

# Session 9

## Relational Databases (SQL III)



# Last week: SQL

id	code	mark
S103	DBS	72
S103	IAI	58
S104	PR1	68
S104	IAI	65
S106	PR2	43
S107	PR1	76
S107	PR2	60
S107	IAI	35

- **SELECT** - return the specified column(s) from a table
- **WHERE** - filters row(s) by condition
- **JOIN... ON... USING...** Join two tables together ON (where) the given condition is true
- **GROUP BY** - creates new entities (rows) by grouping entities and computing specified group (aggregate) facts.
- **HAVING** - filters the result of group by

id	code	mark
S103	DBS	72
S103	IAI	58
S105	PR1	68



# This week: Advanced SQL

Advanced

SELECT, WHERE,  
GROUP BY, HAVING...  
statements

Date and Time  
functions

## Real world examples

How to write longer  
queries...



# Advanced SELECT

Let's consider the way a  
basic SELECT  
SQL statement is executed  
by the computer...

grade		
id	code	mark
S103	DBS	72
S103	IAI	58
S104	PR1	68



```
SELECT id, mark  
FROM grade
```

- 1) Computer starts at first record  
*(where is first? random from your point of view, ordering not guaranteed)*

# Advanced SELECT

Let's consider the way a basic SELECT SQL statement is executed by the computer...

grade		
id	code	mark
S103	DBS	72
S103	IAI	58
S104	PR1	68

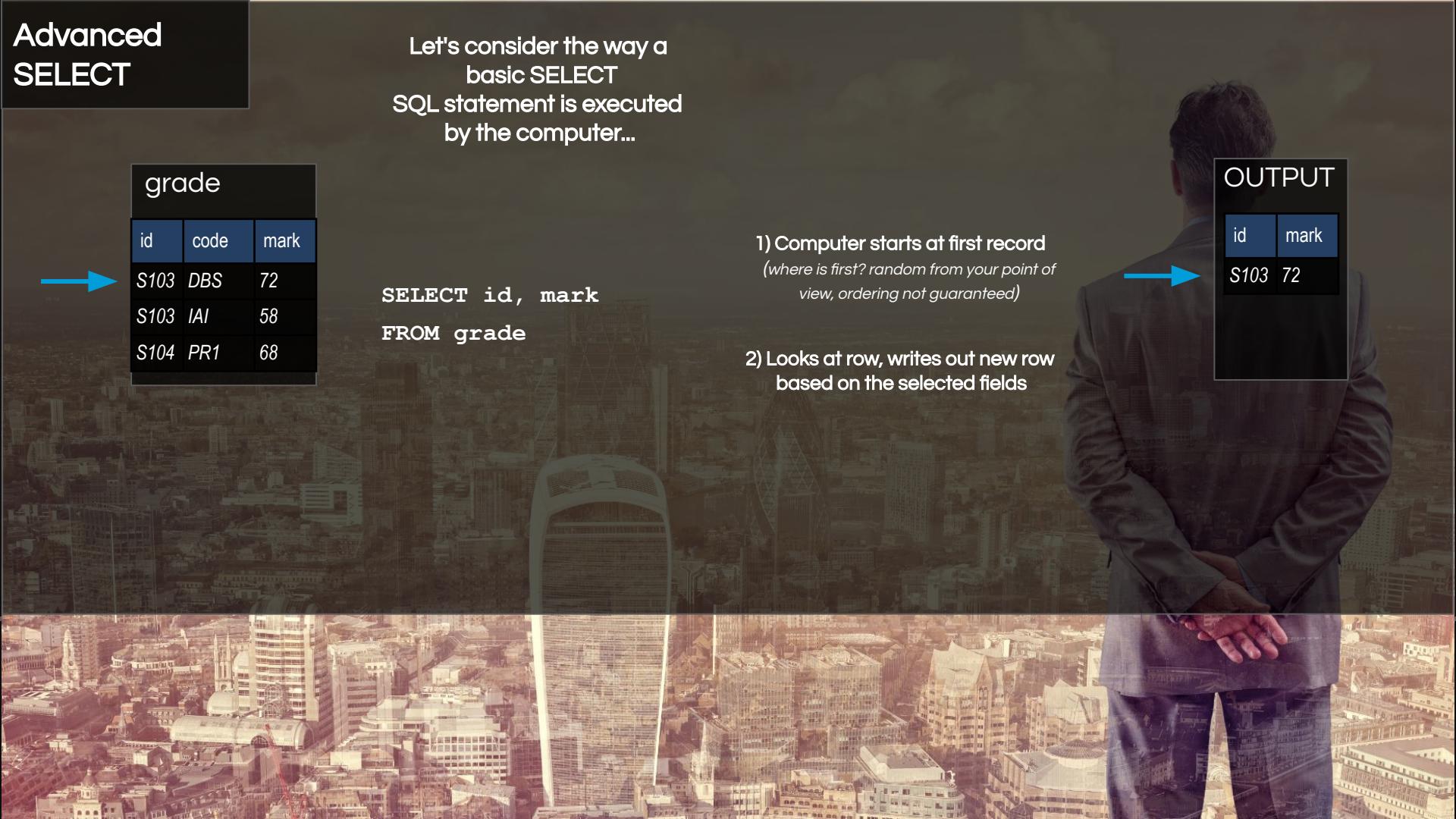
```
SELECT id, mark  
FROM grade
```

1) Computer starts at first record  
*(where is first? random from your point of view, ordering not guaranteed)*

2) Looks at row, writes out new row based on the selected fields

OUTPUT

id	mark
S103	72



# Advanced SELECT

Let's consider the way a basic SELECT SQL statement is executed by the computer...

grade		
id	code	mark
S103	DBS	72
S103	IAI	58
S104	PR1	68

```
SELECT id, mark  
FROM grade
```

1) Computer starts at first record  
(where is first? random from your point of view, ordering not guaranteed)

2) Looks at row, writes out new row based on the selected fields

OUTPUT

id	mark
S103	72
S103	58

# Advanced SELECT

Let's consider the way a basic SELECT SQL statement is executed by the computer...

grade		
id	code	mark
S103	DBS	72
S103	IAI	58
S104	PR1	68

```
SELECT id, mark  
FROM grade
```

## Take home message:

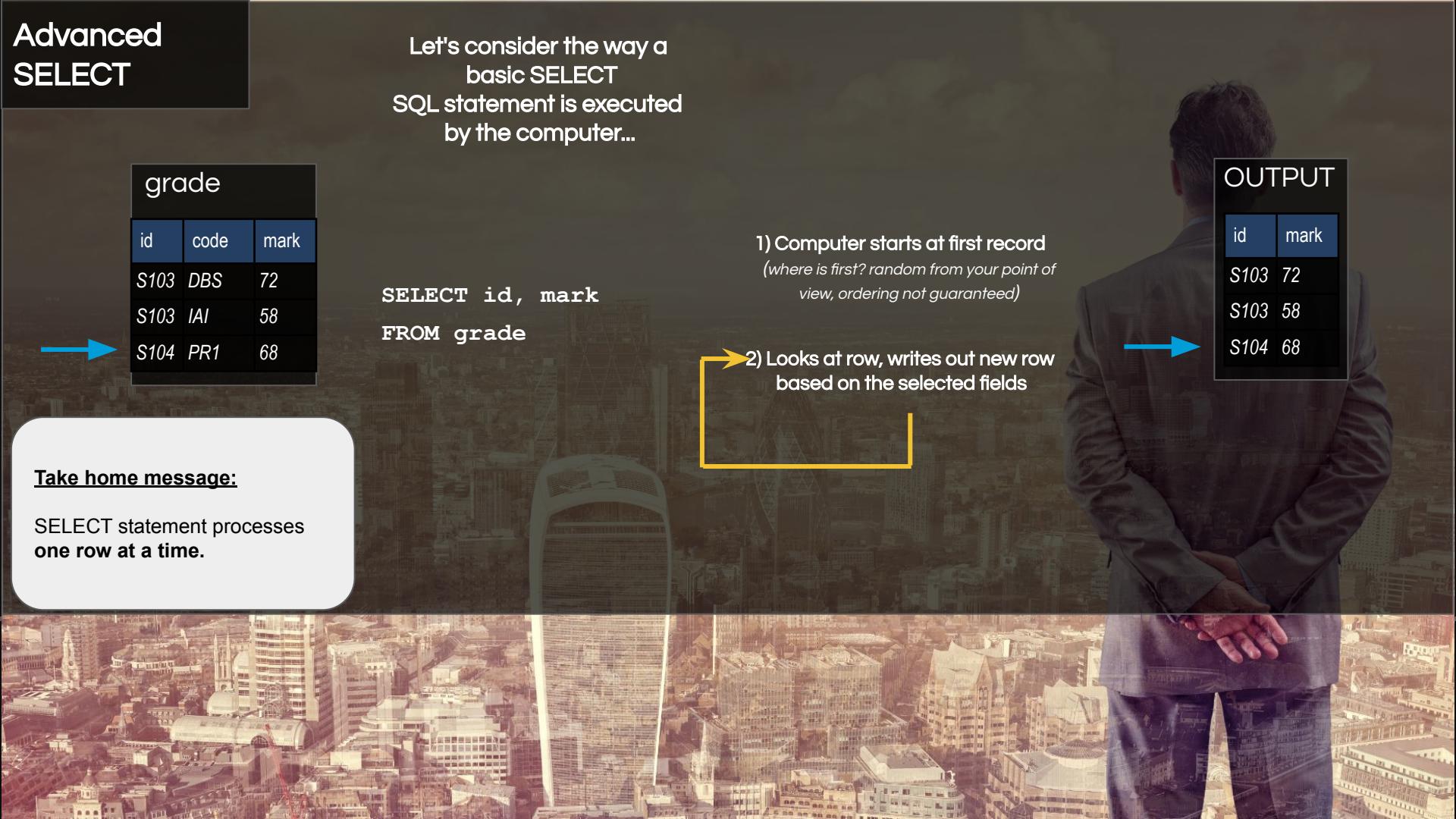
SELECT statement processes one row at a time.

1) Computer starts at first record  
(where is first? random from your point of view, ordering not guaranteed)

2) Looks at row, writes out new row based on the selected fields

## OUTPUT

id	mark
S103	72
S103	58
S104	68



# Advanced SELECT

## Take home message:

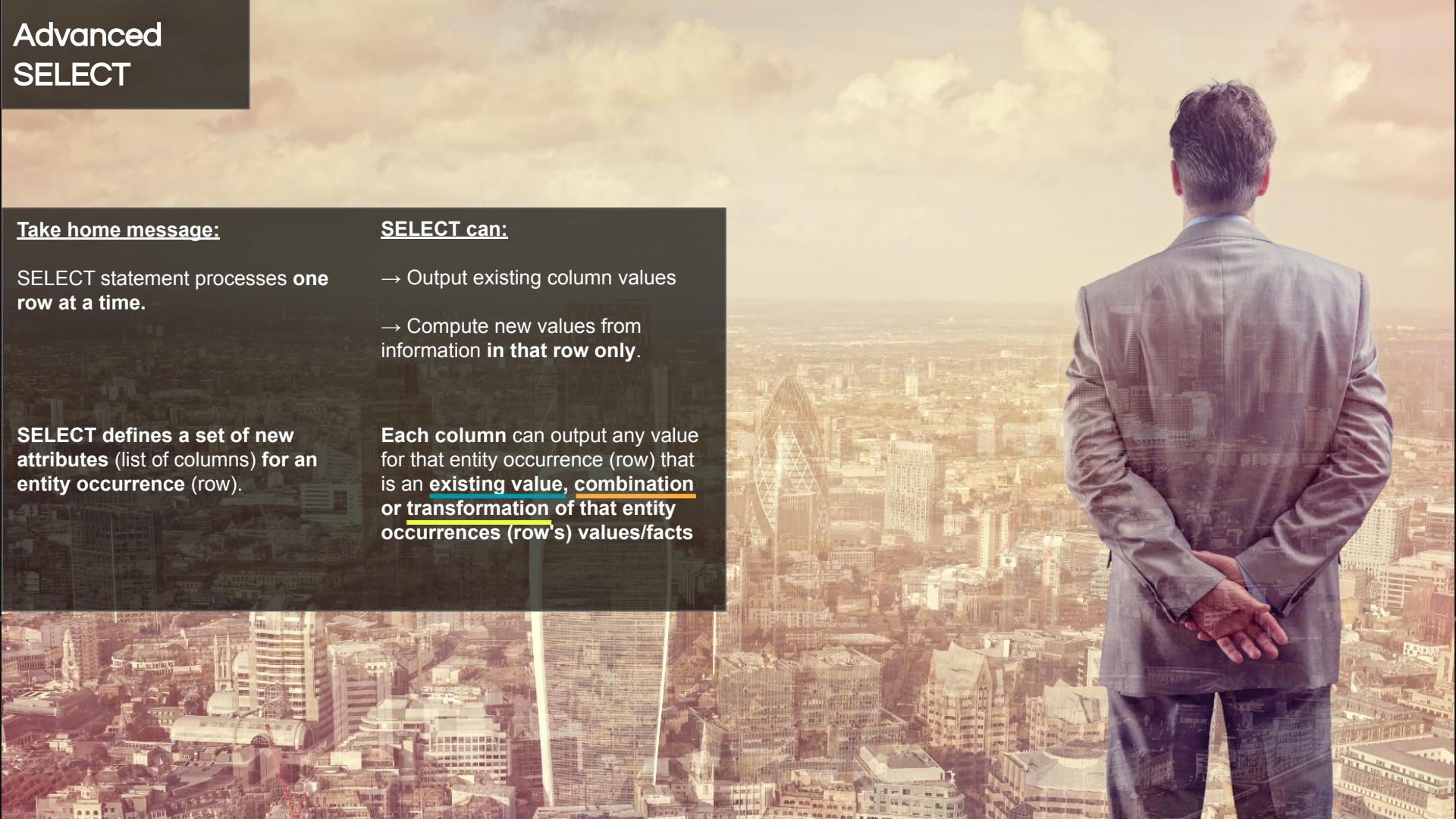
SELECT statement processes **one row at a time.**

**SELECT defines a set of new attributes** (list of columns) for an entity occurrence (row).

## SELECT can:

- Output existing column values
- Compute new values from information **in that row only.**

**Each column** can output any value for that entity occurrence (row) that is an **existing value, combination or transformation of that entity occurrences (row's) values/facts**



