

# Lecture 11-1: Advanced Query Techniques

DATA 503: Fundamentals of Data Engineering

Lucas P. Cordova, Ph.D.

2026-03-30

This lecture covers advanced SQL query techniques including subqueries (in WHERE, FROM, and SELECT clauses), subquery expressions (IN and EXISTS), LATERAL joins, Common Table Expressions, cross tabulations with crosstab(), and CASE expressions. We apply these patterns to Census, survey, and temperature data. Based on Chapter 13 of Practical SQL, 2nd Edition.

## Table of contents

<b>1</b>	<b>What You'll Learn Today</b>	<b>2</b>
<b>2</b>	<b>Part 1: Setting Up</b>	<b>2</b>
<b>3</b>	<b>Part 2: Subqueries</b>	<b>3</b>
<b>4</b>	<b>Filtering with Subqueries in WHERE</b>	<b>6</b>
<b>5</b>	<b>Creating Derived Tables with Subqueries</b>	<b>8</b>
<b>6</b>	<b>Generating Columns with Subqueries</b>	<b>10</b>
<b>7</b>	<b>Part 3: Subquery Expressions</b>	<b>12</b>
<b>8</b>	<b>IN and EXISTS</b>	<b>12</b>
<b>9</b>	<b>Part 4: LATERAL Joins</b>	<b>14</b>
<b>10</b>	<b>Using Subqueries with LATERAL</b>	<b>14</b>
<b>11</b>	<b>Part 5: Common Table Expressions (CTEs)</b>	<b>16</b>
<b>12</b>	<b>CTEs with WITH</b>	<b>17</b>

<b>13 Part 6: Cross Tabulations</b>	<b>21</b>
<b>14 crosstab()</b>	<b>21</b>
<b>15 Part 7: CASE Expressions</b>	<b>24</b>
<b>16 Reclassifying Values with CASE</b>	<b>24</b>
<b>17 Part 8: Summary</b>	<b>28</b>

## 1 What You'll Learn Today

### 1.1 Learning Objectives

By the end of this lecture, you will be able to:

1. Write subqueries in WHERE, FROM, and SELECT clauses to filter, preprocess, and compute values
2. Distinguish between correlated and uncorrelated subqueries and explain their performance implications
3. Use IN and EXISTS subquery expressions to test membership across tables
4. Apply LATERAL joins for per-row subquery evaluation
5. Construct Common Table Expressions (CTEs) for readable, maintainable queries
6. Generate cross tabulations with `crosstab()` and reclassify data using CASE

...

**Course Connection:** These objectives support our course goals of writing complex analytical SQL and building data transformation pipelines that are both correct and maintainable.

## 2 Part 1: Setting Up

### 2.1 Database Setup

#### 2.1.1 Create and Populate the Database

Create a fresh database for this lecture:

```
1 CREATE DATABASE advanced_queries;
```

Download the CSV files from [advanced\\_queries\\_csvs.zip](#) ↓ and unzip them into the /tmp (macOS) or C:\Users\Public (Windows) directory.

Connect to advanced\_queries database, then run the setup script provided: [advanced\\_queries.sql](#) ↓

## 2.2 Data Overview and Expected Counts

The script creates and loads the following tables:

Table	Rows	Description
us_counties_pop_est_2019	3,142	County population estimates
cbp_naics_72_establishments	3,074	Food/accommodation businesses per county
employees	6	Small HR demo table
teachers	6	Small demo table (from Ch 7)
ice_cream_survey	200	Office ice cream flavor preferences
temperature_readings	1,077	Daily temps at three stations

## 2.3 Verify Your Setup

```
1 SELECT count(*) FROM us_counties_pop_est_2019;      -- 3,142
2 SELECT count(*) FROM cbp_naics_72_establishments;  -- 3,074
3 SELECT count(*) FROM employees;                    -- 6
4 SELECT count(*) FROM ice_cream_survey;             -- 200
5 SELECT count(*) FROM temperature_readings;         -- 1,077
```

...

