SET t1.name = temp.name  
WHERE t1.id = 1;**

Update based on select from another table /  
dataframe

## Delete

Delete in SQL vs Pandas



```
df = df.loc[df['col_1'] > 9]
```

**DELETE FROM tab**  
**WHERE col\_1 > 9;**

Delete rows with condition

```
df1.drop(df1[(df1.id.isin(df2.id) &  
(df1.id==1))].index)
```

**DELETE t1**  
**FROM df1 as t1**  
**JOIN df2 as t2 ON t1.id = t2.id**  
**WHERE t2.id = 1;**

Delete rows with condition based on another  
table / dataframe

## Insert

Insert in SQL vs Pandas



```
data = {'col_1': 1, 'col_2': '11'}  
df = df.append(data, ignore_index =  
True)
```

**INSERT INTO tab(col\_1, col\_2)**  
**VALUES (1, '11');**

Add new data/rows to a table/dataframe



Operation	Pandas	SQL
Read CSV	<code>pd.read_csv(file)</code>	LOAD DATA INFILE 'data.csv' INTO TABLE <code>table</code> FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 ROWS;
Print first 10 (or k) rows	<code>df.head(10)</code>	SELECT * FROM <code>table</code> LIMIT 10;
Dimensions	<code>df.shape</code>	SELECT count(*) FROM <code>table</code> ;
Datatype	<code>df.dtypes</code>	DESCRIBE <code>table</code> ;
Filter Data	<code>df[df.column&gt;10]</code>	SELECT * FROM <code>table</code> where <code>column&gt;10</code> ;
Select column(s)	<code>df.column</code>	SELECT <code>column</code> FROM <code>table</code> ;
Sort	<code>df.sort_values("column")</code>	SELECT * FROM <code>table</code> ORDER BY <code>column</code> ;
Fill NaN	<code>df.column.fillna(0)</code>	UPDATE <code>table</code> SET <code>column=0</code> WHERE <code>column</code> IS NULL;
Join	<code>pd.merge(df1, df2, on = "col", how = "inner")</code>	SELECT * FROM <code>table1</code> JOIN <code>table2</code> ON ( <code>table1.col = table2.col</code> );
Concatenate	<code>pd.concat((df1, df2))</code>	SELECT * FROM <code>table1</code> UNION ALL <code>table2</code> ;
Group	<code>df.groupby("column"). agg_col.mean()</code>	SELECT <code>column</code> , avg( <code>agg_col</code> ) FROM <code>table</code> GROUP BY <code>column</code> ;
Unique values	<code>df.column.unique()</code>	SELECT DISTINCT <code>column</code> FROM <code>table</code> ;
Rename column	<code>df.rename(columns = {"old_name": "new_name"})</code>	ALTER TABLE <code>table</code> RENAME COLUMN <code>old_name</code> TO <code>new_name</code> ;
Delete column	<code>df.drop(columns = ["column"])</code>	ALTER TABLE <code>table</code> DROP COLUMN <code>column</code> ;

# Pandas equivalent of SQL

Table below shows the most used equivalents between SQL and Pandas:

CATEGORY	PANDAS	SQL
Data Structure	DataFrames, Series	Tables, Row, Column
Querying	.loc , .iloc , boolean indexing	SELECT, WHERE
Filtering	.query() , boolean indexing	WHERE
Sorting	.sort_values() , .sort_index()	ORDER BY
Grouping	.groupby() , .agg()	GROUP BY, AGGREGATE
Joins	.merge() , .join(), .concat()	JOIN, UNION
Aggregations	.agg() , .apply()	AGGREGATE
Data Transformation	.apply() , .map() , .replace()	replace( <span>column</span> , 'old', 'new')
Count Unique	.nunique(), .agg(['count', 'nunique'])	count(distinct)
Data Cleaning	.dropna() , .fillna()	IS NULL; IS NOT NULL;
Data Type Conversion	.astype()	CAST



EXIT 12

Precise  
Data  
Queries



SELECT \*  
FROM  
everything

EAST ST WEST

Me, an SQL "expert"