# Advanced SELECT

## Take home message:

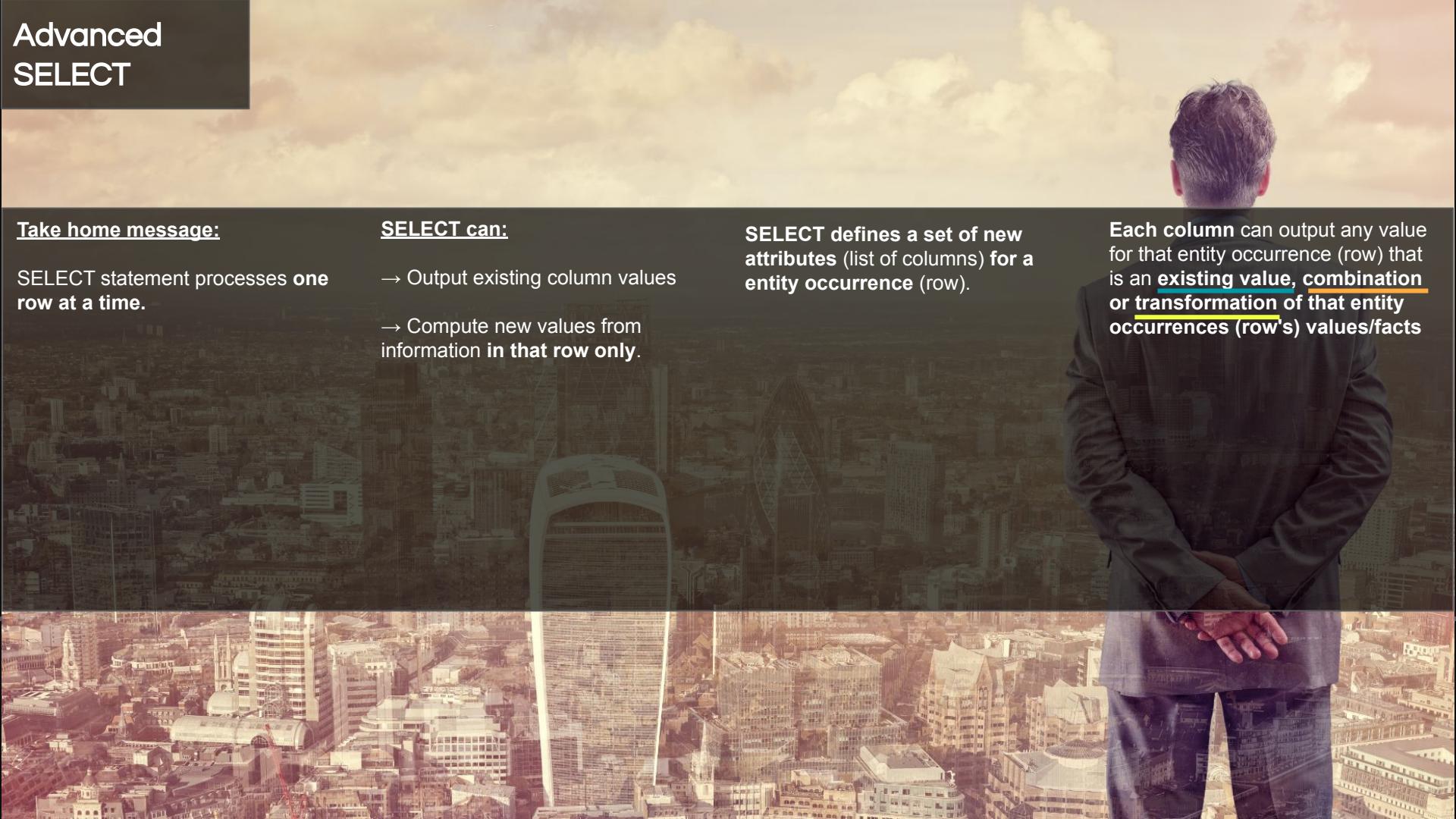
SELECT statement processes **one row at a time**.

## SELECT can:

- Output existing column values
- Compute new values from information **in that row only**.

**SELECT defines a set of new attributes** (list of columns) **for a entity occurrence** (row).

**Each column** can output any value for that entity occurrence (row) that is an **existing value, combination or transformation of that entity occurrences** (row's) values/facts



# Advanced SELECT

## Take home message:

SELECT statement processes **one row at a time.**

Can output any value for that entity occurrence (row) that is an existing value, combination or transformation of that entity occurrences (row's) values/facts

### Example:

Existing  
value

grade

	grade		
	id	code	mark1
→	S103	DBS	72
	S103	IAI	58
	S104	PR1	68
			0

```
SELECT id, mark1  
FROM grade
```

OUTPUT

id	mark1
S103	72



# Advanced SELECT

## Take home message:

SELECT statement processes **one row at a time.**

Can output any value for that entity occurrence (row) that is an existing value, combination or transformation of that entity occurrences (row's) values/facts

Example:  
**Combination** →

grade			
id	code	mark1	mark2
S103	DBS	72	70
S103	IAI	58	61
S104	PR1	68	0

```
SELECT code || ':' || mark1::STRING AS mylabel  
FROM grade
```

OUTPUT

mylabel  
DBS: 72

# Advanced SELECT

## Take home message:

SELECT statement processes **one row at a time.**

Can output any value for that entity occurrence (row) that is an existing value, combination or transformation of that entity occurrences (row's) values/facts

Example:  
**Combination** →

grade

	grade			
	id	code	mark1	mark2
S103	DBS	72	70	
S103	IAI	58	61	
S104	PR1	68	0	

```
SELECT code || ':' || mark1::STRING AS mylabel  
FROM grade
```

|| is the concatenation operator  
(we've seen this before!)

::STRING converts the mark  
(which was an INTEGER) to  
text before concatenation.

OUTPUT

mylabel  
DBS: 72

Also seen this before too!  
If you omit this, the computer will  
assume you meant this in this  
case. Not always the case though.  
So good practice to include it.

# Advanced SELECT

## Take home message:

SELECT statement processes one row at a time.

Can output any value for that entity occurrence (row) that is an existing value, combination or transformation of that entity occurrences (row's) values/facts

### Why 2.0 not 2?

Not strictly required for SparkSQL. For others, i.e. PostgreSQL you do... Why?

Since mark1 and mark2 are INTEGER fields, if you divide by 2 (another INTEGER) the computer will do INTEGER division. I.e.  $1/2 = 0$  We make one a numeric to ensure normal division occurs.

## Example:

Existing value and Combination →

grade			
id	code	mark1	mark2
S103	DBS	72	70
S103	IAI	58	61
S104	PR1	68	0

```
SELECT id, (mark1 + mark2) / 2.0 AS avgmark  
FROM grade
```

## OUTPUT

id	avgmark
S103	71.0
S103	59.5

# Advanced SELECT

## Take home message:

SELECT statement processes one row at a time.

**REMEMBER: Each entity occurrence is processed separately.**

i.e. we go line-by-line

Example:

Existing value and Combination

grade

	id	code	mark1	mark2
→	S103	DBS	72	70
	S103	IAI	58	61
	S104	PR1	68	0

```
SELECT id, (mark1 + mark2) / 2.0 AS avgmark
FROM grade
```

OUTPUT

	id	avgmark
→	S103	71.0
	S103	59.5

# Advanced SELECT

## Take home message:

SELECT statement processes one row at a time.

Can output any value for that entity occurrence (row) that is an existing value, combination or transformation of that entity occurrences (row's) values/facts

**Transformations** are via **FUNCTIONS**.

Functions take a set of values and parameters and return a new value.

postgres define a number of functions, or you can write your own.

A common type of builtin functions are formatting functions.

i.e. we have a date, and want to format it differently.

Again, we've seen this before. There are many more transformations.

student

Example:

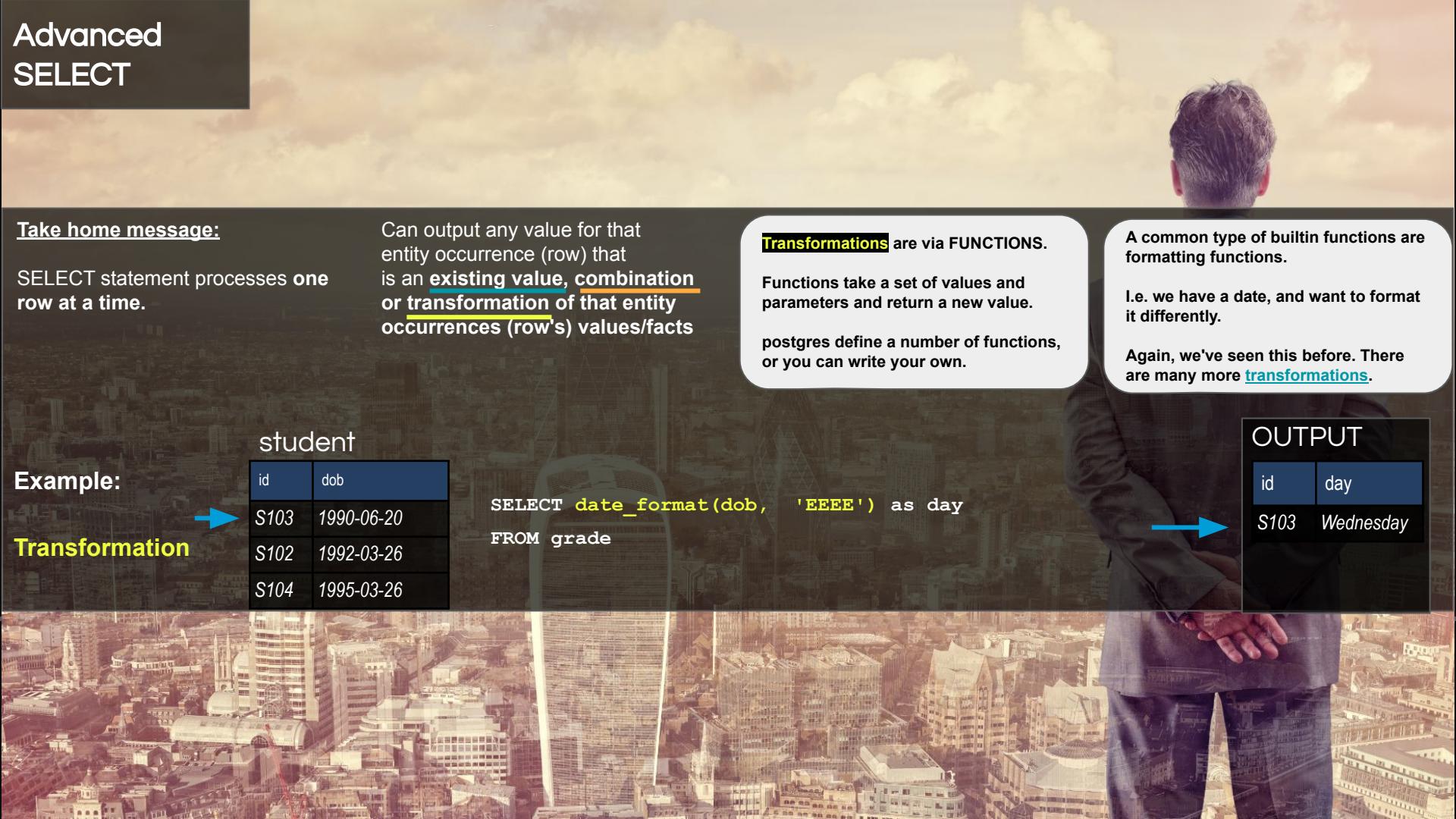
Transformation

	id	dob
→	S103	1990-06-20
	S102	1992-03-26
	S104	1995-03-26

```
SELECT date_format(dob, 'EEEE') as day  
FROM grade
```

OUTPUT

	id	day
→	S103	Wednesday



# Advanced SELECT

But I don't know what functions there are... If you think one should exist, then:  
→ we try google... e.g. "[SparkSQL format date](#)"  
→ we go directly to the [documentation for SparkSQL](#)

## Take home message:

SELECT statement processes one

Spark SQL Built-In Functions

Search docs

Functions

!  
!=  
%  
&  
\*  
+  
-  
/  
<  
<=

## Transformation

date\_format

date\_format(timestamp, fmt) - Converts timestamp to a value of string in the format specified by the date format fmt.

### Arguments:

- timestamp - A date/timestamp or string to be converted to the given format.
- fmt - Date/time format pattern to follow. See [Datetime Patterns](#) for valid date and time format patterns.

### Examples:

```
> SELECT date_format('2016-04-08', 'y');  
2016
```

Since: 1.5.0

Can output any value for that entity occurrence (row) that is an **existing value combination**

**Transformations** are via **FUNCTIONS**.

Functions take a set of values and parameters and return a new value.

postgres define a number of functions, or you can write your own.

A common type of builtin functions are formatting functions.

i.e. we have a date, and want to format it differently

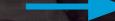
y  
cts

date\_format(dob, 'EEEE') as day

OUTPUT

id	day
S103	Wednesday

S102	1992-03-26
S104	1995-03-26



# Advanced SELECT

## Take home message:

SELECT statement processes one row at a time.

Can output any value for that entity occurrence (row) that is an existing value, combination or transformation of that entity occurrences (row's) values/facts

## Example:

student

id	dob
S103	1990-06-20
S102	1992-03-26
S104	1995-03-26

## Transformation

```
SELECT date_format(dob, 'EEEE') as day  
FROM grade
```

## Datetime Patterns for Formatting and Parsing

There are several common scenarios for datetime usage in Spark:

- CSV/JSON datasources use the pattern string for parsing and formatting datetime content.
- Datetime functions related to convert StringType to/from DateType or TimestampType. For example, unix\_timestamp, date\_format, to\_unix\_timestamp, from\_unixtime, to\_date, to\_timestamp, from\_utc\_timestamp, to\_utc\_timestamp, etc.

Spark uses pattern letters in the following table for date and timestamp parsing and formatting:

Symbol	Meaning	Presentation	Examples
G	era	text	AD; Anno Domini
y	year	year	2020; 20
D	day-of-year	number(3)	189
M/L	month-of-year	month	7; 07; Jul; July
d	day-of-month	number(2)	28
Q/q	quarter-of-year	number/text	3; Q3; Q3: 3rd quarter
E	day-of-week	text	Tue; Tuesday
F	aligned day of week in month	number(1)	
z	am/pm offset	am/pm	
Z	zone-offset	offset-Z	-0800; -08:00;
'	escape for text	delimiter	
"	single quote	literal	
[	optional section start		
]	optional section end		

The count of pattern letters determines the format.

- Text: The text style is determined based on the number of pattern letters used. Less than 4 pattern letters will use the short text form, typically an abbreviation, e.g. day-of-week Monday might output "Mon". Exactly 4 pattern letters will use the full text form, typically the full description, e.g. day-of-week Monday might output "Monday". 5 or more letters will fail.
- Number(n): The n here represents the maximum count of letters this type of datetime pattern can be used.
  - In formatting, if the count of letters is one, then the value is output using the minimum number of digits and without padding otherwise, the count of digits is used as the width of the output field, with the value zero-padded as necessary.
  - In parsing, the exact count of digits is expected in the input field.
- Number/Text: If the count of pattern letters is 3 or greater, use the Text rules above. Otherwise use the Number rules above.
- Fraction: Use one or more (up to 9) contiguous 'S' characters, e.g SSSSSS, to parse and format fraction of second. For parsing, the acceptable fraction length can be [1, the number of contiguous 'S']. For formatting, the fraction length would be padded to the number of contiguous 'S' with zeros. Spark supports datetime of micro-of-second

why  
4x  
Es?

# Advanced SELECT

```
SELECT id,  
       'Born on a ' || date_format(dob, 'EEEE') || ', started studying on a ' || date_format(start_date, 'EEEE') as info  
FROM student
```

**Example:**  
**Existing value,**  
**combination &**  
**transformation**

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

OUTPUT

id	info
S103	Born on a Wednesday, started studying on a Sunday

# Advanced SELECT

**REMEMBER:** Each entity occurrence is processed separately.

i.e. we go line-by-line

```
SELECT id,  
       'Born on a ' || date_format(dob, 'EEEE') || ', started studying on a ' || date_format(start_date, 'EEEE') as info  
FROM student
```

Example:  
**Existing value,  
combination &  
transformation** →

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

OUTPUT

id	info
S103	Born on a Wednesday, started studying on a Monday
S102	Born on a Thursday, started studying on a Wednesday

# Advanced SELECT

**REMEMBER:** Each entity occurrence is processed separately.

i.e. we go line-by-line

```
SELECT id,  
       'Born on a ' || date_format(dob, 'EEEE') || ', started studying on a ' || date_format(start_date, 'EEEE') as info  
FROM student
```

**Example:**  
**Existing value,**  
**combination &**  
**transformation** →

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

OUTPUT

id	info
S103	Born on a Wednesday, started studying on a Monday
S102	Born on a Thursday, started studying on a Wednesday
S104	Born on a Sunday, started studying on a Wednesday

# Advanced WHERE

OK, so that was SELECT.

What about WHERE?

## Recall:

- Filters (keeps or omits) rows based on a given condition.

Same as SELECT, rows are processed independently.