If any count is wrong, re-run the setup. We need all of these today.

## 3 Part 2: Subqueries

...

A query inside a query. It's queries all the way to town.

...

## 3.1 What Is a Subquery?

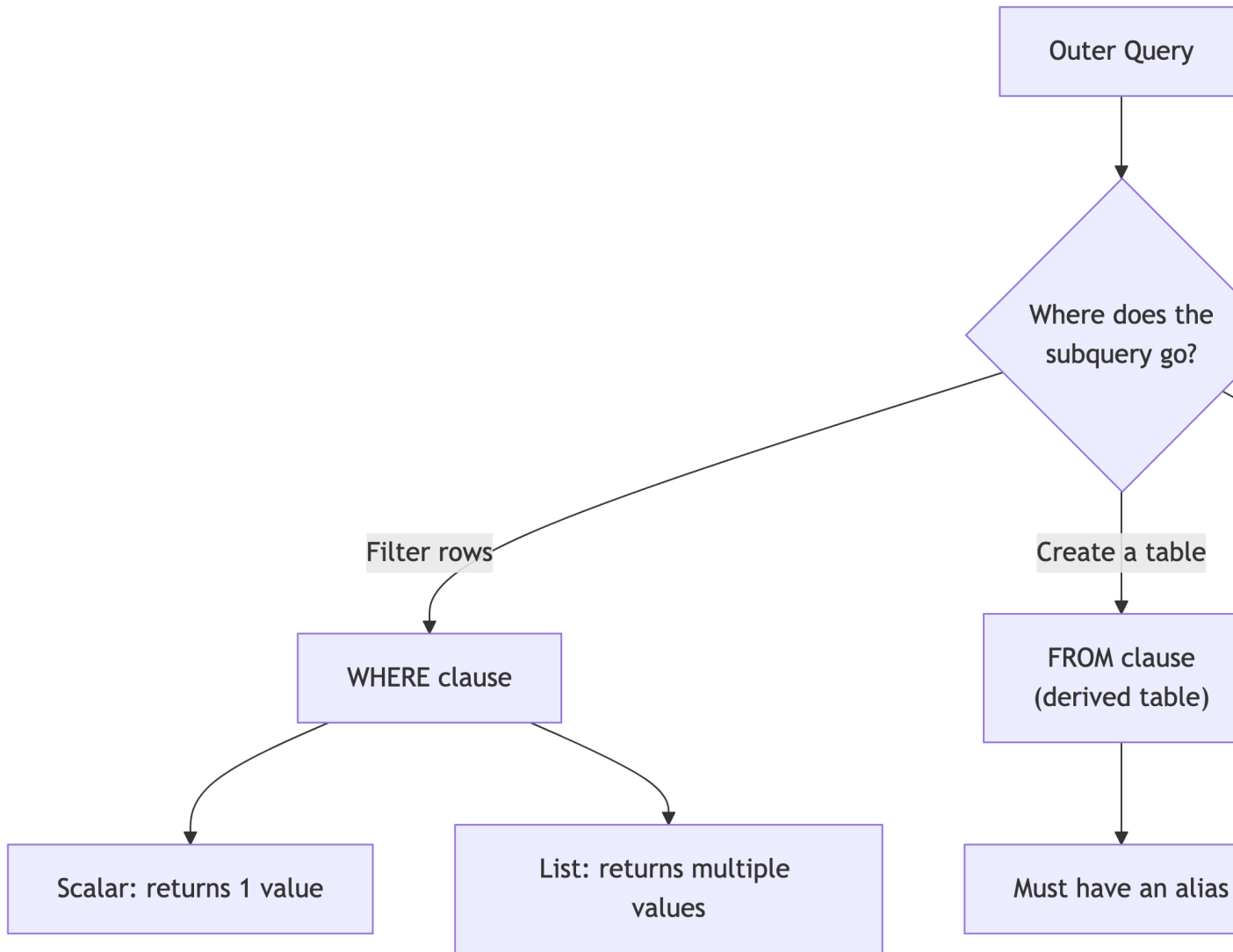
### 3.1.1 The Big Picture

A **subquery** is a query nested inside another query, enclosed in parentheses.

Subqueries can:

- Generate values for filtering (WHERE)
- Create temporary tables (FROM)
- Compute columns (SELECT)

### 3.2 Subquery Process



### 3.3 Subquery Types

There are two types of subqueries:

1. **Correlated:** references the outer query, runs once per row
2. **Uncorrelated:** runs once, independent of outer query

...

Subqueries can also be scalar or non-scalar.

- **Scalar:** returns a single value
- **Non-scalar:** returns multiple values (list)

## 4 Filtering with Subqueries in WHERE

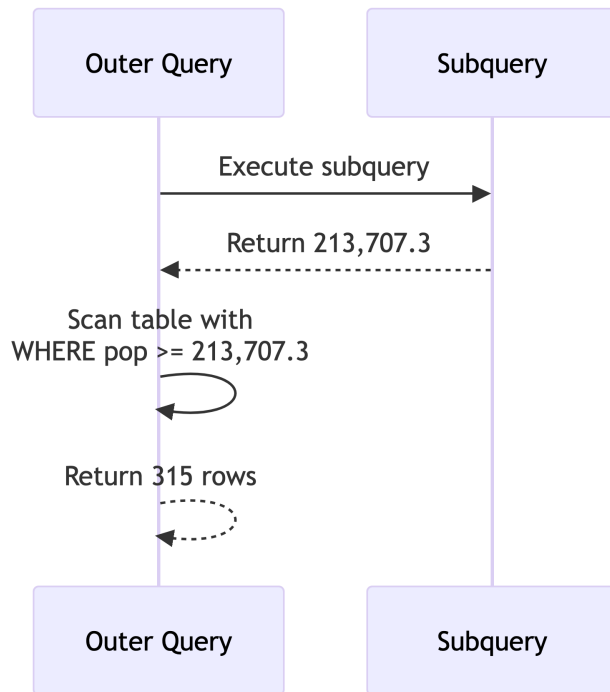
### 4.1 Scalar Subquery in WHERE

The most common pattern: use a subquery to compute a threshold, then filter on it. **Run this.** Inner query first, outer query second. What is this telling us

```
1 SELECT county_name,
2     state_name,
3     pop_est_2019
4 FROM us_counties_pop_est_2019
5 WHERE pop_est_2019 >= (
6     SELECT percentile_cont(.9) WITHIN GROUP (ORDER BY pop_est_2019)
7     FROM us_counties_pop_est_2019
8 )
9 ORDER BY pop_est_2019 DESC;
```

- The inner query computes the 90th percentile of county populations. The outer query returns only counties at or above that threshold.
- The subquery executes first, returns a single value (213,707.3), and the outer query uses it like a constant. This is a **scalar subquery** and it's **uncorrelated**: it runs exactly once.

## 4.2 How It Executes



Key points:

- The subquery runs **once** before the outer query starts scanning
- The result (213,707.3) acts like a constant
- 315 rows returned (about 10% of 3,142)

**Note:** Using `percentile_cont()` in a subquery only works with a single input. Passing an array would return an array, and the `>=` comparison would fail.

## 4.3 Subquery with DELETE

Subqueries work with DML too. Let's create a copy and prune it **Run this.:**

```
1 CREATE TABLE us_counties_2019_top10 AS
2 SELECT * FROM us_counties_pop_est_2019;
3
4 DELETE FROM us_counties_2019_top10
5 WHERE pop_est_2019 < (
6     SELECT percentile_cont(.9) WITHIN GROUP (ORDER BY pop_est_2019)
7     FROM us_counties_2019_top10
```

```

8 );
9
10 SELECT count(*) FROM us_counties_2019_top10;

```

- You started with 3,142 counties and now have ~315 (the top 10%).
- In a pipeline context: this is how you build filtered summary tables. You copy the source, delete what you don't need, keep the rest for downstream consumers.