## **SELECT:**

- we are deciding what to output given the facts for an entity occurrence

## **WHERE:**

- we are deciding whether to *keep* the occurrence entity given the facts for that entity occurrence

A WHERE condition can therefore include:

- existing value,
- combination or
- transformations

of that entity occurrences (row's values/facts).

As with SELECT, it can also be a mix of these.

# Advanced WHERE

```
SELECT *  
FROM student  
WHERE date_format(start_date, 'EEEE') = 'Wednesday' AND NOT id = 'S102'
```

Example:  
**Existing value**  
&  
**Transformation**

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

id	dob	start_date
S103	1990-06-20	2001-07-02

2001-07-02 is NOT a Wednesday

# Advanced WHERE

**REMEMBER:** Each entity occurrence is processed separately.

I.e. we go line-by-line

```
SELECT *  
FROM student  
WHERE date_format(start_date, 'EEEE') = 'Wednesday' AND NOT id = 'S102'
```

Example:  
**Existing value**  
&  
**Transformation**

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

id	dob	start_date
S103	1990-06-20	2001-07-02
<b>S102</b>	1992-03-26	2010-06-30

Record is S102

# Advanced WHERE

**REMEMBER:** Each entity occurrence is processed separately.

I.e. we go line-by-line

```
SELECT *  
FROM student  
WHERE date_format(start_date, 'EEEE') = 'Wednesday' AND NOT id = 'S102'
```

Example:  
**Existing value**  
&  
**Transformation** →

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

Record is not S102 & 2011-08-03  
is a Wednesday!

# Advanced GROUP BY

OK, so that was SELECT.  
And WHERE.

What about GROUP BY?

## Recall:

→ GROUP BY specifies the column for which all distinct values will become bucket labels

More generally, however,  
GROUP BY defines the measure for which entity occurrences (rows) will be considered the same for the purpose of grouping.

As such, per occurrence (row), this measure can be:

- existing value,
- combination or
- transformations

EXACTLY THE SAME THINGS AS:  
→ SELECT  
→ WHERE

# Advanced GROUP BY

```
SELECT date_format(start_date, 'EEEE') as day, COUNT(*) AS ct
FROM student
GROUP BY date_format(start_date, 'EEEE')
```

Example:  
Transformation

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03



# Advanced GROUP BY

**Recall:** When using GROUP BY  
SELECT can only have bucket labels  
or  
results of aggregate functions.

```
SELECT date_format(start_date, 'EEEE') as day, COUNT(*) as ct
FROM student
GROUP BY date_format(start_date, 'EEEE')
```

student

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

**Example:  
Transformation**



# Advanced GROUP BY

**Recall:** When using GROUP BY  
SELECT can only have bucket labels  
or  
results of aggregate functions.

```
SELECT date_format(start_date, 'EEEE') as day, COUNT(*) as ct
FROM student
GROUP BY date_format(start_date, 'EEEE')
```

**Example:**  
**Transformation**

id	dob	start_date
S103	1990-06-20	2001-07-02
S102	1992-03-26	2010-06-30
S104	1995-03-26	2011-08-03

**NOTE** the difference:

## Aggregate functions:

Take a group of entity occurrences (rows), return a value summarizing the group.

## Transformation functions:

Transform facts/values for a single entity occurrence (row)

# Advanced GROUP BY

```
SELECT date_format(start_date, 'EEEE') as day, COUNT(*) as ct
FROM student
GROUP BY date_format(start_date, 'EEEE')
```

Example:  
Transformation

id	dob	start_date	day
S103	1990-06-20	2001-07-02	Monday
S102	1992-03-26	2010-06-30	Wednesday
S104	1995-03-26	2011-08-03	Wednesday



Monday	S103	1990-06-20	2001-07-02
--------	------	------------	------------



Wednesday	S102	1992-03-26	2010-06-30
	S104	1995-03-26	2011-08-03



day	ct
Monday	1
Wednesday	2

# Advanced HAVING

OK, so that was SELECT.

And WHERE.  
And GROUP BY.

What about HAVING?

Recall:

→ HAVING simply filters the rows  
after the GROUP BY has run

WHERE & HAVING therefore both:

→ decide whether to *keep* the  
occurrence entity given the facts for  
that entity occurrence (row)

WHERE:

entity occurrence = rows in original tables

HAVING

entity occurrence = rows created as a  
result of a GROUP BY part of a query

# Advanced HAVING

OK, so that was SELECT.  
And WHERE.  
And GROUP BY.

What about HAVING?

Recall:

→ HAVING simply filters the rows  
after the GROUP BY has run

WHERE & HAVING therefore both:

→ decide whether to *keep* the  
occurrence entity given the facts for  
that entity occurrence (row)

WHERE:

entity occurrence = rows in original tables

HAVING

entity occurrence = rows created as a  
result of a GROUP BY part of a query

As with everything else, it  
can also be a mix of these.

A HAVING condition can therefore  
include:

→ existing value  
(in the table after the GROUP BY is applied),

→ combination or

→ transformations

of that entity occurrences (row's  
values/facts).

**REMEMBER:** the row are the rows  
in the table formed AFTER the  
group by!!!!!!



*Data at Scale*

Dr Georgiana Nica-Avram

# Advanced HAVING

```
SELECT date_format(start_date, 'EEEE') as day, COUNT(*)  
FROM student  
GROUP BY date_format(start_date, 'EEEE')  
HAVING COUNT(*) + 10 = 11
```

student

**Example:**  
**Combination**  
**(Transformation)**

	<b>id</b>	<b>dob</b>	<b>start_date</b>	<b>day</b>
	S103	1990-06-20	2001-07-02	Monday
	S102	1992-03-26	2010-06-30	Wednesday
	S104	1995-03-26	2011-08-03	Wednesday



Monday

	<b>id</b>	<b>dob</b>	<b>start_date</b>
	S103	1990-06-20	2001-07-02



Wednesday

	<b>id</b>	<b>dob</b>	<b>start_date</b>
	S102	1992-03-26	2010-06-30
	S104	1995-03-26	2011-08-03

day  
Monday

<b>day</b>	<b>count(*)</b>
Monday	1
Wednesday	2

# Summary

Advanced: SELECT,  
WHERE, GROUP BY, HAVING

→ All require the specification of per entity occurrence measures to use

This *may* be existing attributes.

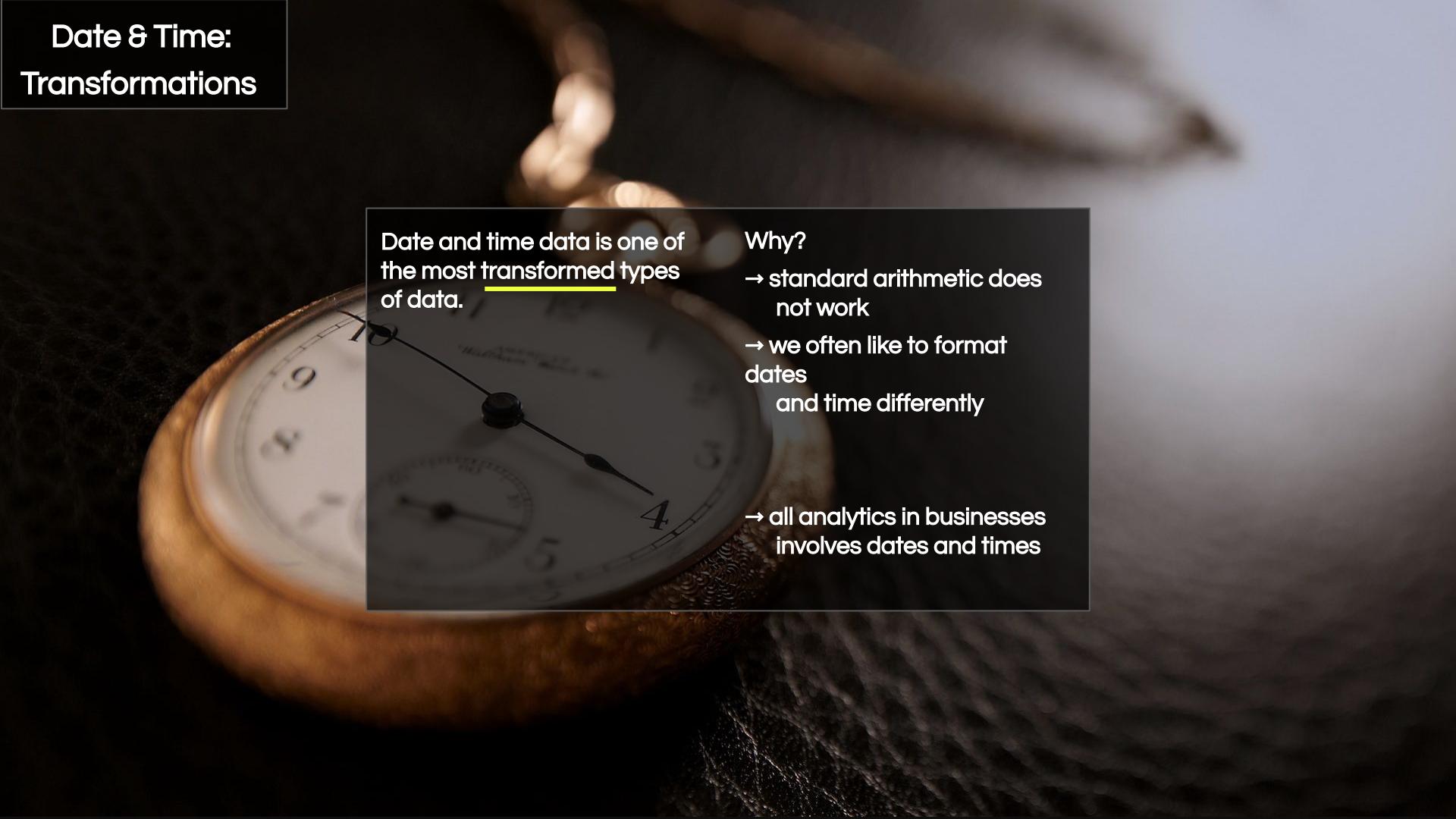
Or they may be new ones that are:

- combinations, or
- transformations



## Date & Time:

### Transformations



Date and time data is one of the most transformed types of data.

Why?

- standard arithmetic does not work
- we often like to format dates and time differently
- all analytics in businesses involves dates and times

# Date & Time:

## Transformations

How do I find out what functions I can use for problem X?

Answer 1:

- Google to find the documentation
- Google to find a tutorial / hints in solutions on stackoverflow

→ DO NOT try and copy paste a solution to a different problem and change things "until it works"

→ You will not learn, be uncertain if you have done it right. Likely a poor solution.

→ Use it as a hint to look up the documentation for functions (transformations) you don't understand.

→ Once you understand the functions you can CORRECTLY write your own

# Date & Time:

## Transformations

How do I find out what functions I can use for problem X?

Answer 2:  
Directly look up the documentation  
**[recommended]**

→ use the search box on the page

SparkSQL Functions

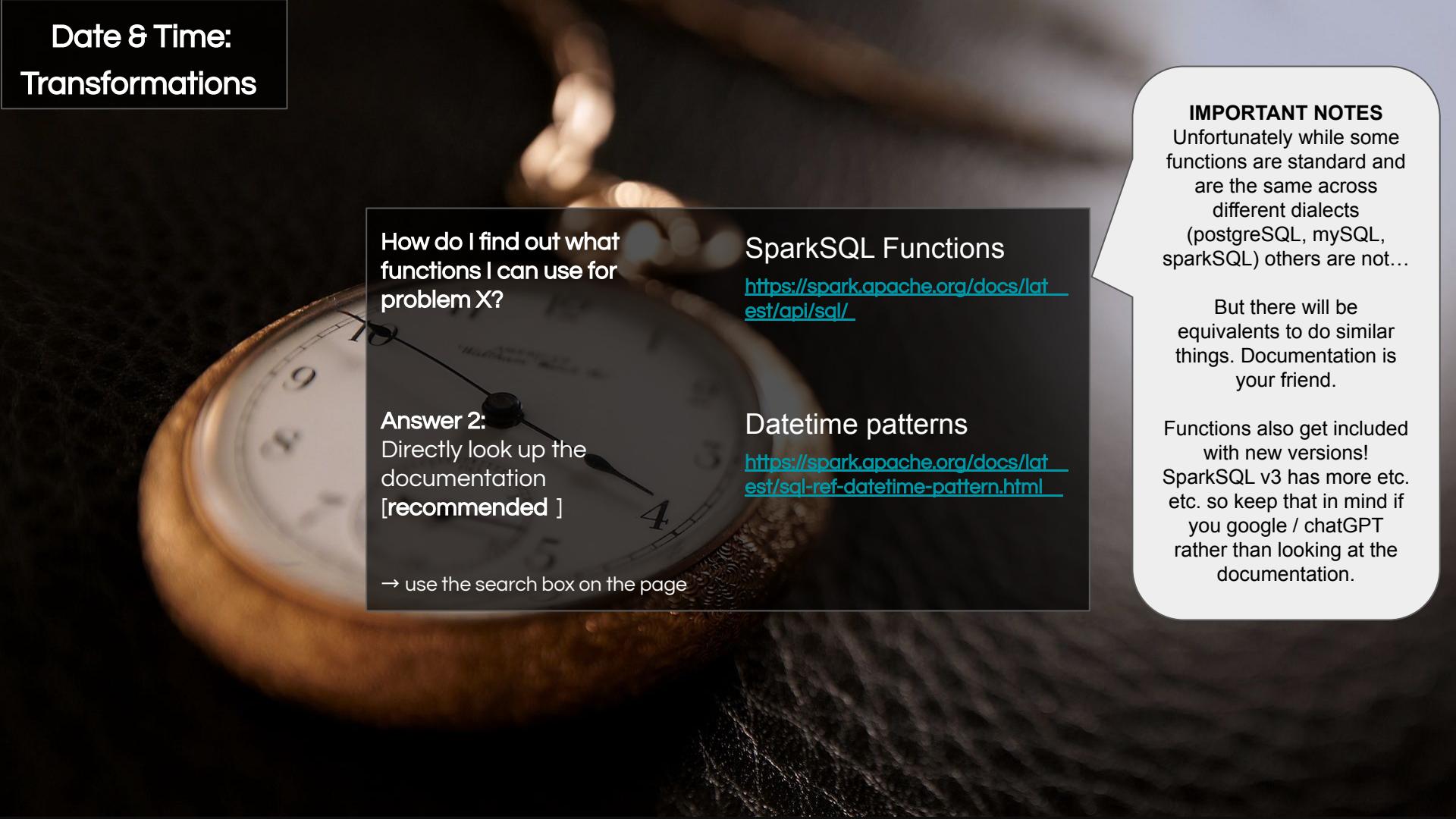
<https://spark.apache.org/docs/latest/api/sql/>

Datetime patterns

<https://spark.apache.org/docs/latest/sql-ref-datetime-pattern.html>

# Date & Time:

## Transformations



How do I find out what functions I can use for problem X?