## 4.4 Quick Check: Subquery Execution

When does an uncorrelated subquery in a WHERE clause execute?

- **A.** Once before the outer query starts
- **B.** Once per row of the outer query
- **C.** After the outer query finishes
- **D.** It depends on the table size

...

**A.** The database computes the inner value once, then scans the outer table. Correlated subqueries (coming up soon) are different: they run once *per row*. That distinction matters for performance.

## 5 Creating Derived Tables with Subqueries

### 5.1 Subquery as a Table in FROM

A subquery in the FROM clause creates a **derived table**: a temporary result set you can query like a regular table (**run this**):

```

1 SELECT round(calcs.average, 0) AS average,
2         calcs.median,
3         round(calcs.average - calcs.median, 0) AS median_average_diff
4 FROM (
5     SELECT avg(pop_est_2019) AS average,
6            percentile_cont(.5)
7            WITHIN GROUP (ORDER BY pop_est_2019)::numeric AS median
8     FROM us_counties_pop_est_2019
9 ) AS calcs;

```

- The inner query computes two aggregates. The outer query calculates the difference. The AS calcs alias is **required**: PostgreSQL demands that derived tables have names.

- The average is ~104,468 and the median is ~25,726. That 78,742 gap tells us some massive counties are inflating the average. This is why we always check both.

## 5.2 Joining Two Derived Tables

This is powerful for combining aggregations from different tables **Run this**. What are your observations on what this is telling us?

```

1  SELECT census.state_name AS st,
2         census.pop_est_2018,
3         est.establishment_count,
4         round((est.establishment_count / census.pop_est_2018::numeric) * 1000, 1)
5         AS estabs_per_thousand
6  FROM
7      (SELECT st, sum(establishments) AS establishment_count
8       FROM cbp_naics_72_establishments
9       GROUP BY st) AS est
10 JOIN
11     (SELECT state_name, sum(pop_est_2018) AS pop_est_2018
12      FROM us_counties_pop_est_2019
13      GROUP BY state_name) AS census
14 ON est.st = census.state_name
15 ORDER BY estabs_per_thousand DESC;
```

## 5.3 Key Points on Previous Query

- Each derived table aggregates to the state level separately, then they join on state name.
- Tourism hotspots like DC, Montana, and Vermont top the list.
- Mississippi and Kentucky are at the bottom.

## 5.4 Quick Check: Derived Table Rules

What happens if you forget the AS alias on a derived table?

- **A.** PostgreSQL infers a name automatically
- **B.** You get a syntax error
- **C.** It works but you can't reference the columns
- **D.** The query runs but returns wrong results

...

**B.** PostgreSQL requires every derived table to have an alias. `AS calcs`, `AS est`, `AS census`: pick a name. Without it, the parser rejects the query before it even tries to run.

## 6 Generating Columns with Subqueries

### 6.1 Column-Level Subqueries

You can put a subquery directly in the `SELECT` list to add a computed column:

```
1 SELECT county_name,  
2         state_name AS st,  
3         pop_est_2019,  
4         (SELECT percentile_cont(.5) WITHIN GROUP (ORDER BY pop_est_2019)  
5          FROM us_counties_pop_est_2019) AS us_median  
6 FROM us_counties_pop_est_2019;
```

...

Every row gets the same median value (25,726) alongside its own population. This lets you compare each county to the national benchmark without a `JOIN`.

...

On its own, that repeating value isn't thrilling. Let's make it useful.