Answer 2:

Directly look up the documentation  
**[recommended]**

→ use the search box on the page

### SparkSQL Functions

<https://spark.apache.org/docs/latest/api/sql/>

### Datetime patterns

<https://spark.apache.org/docs/latest/sql-ref-datetime-pattern.html>

### IMPORTANT NOTES

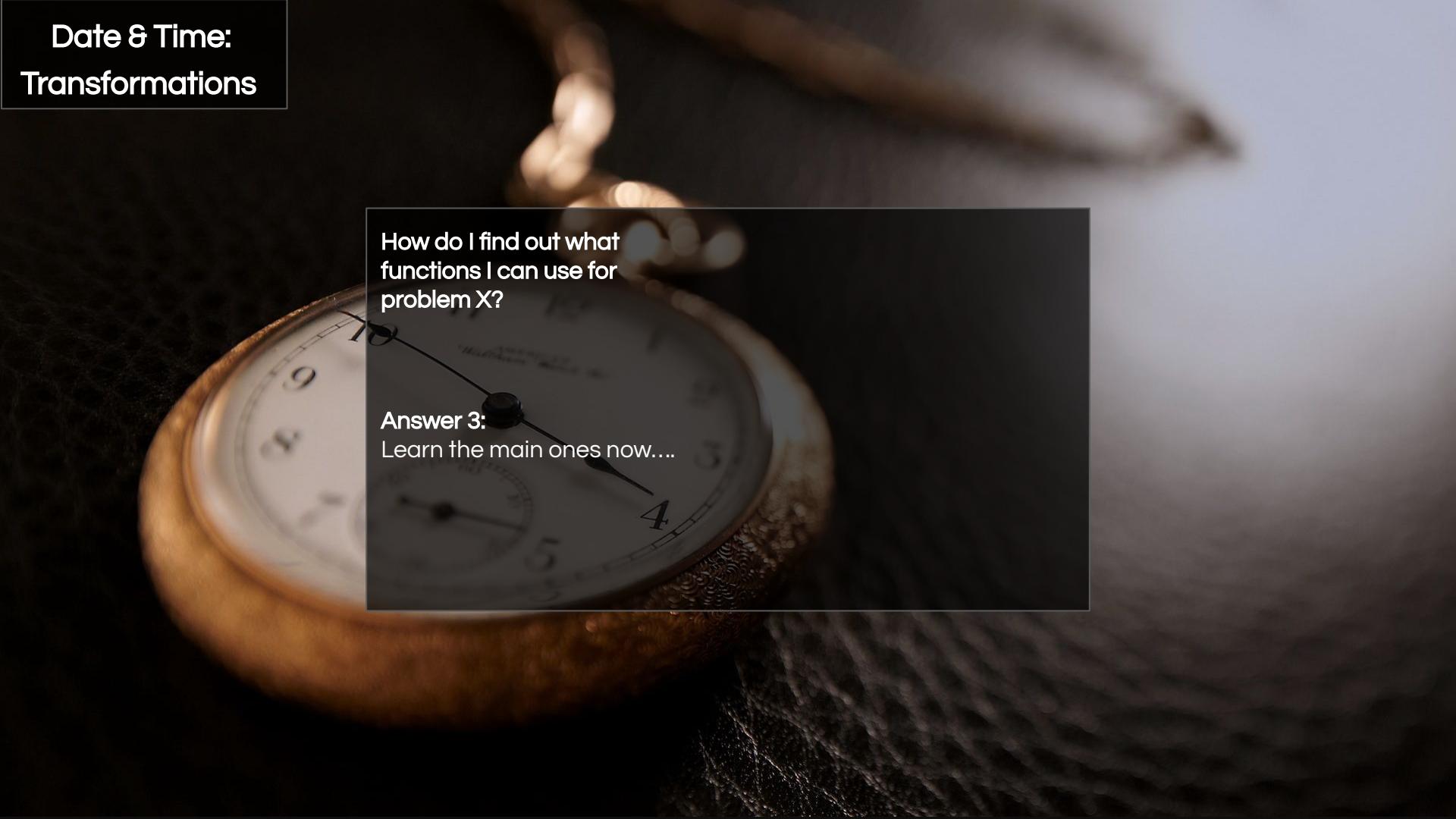
Unfortunately while some functions are standard and are the same across different dialects (postgreSQL, mySQL, sparkSQL) others are not...

But there will be equivalents to do similar things. Documentation is your friend.

Functions also get included with new versions! SparkSQL v3 has more etc. etc. so keep that in mind if you google / chatGPT rather than looking at the documentation.

## Date & Time:

## Transformations



How do I find out what functions I can use for problem X?

**Answer 3:**  
Learn the main ones now....

# Date & Time:

## Transformations

The core data types we have are:  
TIMESTAMP  
DATE

### The main date/time functions

DATE\_FORMAT(...)  
DATE\_DIFF(...)  
DATE\_ADD(...)  
DATE\_SUB(...)  
ADD\_MONTHS(...)  
MONTHS\_BETWEEN(...)  
CURRENT\_TIMESTAMP()  
EXTRACT(...)  
DATE\_TRUNC(...)

We already know  
this one!

## Difference (via functions)

### What it achieves: Computes the difference between two dates/times.

Function	Description	Example	Result
<code>DATE_DIFF(endDate, startDate)</code>	Compute the difference between TIMESTAMP or DATE objects. Returns <b>the number of days</b> .	<pre>SELECT DATE_DIFF('2007-06-06 10:00:00'::TIMESTAMP,           '1706-03-30 09:00:00'::TIMESTAMP);</pre>	<b>11006</b> <b>(type = INTEGER)</b>
<code>MONTHS_BETWEEN(endDate,startDate)</code>	Compute the difference between TIMESTAMP or DATE objects. Returns <b>the number of months</b> .	<pre>months_between('2007-03-02'::DATE,                '2007-01-01'::DATE));</pre>	<b>2.03225806</b> <b>(type = DOUBLE)</b>
<code>DATE_SUB(startDate, numDays)</code>	Subtract a given number of days from a DATE. Returns a DATE  WARNING: Using a TIMESTAMP here will currently cause a silent error!	<pre>SELECT date_sub('2016-07-30'::DATE, 1); SELECT date_sub('2016-07-30 23:00:00'::TIMESTAMP,                1);</pre>	<b>2016-07-29 (type = DATE)</b> <b>NULL (silent error - do not use!!)</b>

## Difference (via functions)

**What it achieves:**  
Computes the difference between two dates/times.

Function	Description	Example	Result
<code>DATE_DIFF(endDate, startDate)</code>	Compute the difference between TIMESTAMP or DATE objects. Returns <b>the number of days</b> .	<pre>SELECT DATE_DIFF('2007-06-06 10:00:00'::TIMESTAMP,            '1706-03-30 09:00:00'::TIMESTAMP);</pre>	<b>11006</b> <b>(type = INTEGER)</b>
<code>MONTHS_BETWEEN(endDate, startDate)</code>	Compute the difference between TIMESTAMP or DATE objects. Returns <b>the number of months</b> .	<pre>months_between('2007-03-02'::DATE,                '2007-01-01'::DATE));</pre>	<b>2.03225806</b> <b>(type = DOUBLE)</b>
<code>DATE_SUB(startDate, numDays)</code>	Subtract a given number of days from a DATE. Returns a DATE  WARNING: Using a TIMESTAMP here will currently cause a silent error!	<pre>SELECT date_sub('2016-07-30'::DATE, 1); SELECT date_sub('2016-07-30 23:00:00'::TIMESTAMP,                1);</pre>	<b>2016-07-29 (type = DATE)</b> <b>NULL (silent error - do not use!!)</b>

List the age (in years) of all customers

### Real world use:

```
SELECT cust_id, (MONTHS_BETWEEN(CURRENT_TIMESTAMP(), dob) /12)::INT as age
FROM customer
```

Casting to INT forces a round down (simply removes the decimal component) - which is what we want in this case.

## Difference (via subtract operator)

## What it achieves: Computes the difference between two dates/times.

Function	Description	Example	Result
<code>TIMESTAMP - TIMESTAMP</code>	<p>Compute the difference between TIMESTAMP.      Returns <b>an "INTERVAL DAY TO SECONDS"</b>.      Cast to BIGINT to get seconds between the timestamps.</p> <p>or DATE objects. Returns <b>an INTERVAL DAY</b>.      This is a special type</p>	<pre>SELECT '2007-06-06 10:00:00'::TIMESTAMP       - '1706-03-30 09:00:00'::TIMESTAMP;  SELECT ('2007-06-06 10:00:00'::TIMESTAMP -        '1706-03-30 09:00:00'::TIMESTAMP)::BIGINT</pre>	<b>Unreadable Object</b> <b>(type = INTERVAL DAY TO SECONDS)</b>  <b>9504522000</b> <b>(type = BIGINT)</b>
<code>DATE - DATE</code>	<p>Compute the difference between DATE.      Returns <b>an "INTERVAL DAY"</b>.      Cast to INTEGER to get days between the timestamps.</p>	<pre>SELECT '2007-06-06 10:00:00'::DATE       - '1706-03-30 09:00:00'::DATE;  SELECT ('2007-06-06 10:00:00'::DATE -        '1706-03-30 09:00:00'::DATE)::INTEGER</pre>	<b>Unreadable Object</b> <b>(type = INTERVAL DAY TO SECONDS)</b>  <b>11006</b> <b>(type = INTEGER)</b>
<code>TIMESTAMP - INTERVAL</code> <code>DATE - INTERVAL</code>  <b>NOTE: BE CAREFUL OF THE RETURN TYPE!</b>	<p>Subtract an interval of time from either a TIMESTAMP or DATE.      Returns a <b>DATE</b> (if originally DATE and subtracting days, months or years)  <b>or TIMESTAMP</b> otherwise.</p>	<pre>SELECT '2016-07-30 10:00:00'::TIMESTAMP -        INTERVAL 30 MINUTES  SELECT '2016-07-30'::DATE - INTERVAL 30 MINUTES  SELECT '2016-07-30'::DATE - INTERVAL 2 DAYS</pre>	<b>2016-07-30 09:30:00 (type = TIMESTAMP)</b>  <b>2016-07-29 23:30:00 (type = TIMESTAMP)</b>  <b>2016-07-28 (type = DATE)</b>

# Addition

## What it achieves:

Adds an INTERVAL to a  
DATE/TIME/INTERVAL/TIMESTAM  
P

Function	Description	Example	Result
<code>TIMESTAMP + INTERVAL</code>	Adds an INTERVAL of time to a TIMESTAMP or DATE	<pre>SELECT '2007-06-06 10:00:00'::TIMESTAMP + INTERVAL 10 days;</pre>	<b>2007-06-16 10:00:00 (type = TIMESTAMP)</b>
<code>DATE + INTERVAL</code>	<b>RETURNS a DATE</b> if starting with DATE and adding days, months, years <b>TIMESTAMP</b> otherwise	<pre>SELECT '2007-06-06'::TIMESTAMP + INTERVAL 10 days;</pre>	<b>2007-06-16 (type = DATE)</b>
<code>DATE_ADD(startDate, numDays)</code>	Add a given number of days from a DATE or TIMESTAMP. Returns a DATE (TIMESTAMP will be truncated!)  NOTE: Try not to use a TIMESTAMP, cast to DATE to show you know what is happening to those who read your code in the future!	<pre>SELECT date_add('2016-07-29'::DATE, 1); SELECT date_add('2016-07-29 23:00:00'::TIMESTAMP, 1);</pre>	<b>2016-07-30 (type = DATE)</b> <b>2016-07-30 (type = DATE)</b>

## CURRENT\_TIMESTAMP()

### What it achieves:

Gets the current time and date.

Want just the date? Cast!

CURRENT\_TIMESTAMP()::DATE =  
2018-10-19

Function	Description	Example	Result
<code>CURRENT_TIMESTAMP()</code>	Returns the current date and time. <b>RETURNS a TIMESTAMP</b>	<pre>SELECT CURRENT_TIMESTAMP();</pre>	<b>2023-10-21T15:11:46.071+0000 (type = TIMESTAMP)</b>

# EXTRACT

## What it achieves:

Extracts and returns PART of a  
DATE or TIMESTAMP

Function	Description	Example	Result
<code>EXTRACT(field FROM timestamp)</code>	Get subfield. Options: year, yearofweek, quarter, month, week, day, dayofweek, dayofweek_iso, dow, hour, minute, second See: <a href="https://spark.apache.org/docs/latest/api/sql/#extract">https://spark.apache.org/docs/latest/api/sql/#extract</a>	<code>extract(month from '2007-05-01'::TIMESTAMP)</code>	<b>5</b>
		<code>extract(dow from '2007-05-01'::DATE)</code>	<b>3</b>

# DATE\_TRUNC

Values of type DATE are cast automatically to TIMESTAMP

## What it achieves:

Truncates a DATE or TIMESTAMP to a given precision by setting lower precision elements to zero (or one for day and month).

Function	Description	Example	Result
<code>DATE_TRUNC(field, timestamp)</code>	Sets elements from source that are less than the precision specified in field to zero (or one for day & month). Options for field: microseconds, milliseconds, second, minute, hour, day, week, month, quarter, year See: <a href="https://spark.apache.org/docs/latest/api/sql/#date_trunc">https://spark.apache.org/docs/latest/api/sql/#date_trunc</a>	<pre>SELECT DATE_TRUNC('month', '2007-05-15'::DATE)</pre>	<b>2007-05-01 00:00:00</b>
		<pre>SELECT DATE_TRUNC('hour', '2007-05-01 09:30:05'::TIMESTAMP)</pre>	<b>2007-05-01 09:00:00</b>

# DATE\_TRUNC

Values of type DATE are cast automatically to TIMESTAMP

## What it achieves:

Truncates a DATE or TIMESTAMP to a given precision by setting lower precision elements to zero (or one for day and month).

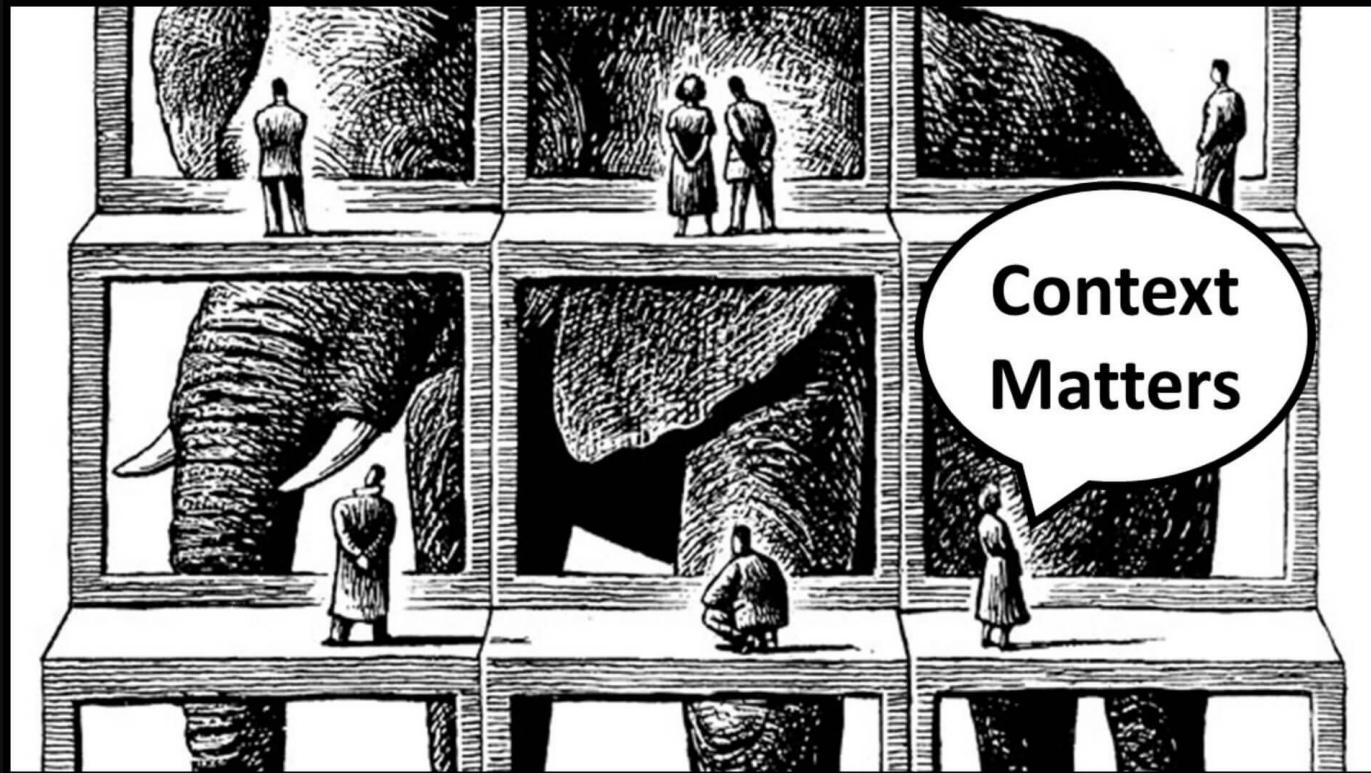
## Real world use:

```
SELECT DATE_TRUNC('month', receipt_timestamp) as year_month,  
       SUM(price) as total_sales  
  FROM line_item  
 GROUP BY DATE_TRUNC('month', receipt_timestamp)
```

Function	Description	Example	Result
<code>DATE_TRUNC(field, timestamp)</code>	Sets elements from source that are less than the precision specified in field to zero (or one for day & month). Options for field: microseconds, milliseconds, second, minute, hour, day, week, month, quarter, year See: <a href="https://spark.apache.org/docs/latest/api/sql/#date_trunc">https://spark.apache.org/docs/latest/api/sql/#date_trunc</a>	<pre>SELECT DATE_TRUNC('month', '2007-05-15'::DATE)</pre>	<b>2007-05-01 00:00:00</b>
		<pre>SELECT DATE_TRUNC('hour', '2007-05-01 09:30:05'::TIMESTAMP)</pre>	<b>2007-05-01 09:00:00</b>

Return the total sales per month.

*Let's put this all in  
practice...*



*Data at Scale*

Dr Georgiana Nica-Avram

*Let's put this all in  
practice...*



## **Let's put this all in practice...**

Image you've just been employed as a data analyst by UNREALISTIC CORP. The company is a startup which sells coffee products by delivery over the Internet.

UNREALISTIC CORP. has employed you to help develop their customer loyalty program. The team developing the program has asked for answers to the following questions:

- 1) *List the customers by name who shop with us more than once a month?*
- 2) *How many customers shop with us more than once a month?*

## **Let's put this all in practice...**

- 1) *List the customers by name who shop with us more than once a month?*
  
- 2) *How many customers shop with us more than once a month?*

To get you started UNREALISTIC CORP. provides you with a full set of documentation about their database.

They also introduce you to their data quality team which has infinite resources who can 100% guarantee there are no errors in the database.

# The documentation



## customers

Field Name	Data Type	Constraint	Default	Description	Example
customer_id	INTEGER	Primary Key	auto-incremet	Unique identifier generated for each customer	8734
name	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
order_id	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
order_ts	TIMESTAMP	NOT NULL	current_timestamp()	Timestamp of when the order was placed	2017-01-01 13:01:02
customer_id	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## items

Field Name	Data Type	Constraint	Default	Description	Example
item_id	INTEGER	Primary Key	auto-increment	Unique identifier generated for each item	545
description	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
id	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
order_id	INTEGER	Foreign Key ( orders.order_id ), NOT NULL		The order_id for the order which this line item belongs.	123456
item_code	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
cost	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
qty	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

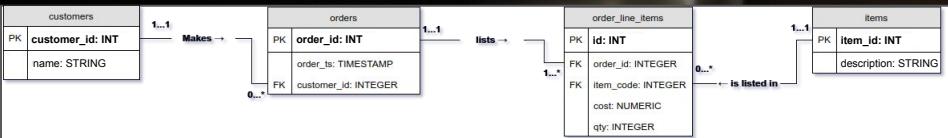
**Let's put this all in  
practice...**

- 1) *List the customers by name who shop with us more than once a month while they were active customers?*
  
- 2) *How many customers shop with us more than once a month?*

**Since we can guarantee there are no errors in the data, we can get started on the queries!**

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

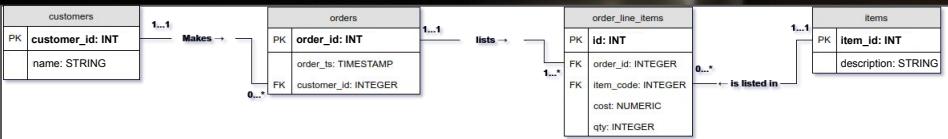
## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

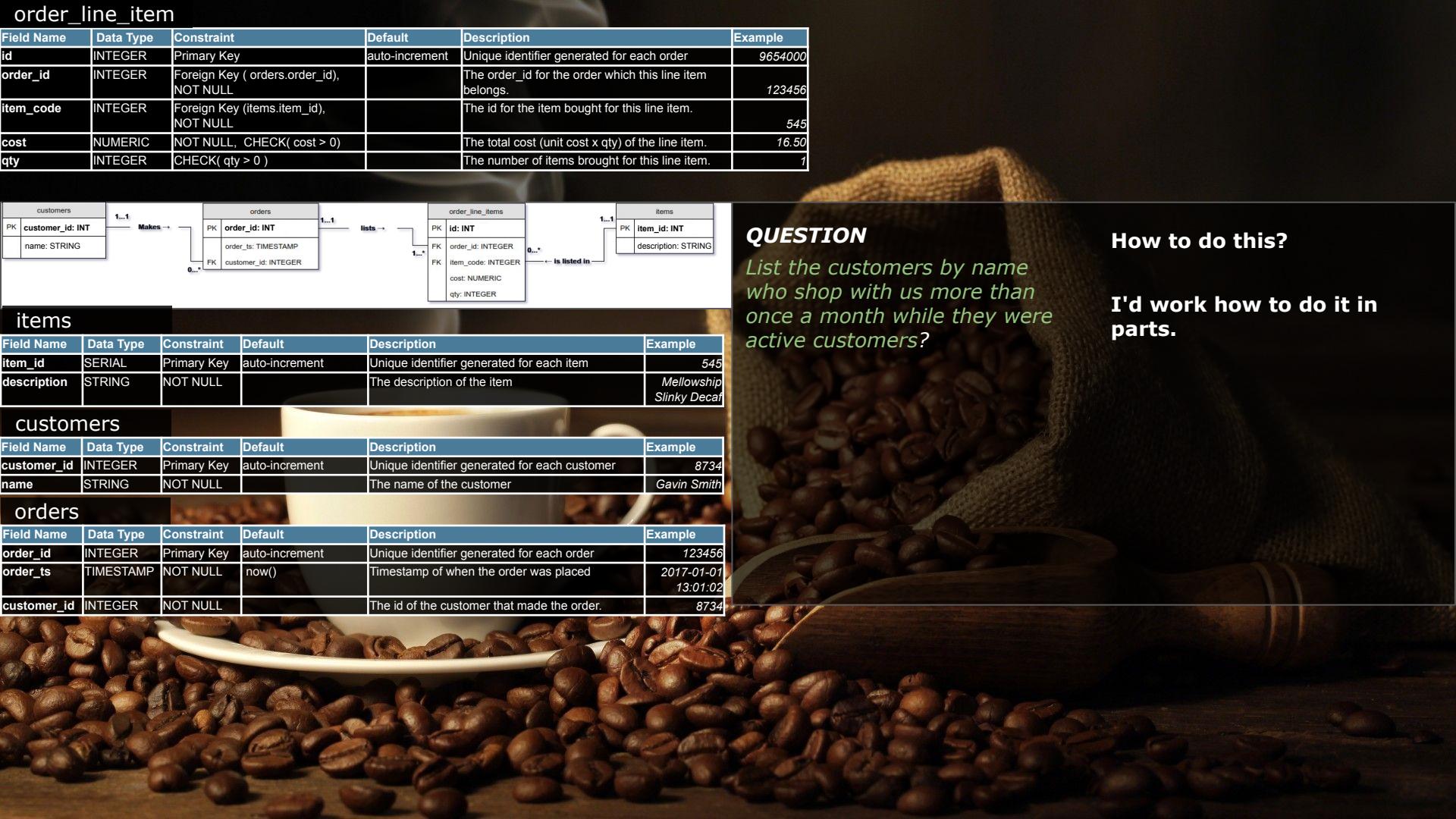
Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

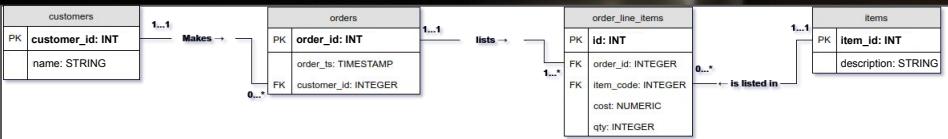
## How to do this?

I'd work how to do it in parts.



## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0)		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## QUESTION

List the customers by name who **shop with us more than once a month** while they were active customers?

## How to do this?

I'd work how to do it in parts.

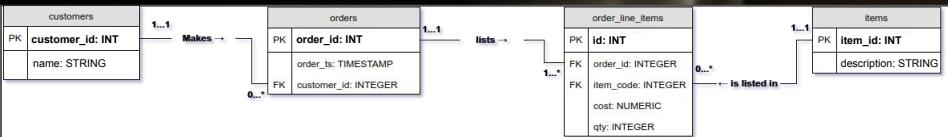
SELECT

FROM **customers**

WHERE

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## QUESTION

List the customers by name who **shop with us more than once a month** while they were active customers?

SELECT

FROM customers, orders

WHERE

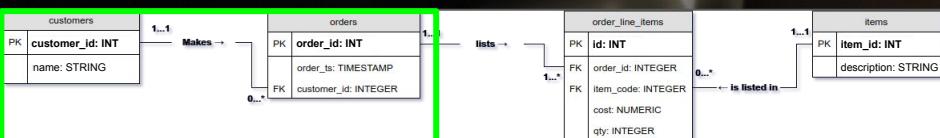
Ah! So we have a join to deal with!

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

We have a one-to-many join.

This is the **most common** by far.

SELECT

```

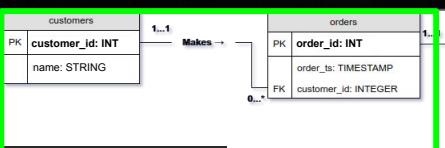
  FROM orders
  JOIN customers
  ON orders.customer_id =
  customers.customer_id
  
```

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

We have a one-to-many join.

This is the **most common** by far.

One interpretation is:

Each **row in orders** gets **extended** with (the corresponding information from the **customer** table) **name** (each row in orders matches exactly one row in customer).

SELECT

FROM **orders**

JOIN **customers**

```

ON orders.customer_id =
customers.customer_id
  
```

So my logic:

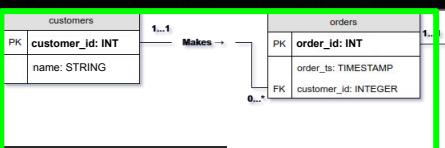
- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

We have a one-to-many join.

This is the **most common** by far.

One interpretation is:

Each **row in orders** gets **extended** with (the corresponding information from the **customer** table) **name** (each row in orders matches exactly one row in customer).

SELECT

FROM **orders**

JOIN **customers**

ON **orders.customer\_id = customers.customer\_id**

So my logic:

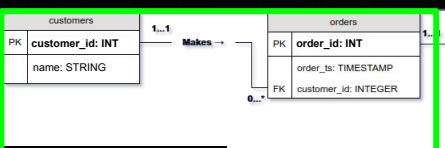
- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

We have a one-to-many join.

This is the **most common** by far.

One interpretation is:

Each **row in orders** gets **extended** with (the corresponding information from the **customer** table) **name** (each row in orders matches exactly one row in customer).

SELECT

FROM **orders**

JOIN **customers**

```

ON orders.customer_id =
customers.customer_id
  
```

So my logic:

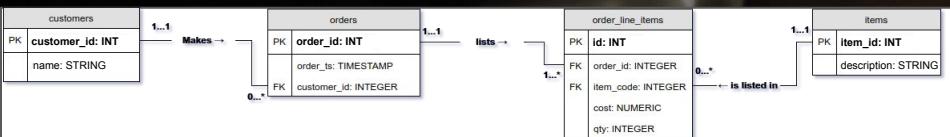
- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

## Let's break this logic down!

```

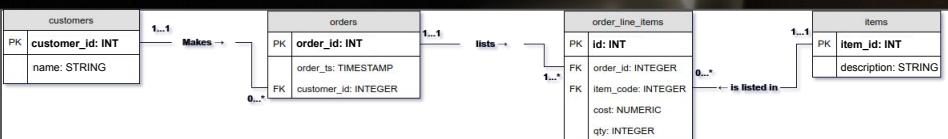
SELECT customer_id, name
FROM orders
JOIN customers
ON orders.customer_id =
customers.customer_id
GROUP BY customer_id, name
    
```

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

## Let's break this logic down!

```

SELECT customer_id, name,
       COUNT(DATE_TRUNC('month', order_ts))
FROM orders
JOIN customers
ON orders.customer_id = customers.customer_id
GROUP BY customer_id, name
    
```

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( <code>orders.order_id</code> ), <code>NOT NULL</code>		The <code>order_id</code> for the order which this line item belongs	123456
<code>item_qty</code>	INTEGER				545
<code>cost_qty</code>	DECIMAL				16.50
<code>item_cost</code>	DECIMAL				1

Ok. So this line got a little more complex.

Let's break it down.

- GROUP BY = buckets, one per customer
- Running the SELECT processes row-by-row (so bucket-by-bucket)

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

### So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute:  
number of months visited
- 3) For each customer compute:  
number of months active
- 4) Then filter:  $(2) / (3) > 1$

Let's break this logic down!

```
SELECT customer_id, name,  
       COUNT(DATE_TRUNC('month', order_ts)),  
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) )  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name
```

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( <code>orders.order_id</code> ), <code>NOT NULL</code>		The <code>order_id</code> for the order which this line item belongs	123456
<code>item_qty</code>	DECIMAL				545
<code>cost_qty</code>	DECIMAL				16.50
<code>customer_id</code>	INTEGER				1

Ok. So this line got a little more complex.

Let's break it down.

- GROUP BY = buckets, one per customer
- Running the SELECT processes row-by-row (so bucket-by-bucket)
- For customer 1:
  - Look in bucket, compute: `MAX(order_ts)`
  - Look in bucket, compute: `MIN(order_ts)`

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

### So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter:  $(2) / (3) > 1$

### Let's break this logic down!

```
SELECT customer_id, name,  
       COUNT(DATE_TRUNC('month', order_ts)),  
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) )  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name
```

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( <code>orders.order_id</code> ), <code>NOT NULL</code>		The <code>order_id</code> for the order which this line item belongs	123456 545 16.50 1

Ok. So this line got a little more complex.

Let's break it down.

- GROUP BY = buckets, one per customer
- Running the SELECT processes row-by-row (so bucket-by-bucket)
- For customer 1:
  - Look in bucket, compute: `MAX(order_ts)`
  - Look in bucket, compute: `MIN(order_ts)`
  - Take these two results (now TIMESTAMPS) and pass it to the `MONTHS_BETWEEN(..)` function to get the number of months between these - the number of months the customer has been active

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

### So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter:  $(2) / (3) > 1$

```
SELECT customer_id, name,  
       COUNT(DATE_TRUNC('month', order_ts)),  
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) )  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name
```

### Let's break this logic down!

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( <code>orders.order_id</code> ), <code>NOT NULL</code>		The <code>order_id</code> for the order which this line item belongs	123456
<code>item_qty</code>	INTEGER				545
<code>cost_qty</code>	DECIMAL				16.50
<code>customer_id</code>	INTEGER				1

Ok. So this line got a little more complex.

Let's break it down.