### 6.2 Subquery in a Calculation

Take it further: compute the difference from the median inline:

```
1 SELECT county_name,  
2         state_name AS st,  
3         pop_est_2019,  
4         pop_est_2019 - (SELECT percentile_cont(.5) WITHIN GROUP (ORDER BY pop_est_2019)  
5                          FROM us_counties_pop_est_2019) AS diff_from_median  
6 FROM us_counties_pop_est_2019  
7 WHERE (pop_est_2019 - (SELECT percentile_cont(.5) WITHIN GROUP (ORDER BY pop_est_2019)  
8                          FROM us_counties_pop_est_2019))  
9        BETWEEN -1000 AND 1000;
```

...

This finds ~78 counties whose population is within 1,000 of the national median.

...

Notice the subquery appears **twice**: once in SELECT, once in WHERE. That redundancy is a code smell.

...

CTEs will fix it later. Hang tight.

### 6.3 Try It: Above-Average Counties

Write a query that returns all counties where `pop_est_2019` is above the national average. Show county name, state, and population. Order by population descending.

How many counties are above average? Hint: it's less than half. Why?

```
1 -- Your query here
...
1 SELECT county_name, state_name, pop_est_2019
2 FROM us_counties_pop_est_2019
3 WHERE pop_est_2019 > (
4     SELECT avg(pop_est_2019) FROM us_counties_pop_est_2019
5 )
6 ORDER BY pop_est_2019 DESC;
```

Far fewer than half! The mean is pulled up by massive counties (LA, Cook, Harris), so most counties fall below average.

### 6.4 Key Points on Previous Query

**Median vs. mean:** it matters every time.

## 7 Part 3: Subquery Expressions

...

Checking membership across tables. The SQL equivalent of “are you on the VIP guest list?”

...

## 8 IN and EXISTS

### 8.1 Setup: Employees and Retirees

```
1 CREATE TABLE retirees (  
2     id int,  
3     first_name text,  
4     last_name text  
5 );  
6  
7 INSERT INTO retirees  
8 VALUES (2, 'Janet', 'King'),  
9         (4, 'Michael', 'Taylor');
```

Now we have two employees who retired. Let’s find them, and find who’s still active.

### 8.2 Generating Values for IN

```
1 SELECT first_name, last_name  
2 FROM employees  
3 WHERE emp_id IN (  
4     SELECT id FROM retirees  
5 )  
6 ORDER BY emp_id;
```

IN checks if `emp_id` matches any value in the subquery result. Simple, readable, uncorrelated.

...

**Warning:** Avoid NOT IN with subqueries. If the subquery returns any NULL values, NOT IN returns zero rows. Always. PostgreSQL’s own wiki recommends using NOT EXISTS instead.

### 8.3 Using EXISTS (Correlated)

```
1 SELECT first_name, last_name
2 FROM employees
3 WHERE EXISTS (
4     SELECT id
5     FROM retirees
6     WHERE id = employees.emp_id
7 );
```

EXISTS returns TRUE if the subquery finds **at least one row**. The subquery references `employees.emp_id` from the outer query, making it **correlated**: it runs once per row.

...

For small tables, IN and EXISTS perform identically. For large tables, EXISTS can be faster because it **short-circuits**: it stops scanning as soon as it finds one match.

### 8.4 NOT EXISTS: The Anti-Join

```
1 SELECT first_name, last_name
2 FROM employees
3 WHERE NOT EXISTS (
4     SELECT id
5     FROM retirees
6     WHERE id = employees.emp_id
7 );
```

**Run this.** Returns employees who are NOT retirees. This is an **anti-join**: “give me everything in table A that has no match in table B.”

...

O==O Pipeline use case: find new records that haven't been processed yet.

```
1 WHERE NOT EXISTS (SELECT 1 FROM processed WHERE processed.id = raw.id)
```

This pattern shows up constantly in ETL workflows.

## 8.5 Quick Check: IN vs EXISTS

For a table with 10 million rows, which is typically faster for membership checks?

- A. IN (always)
- B. EXISTS (always)
- C. EXISTS (usually, because it short-circuits)
- D. They're always identical

...

C. EXISTS stops at the first match. IN materializes the full subquery result into a list. For large datasets, that difference matters. But the optimizer is smart and sometimes rewrites one to the other. Profile, don't guess.