- GROUP BY = buckets, one per customer
- Running the SELECT processes row-by-row (so bucket-by-bucket)
- For customer 1:
  - Look in bucket, compute: `MAX(order_ts)`
  - Look in bucket, compute: `MIN(order_ts)`
  - Take these two results (now TIMESTAMPS) and pass it to the `MONTHS_BETWEEN(..)` function to get the number of months between these - the number of months the customer has been active
  - The customer doesn't have to have come in exactly two months apart to be considered shopped in two months, so we round this difference up.
  - This is done by a new function `CEIL(..)`

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

### So my logic:

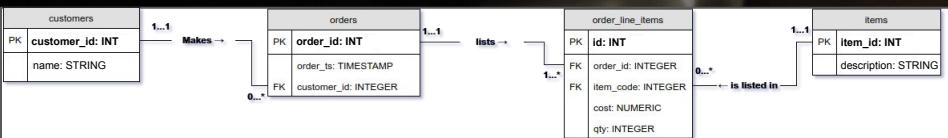
- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter:  $(2) / (3) > 1$

```
SELECT customer_id, name,  
       COUNT(DATE_TRUNC('month', order_ts)),  
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) )  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name
```

### Let's break this logic down!

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id ), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	STRING	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	STRING	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	INTEGER	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

## So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

```

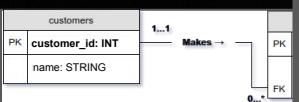
SELECT customer_id, name,
       COUNT(DATE_TRUNC('month', order_ts)) /
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) )
           AS shops_per_month
FROM orders
JOIN customers
ON orders.customer_id = customers.customer_id
GROUP BY customer_id, name
HAVING COUNT(DATE_TRUNC('month', order_ts)) /
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
  
```

## Let's break this logic down!

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
id	INTEGER	Primary Key		each order	9654000
order_id	INTEGER	Foreign Key	NOT NULL	this line item	123456
item_code	INTEGER	Foreign Key	NOT NULL	item.	545
cost	NUMERIC	NOT NULL		line item.	16.50
qty	INTEGER	CHECK(		line item.	1

And we're done!



Field Name	Data Type	Constraint
item_id	SERIAL	Primary Key
description	STRING	NOT NULL

## customers

```
SELECT 1/x, avg(y)
FROM tab
GROUP BY x
HAVING NOT x = 0;
```

Field Name	Data Type	Constraint
order_id	INTEGER	Primary Key
order_ts	TIMESTAMP	NOT NULL
customer_id	INTEGER	NOT NULL

### Side note:

Why can't I just refer to **shops\_per\_month** in the HAVING?

As the HAVING part is computed first.

Why? Consider:

```
SELECT 1/x, avg(y)
FROM tab
GROUP BY x
HAVING NOT x = 0;
```

Computing the SELECT fields before the HAVING would result in a divide by zero error.

customer_id	INT	customer name

customer_id	INT	customer name

customer_id	INT	customer name

customer_id	INT	customer name

customer_id	INT	customer name

customer_id	INT	customer name

## QUESTION

List the customers by name who shop with us more than once a month while they were active customers?

### Let's break this logic down!

#### So my logic:

- 1) Each row: order + customer name
- 2) For each customer compute: number of months visited
- 3) For each customer compute: number of months active
- 4) Then filter: (2) / (3) > 1

```
SELECT customer_id, name,
       COUNT(DATE_TRUNC('month', order_ts)) /
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) )
            AS shops_per_month
  FROM orders
  JOIN customers
    ON orders.customer_id = customers.customer_id
 GROUP BY customer_id, name
 HAVING COUNT(DATE_TRUNC('month', order_ts)) /
       CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
```

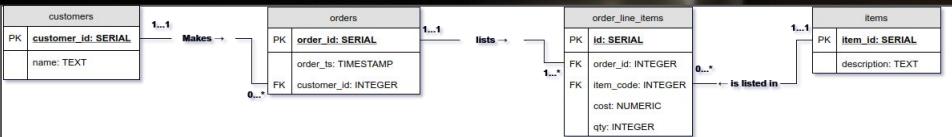
## **One down, one to go!**

- 1) *List the customers by name who shop with us more than once a month while they were active customers?*
  
- 2) *How many customers shop with us more than once a month?*

```
SELECT customer_id, name
FROM orders
JOIN customers
ON orders.customer_id = customers.customer_id
GROUP BY customer_id, name
HAVING COUNT( DATE_TRUNC('month', order_ts)) /
CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
```

## order\_line\_item

Field Name	Data Type	Constraint	Default	Description	Example
<code>id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each order	9654000
<code>order_id</code>	INTEGER	Foreign Key ( orders.order_id), NOT NULL		The order_id for the order which this line item belongs.	123456
<code>item_code</code>	INTEGER	Foreign Key (items.item_id), NOT NULL		The id for the item bought for this line item.	545
<code>cost</code>	NUMERIC	NOT NULL, CHECK( cost > 0 )		The total cost (unit cost x qty) of the line item.	16.50
<code>qty</code>	INTEGER	CHECK( qty > 0 )		The number of items brought for this line item.	1



## items

Field Name	Data Type	Constraint	Default	Description	Example
<code>item_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each item	545
<code>description</code>	TEXT	NOT NULL		The description of the item	Mellowship Slinky Decaf

## customers

Field Name	Data Type	Constraint	Default	Description	Example
<code>customer_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each customer	8734
<code>name</code>	TEXT	NOT NULL		The name of the customer	Gavin Smith

## orders

Field Name	Data Type	Constraint	Default	Description	Example
<code>order_id</code>	SERIAL	Primary Key	auto-increment	Unique identifier generated for each order	123456
<code>order_ts</code>	TIMESTAMP	NOT NULL	now()	Timestamp of when the order was placed	2017-01-01 13:01:02
<code>customer_id</code>	INTEGER	NOT NULL		The id of the customer that made the order.	8734

## QUESTION

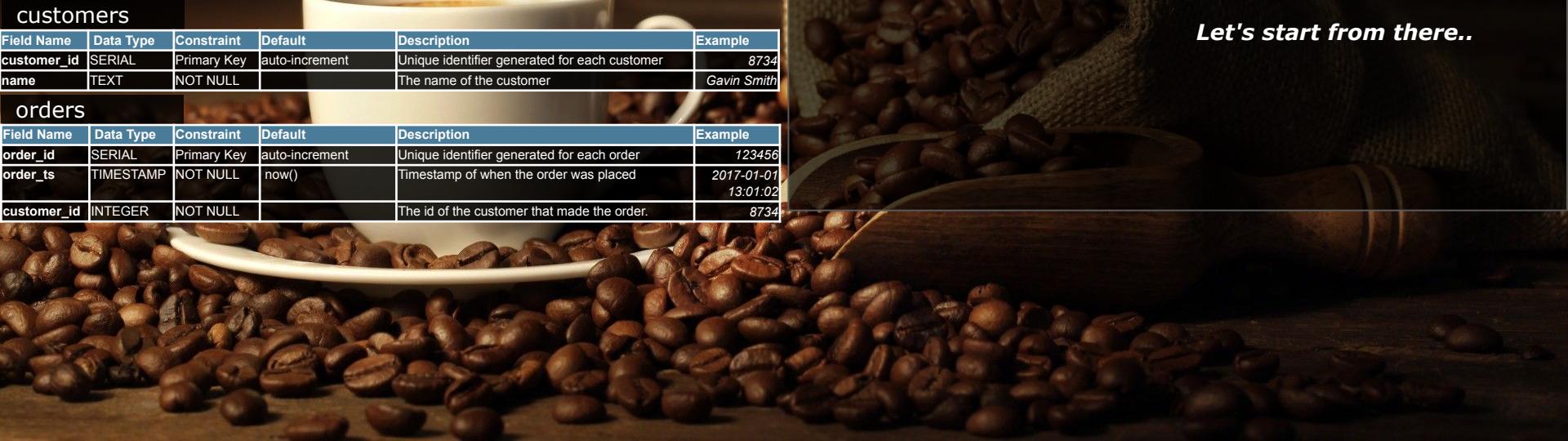
How many customers shop with us more than once a month?

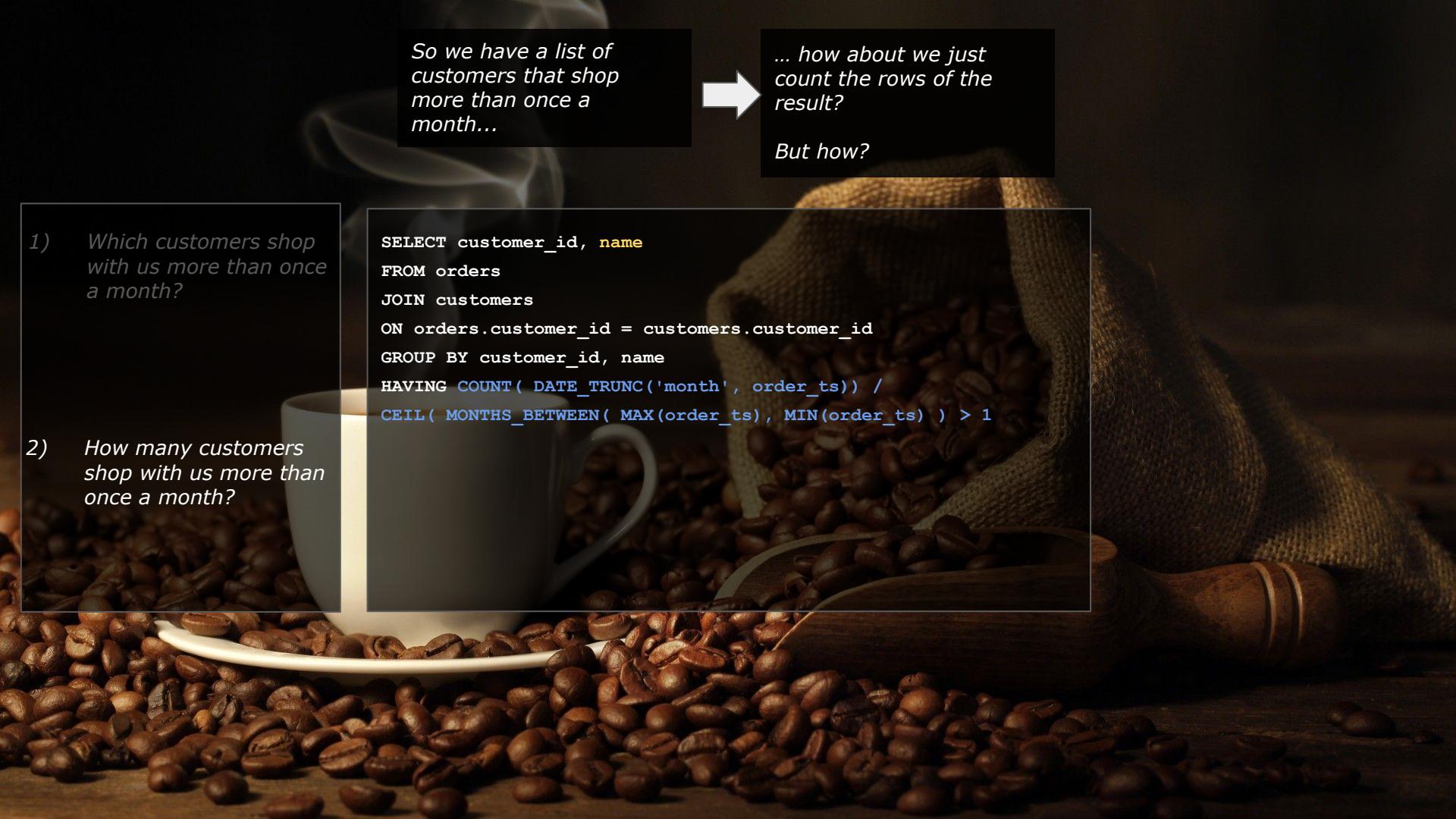
**Wait... this looks very similar to before...**

### Before:

List the customers by name who shop with us more than once a month?

**Let's start from there..**





*So we have a list of customers that shop more than once a month...*

... how about we just count the rows of the result?

*But how?*

1) *Which customers shop with us more than once a month?*

2) *How many customers shop with us more than once a month?*

```
SELECT customer_id, name  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name  
HAVING COUNT( DATE_TRUNC('month', order_ts)) /  
CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
```



So we have a list of customers that shop more than once a month...

... how about we just count the rows of the result?

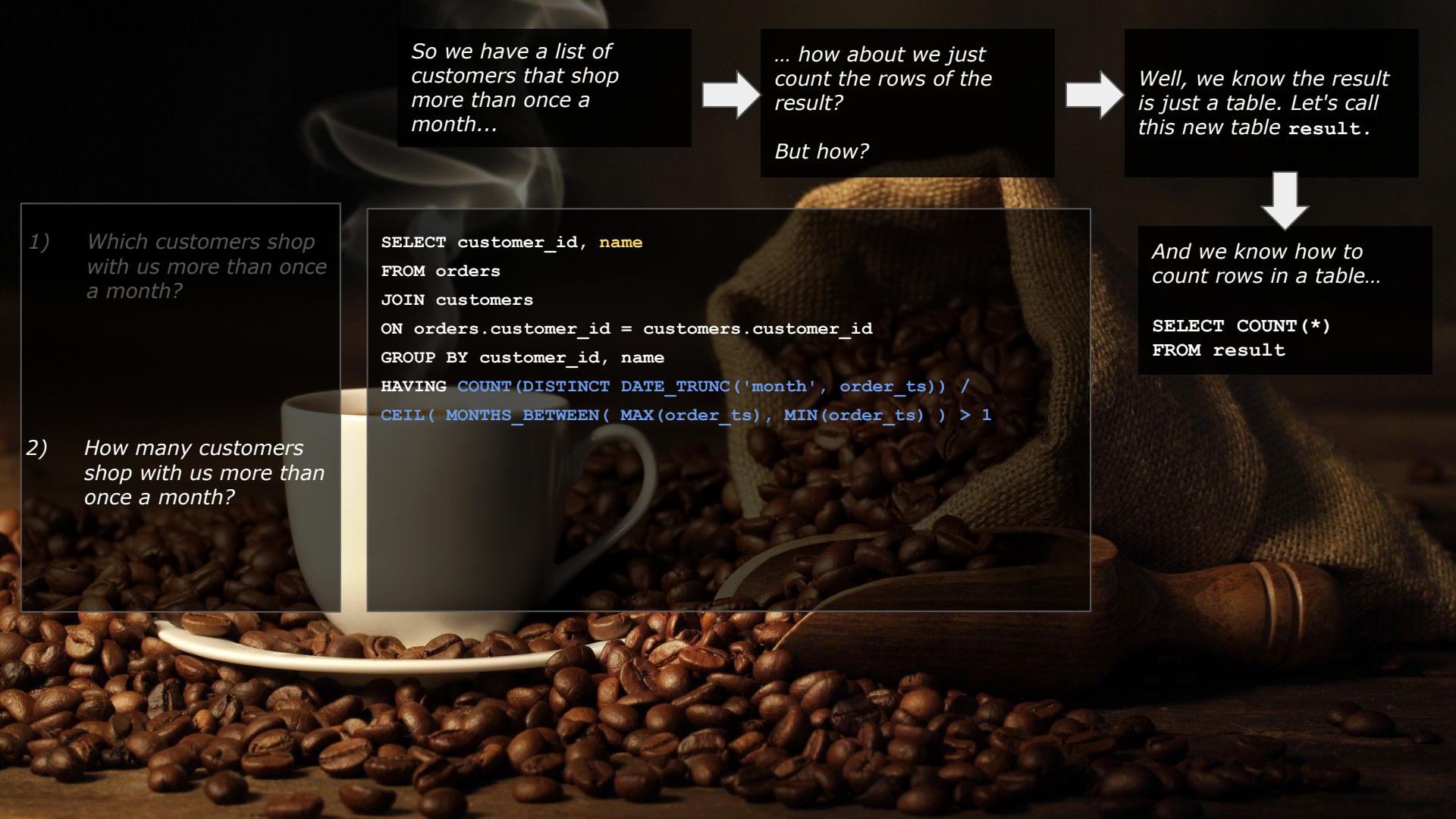
But how?

Well, we know the result is just a table. Let's call this new table **result**.

1) Which customers shop with us more than once a month?

2) How many customers shop with us more than once a month?

```
SELECT customer_id, name  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name  
HAVING COUNT( DATE_TRUNC('month', order_ts)) /  
CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
```

- 
- 1) Which customers shop with us more than once a month?
  - 2) How many customers shop with us more than once a month?

So we have a list of customers that shop more than once a month...

... how about we just count the rows of the result?

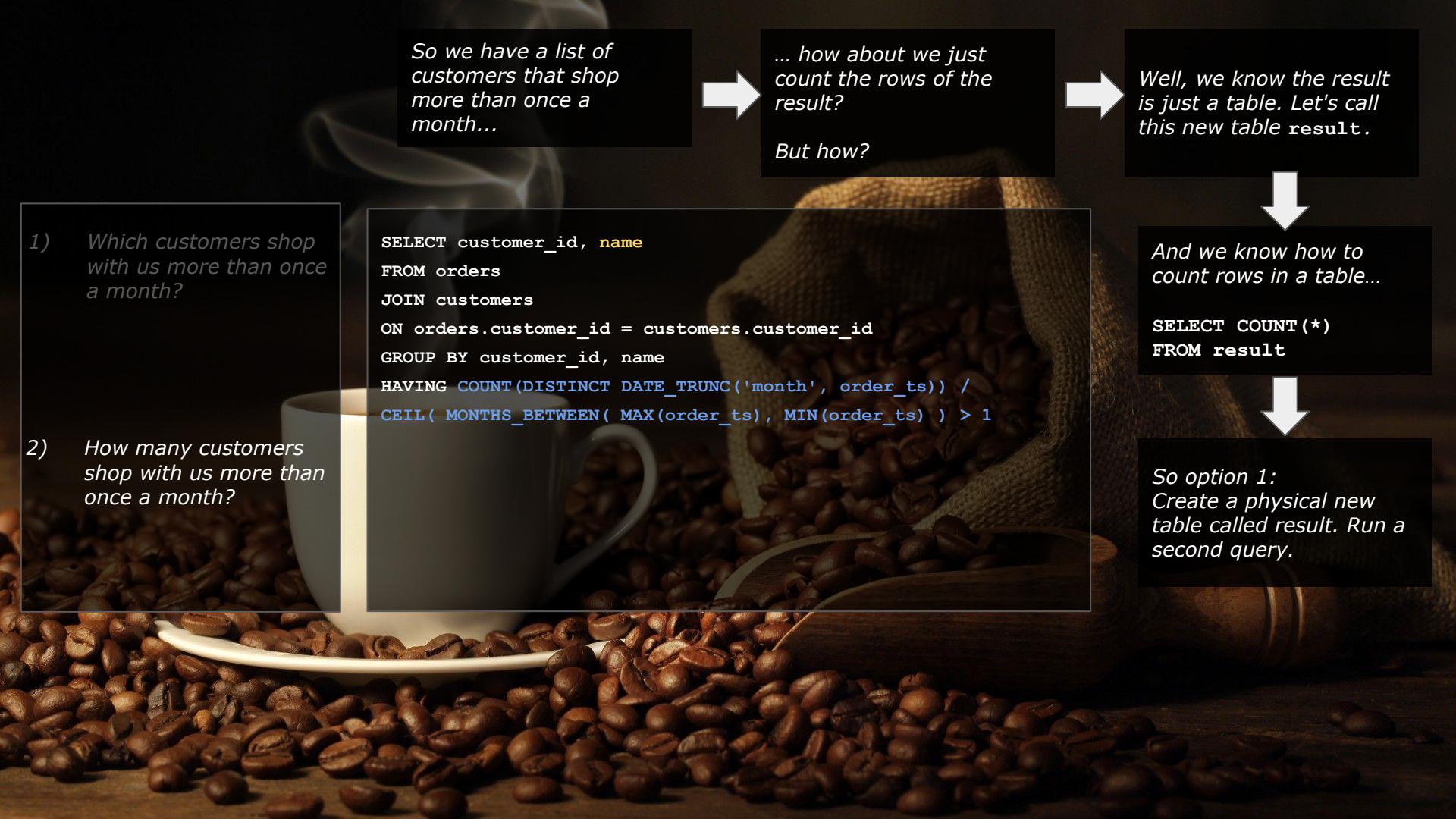
But how?

```
SELECT customer_id, name  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name  
HAVING COUNT(DISTINCT DATE_TRUNC('month', order_ts)) /  
CEIL(MONTHS_BETWEEN(MAX(order_ts), MIN(order_ts))) > 1
```

Well, we know the result is just a table. Let's call this new table **result**.

And we know how to count rows in a table...

```
SELECT COUNT(*)  
FROM result
```

- 
- 1) Which customers shop with us more than once a month?
  - 2) How many customers shop with us more than once a month?

So we have a list of customers that shop more than once a month...

... how about we just count the rows of the result?

But how?

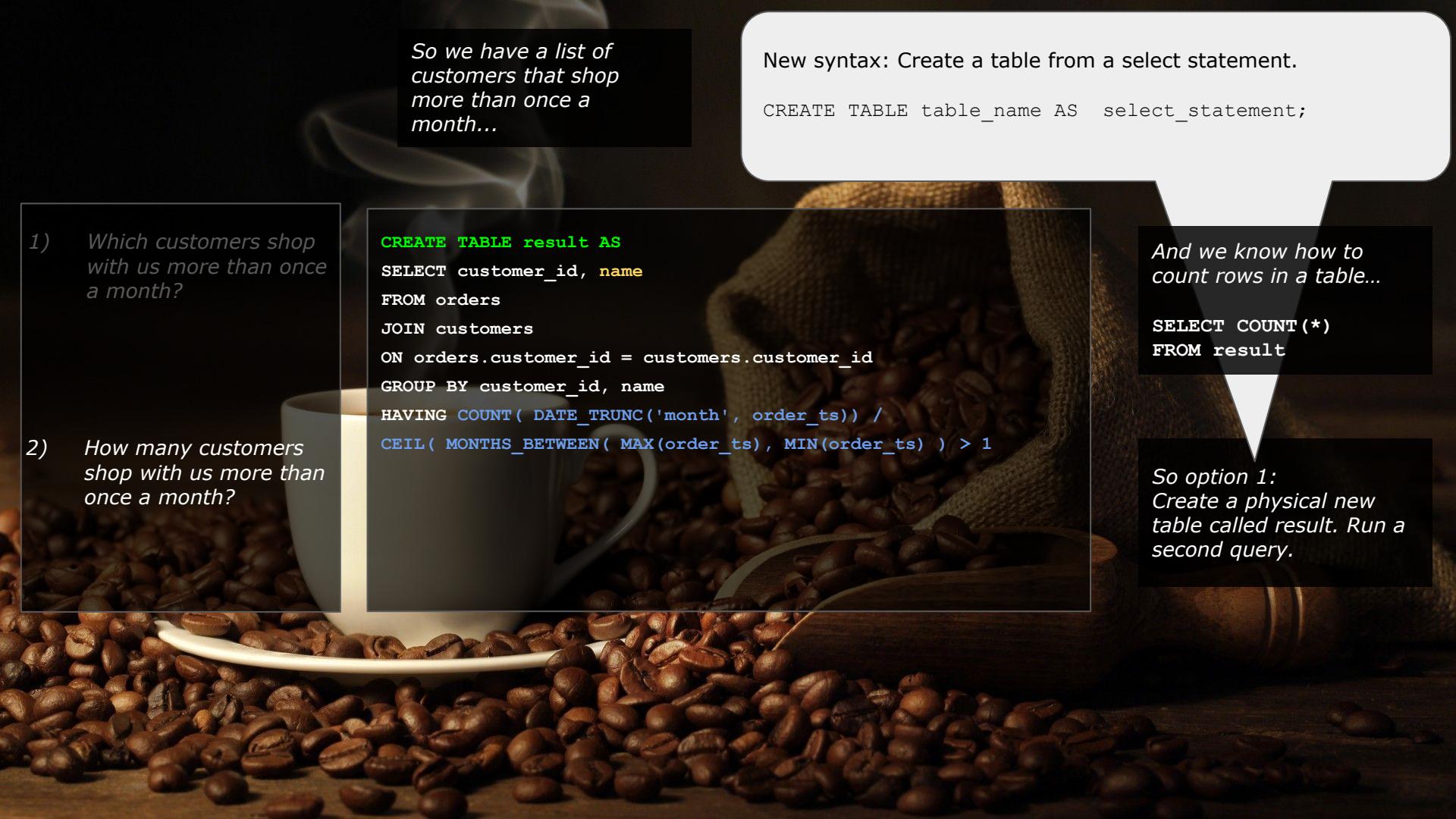
```
SELECT customer_id, name  
FROM orders  
JOIN customers  
ON orders.customer_id = customers.customer_id  
GROUP BY customer_id, name  
HAVING COUNT(DISTINCT DATE_TRUNC('month', order_ts)) /  
CEIL(MONTHS_BETWEEN(MAX(order_ts), MIN(order_ts))) > 1
```

Well, we know the result is just a table. Let's call this new table **result**.

And we know how to count rows in a table...

```
SELECT COUNT(*)  
FROM result
```

So option 1:  
Create a physical new table called **result**. Run a second query.

- 
- 1) Which customers shop with us more than once a month?
  - 2) How many customers shop with us more than once a month?

So we have a list of customers that shop more than once a month...

New syntax: Create a table from a select statement.

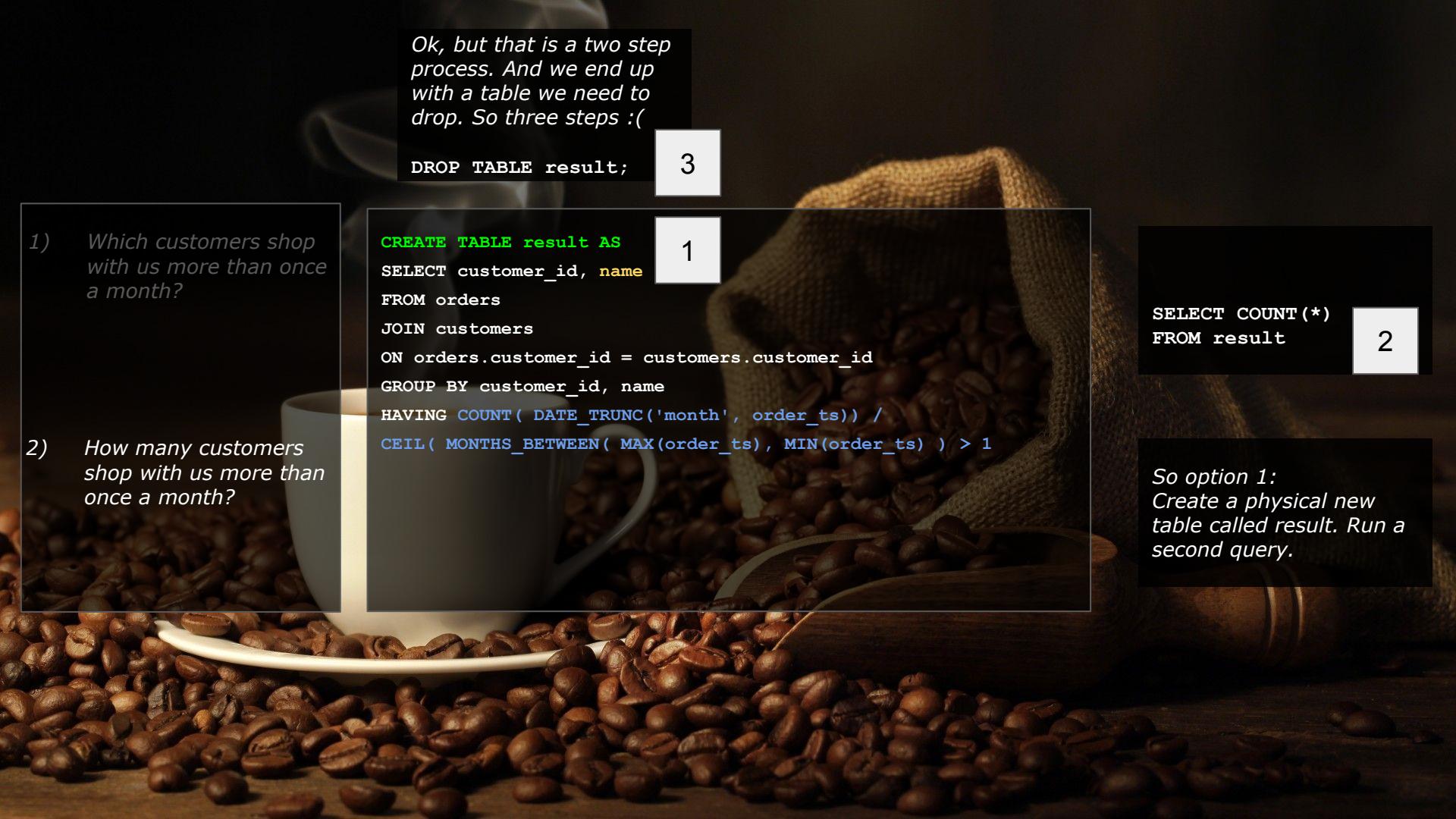
```
CREATE TABLE table_name AS select_statement;
```

```
CREATE TABLE result AS
SELECT customer_id, name
FROM orders
JOIN customers
ON orders.customer_id = customers.customer_id
GROUP BY customer_id, name
HAVING COUNT( DATE_TRUNC('month', order_ts)) /
CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
```

And we know how to count rows in a table...

```
SELECT COUNT(*)
FROM result
```

So option 1:  
Create a physical new table called result. Run a second query.



*Ok, but that is a two step process. And we end up with a table we need to drop. So three steps :(*

DROP TABLE result;

3

- 1) *Which customers shop with us more than once a month?*
- 2) *How many customers shop with us more than once a month?*

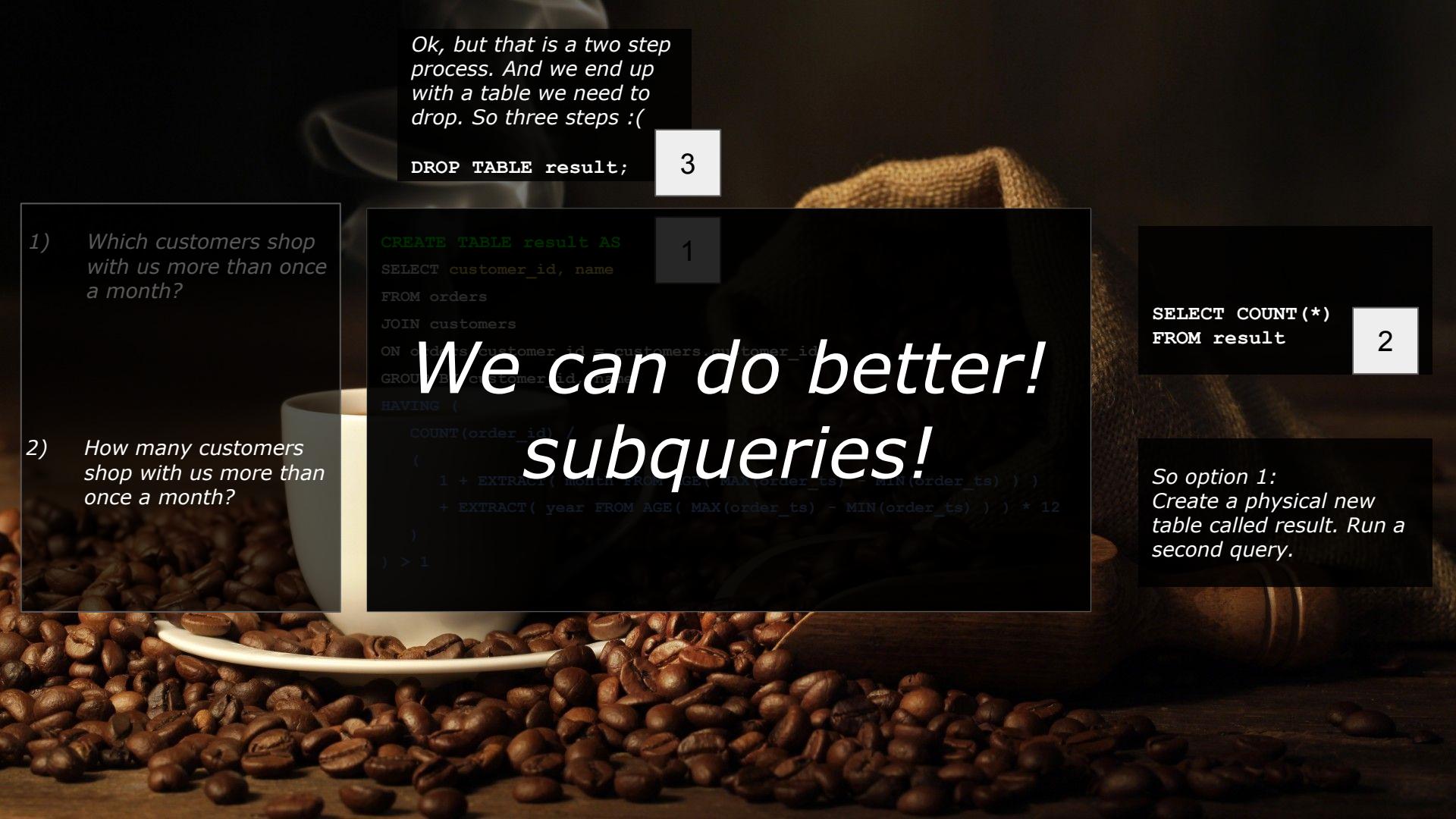
```
CREATE TABLE result AS
SELECT customer_id, name
FROM orders
JOIN customers
ON orders.customer_id = customers.customer_id
GROUP BY customer_id, name
HAVING COUNT( DATE_TRUNC('month', order_ts)) /
CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
```

1

SELECT COUNT(\*)
FROM result

2

*So option 1:  
Create a physical new table called result. Run a second query.*



*Ok, but that is a two step process. And we end up with a table we need to drop. So three steps :(*

DROP TABLE result;

3

- 1) *Which customers shop with us more than once a month?*
- 2) *How many customers shop with us more than once a month?*

```
CREATE TABLE result AS
SELECT customer_id, name
FROM orders
JOIN customers
ON orders.customer_id = customers.customer_id
GROUP BY customer_id
HAVING (
    COUNT(order_id) /
    (
        1 + EXTRACT( month FROM AGE( MAX(order_ts), MIN(order_ts) ) )
        + EXTRACT( year FROM AGE( MAX(order_ts) - MIN(order_ts) ) ) * 12
    )
) > 1
```

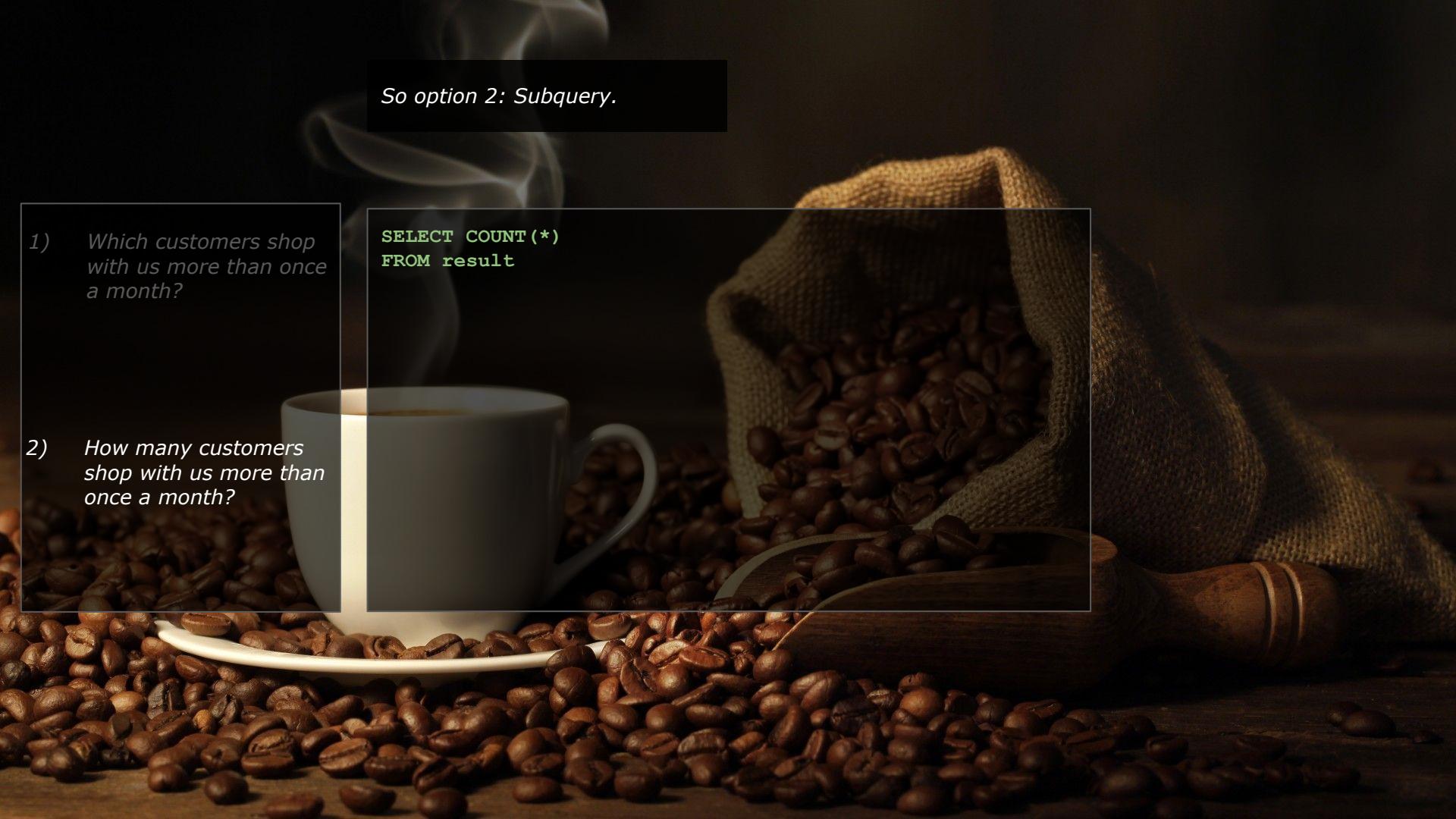
1

SELECT COUNT(\*)
FROM result

2

# We can do better! subqueries!

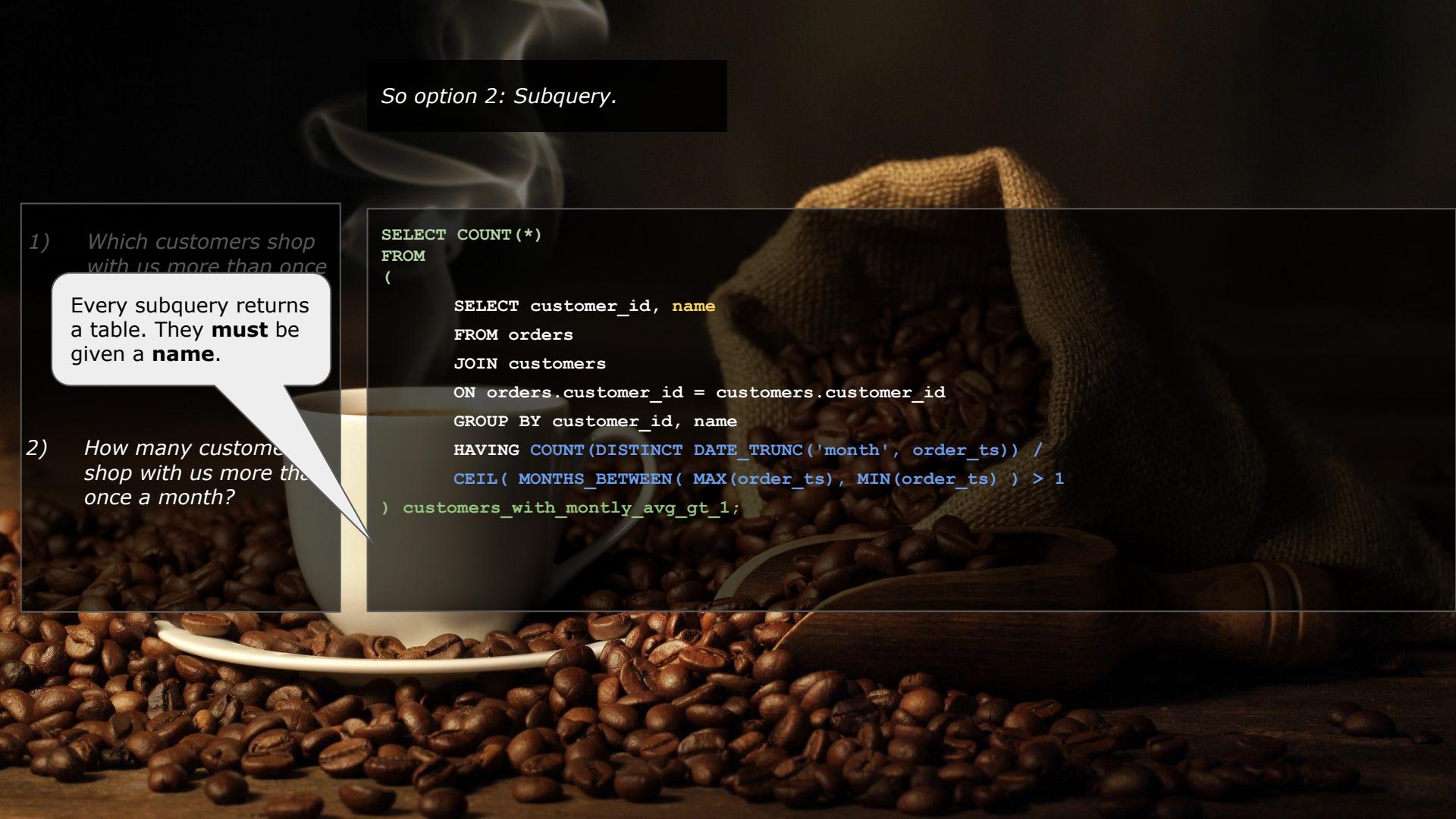
*So option 1:  
Create a physical new  
table called result. Run a  
second query.*



*So option 2: Subquery.*

- 1) *Which customers shop with us more than once a month?*
  
- 2) *How many customers shop with us more than once a month?*

```
SELECT COUNT(*)  
FROM result
```



*So option 2: Subquery.*

- 1) *Which customers shop with us more than once*

Every subquery returns a table. They **must** be given a **name**.

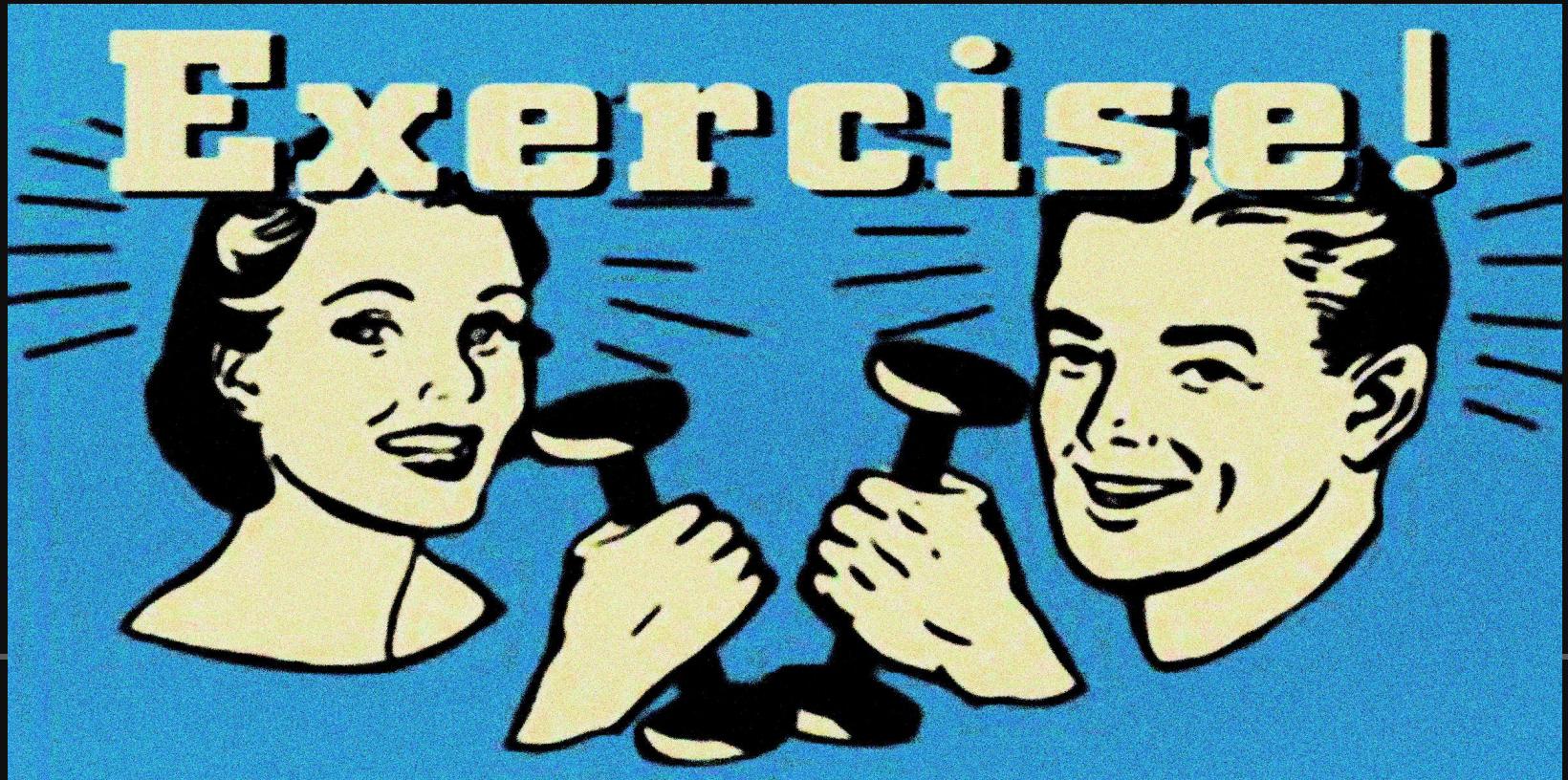
- 2) *How many customers shop with us more than once a month?*

```
SELECT COUNT(*)
FROM
(
    SELECT customer_id, name
    FROM orders
    JOIN customers
    ON orders.customer_id = customers.customer_id
    GROUP BY customer_id, name
    HAVING COUNT(DISTINCT DATE_TRUNC('month', order_ts)) /
        CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
) customers_with_monthly_avg_gt_1;
```

**Well, that was  
easier.  
And we're done!**

- 1) *Which customers shop with us more than once a month?*
  
- 2) *How many customers shop with us more than once a month?*

```
SELECT COUNT(*)
FROM
(
    SELECT customer_id, customers.name
    FROM orders, customers
    WHERE orders.customer_id = customers.customer_id
    GROUP BY customer_id, name
    HAVING COUNT( DATE_TRUNC('month', order_ts)) /
        CEIL( MONTHS_BETWEEN( MAX(order_ts), MIN(order_ts) ) ) > 1
) customers_with_montly_avg_gt_1;
```



NLAB: *Data at Scale*

Dr Georgiana Nica-Avram