## 9 Part 4: LATERAL Joins

...

The most powerful subquery pattern you've (probably) never heard of.

...

## 10 Using Subqueries with LATERAL

### 10.1 LATERAL in FROM: Reusing Calculations

LATERAL lets a subquery in the FROM clause reference columns from preceding tables. Think of it as “for each row, compute this” **Run this.**

```
1 SELECT county_name,
2         state_name,
3         pop_est_2018,
4         pop_est_2019,
5         raw_chg,
6         round(pct_chg * 100, 2) AS pct_chg
7 FROM us_counties_pop_est_2019,
8      LATERAL (SELECT pop_est_2019 - pop_est_2018 AS raw_chg) rc,
9      LATERAL (SELECT raw_chg / pop_est_2018::numeric AS pct_chg) pc
10 ORDER BY pct_chg DESC;
```

...

Each LATERAL subquery can reference columns from earlier in the FROM clause, including results from *other* LATERAL subqueries. `pct_chg` uses `raw_chg` from the first LATERAL. Without LATERAL, you'd need nested subqueries or repeat the calculation.

## 10.2 LATERAL with JOIN: Top-N Per Group

The killer feature. First, let's set up the data:

```
1 ALTER TABLE teachers ADD CONSTRAINT id_key PRIMARY KEY (id);
2
3 CREATE TABLE teachers_lab_access (
4     access_id bigint PRIMARY KEY GENERATED ALWAYS AS IDENTITY,
5     access_time timestamp with time zone,
6     lab_name text,
7     teacher_id bigint REFERENCES teachers (id)
8 );
9
10 INSERT INTO teachers_lab_access (access_time, lab_name, teacher_id)
11 VALUES ('2022-11-30 08:59:00-05', 'Science A', 2),
12         ('2022-12-01 08:58:00-05', 'Chemistry B', 2),
13         ('2022-12-21 09:01:00-05', 'Chemistry A', 2),
14         ('2022-12-02 11:01:00-05', 'Science B', 6),
15         ('2022-12-07 10:02:00-05', 'Science A', 6),
16         ('2022-12-17 16:00:00-05', 'Science B', 6);
```

Run this setup before the next slide.

## 10.3 LATERAL JOIN in Action

```
1 SELECT t.first_name, t.last_name, a.access_time, a.lab_name
2 FROM teachers t
3 LEFT JOIN LATERAL (
4     SELECT *
5     FROM teachers_lab_access
6     WHERE teacher_id = t.id
7     ORDER BY access_time DESC
8     LIMIT 2
9 ) a ON true
10 ORDER BY t.id;
```

...

For each teacher, this returns their 2 most recent lab accesses. `LEFT JOIN` ensures teachers with no lab visits still appear (with `NULLs`). The `ON true` is required syntax since the `WHERE` inside handles the correlation.

...

Think of `LATERAL JOIN` like a for-loop: “for each teacher, run this subquery.” In data engineering: “for each customer, get their last 3 orders.” “For each sensor, get the 5 most recent readings.”

## 10.4 Try It: Top-1 Per Teacher

Modify the `LATERAL` query to return only each teacher’s **most recent** lab access (not two). Teachers with no lab access should still appear (show `NULLs`).

```
1 -- Your query here
...
1 SELECT t.first_name, t.last_name, a.access_time, a.lab_name
2 FROM teachers t
3 LEFT JOIN LATERAL (
4     SELECT *
5     FROM teachers_lab_access
6     WHERE teacher_id = t.id
7     ORDER BY access_time DESC
8     LIMIT 1
9 ) a ON true
10 ORDER BY t.id;
```

Change `LIMIT 2` to `LIMIT 1`. That’s it. The `LEFT JOIN` keeps all teachers regardless. Simple, powerful, clean.

## 11 Part 5: Common Table Expressions (CTEs)

...

Named temporary result sets that make complex queries readable. Your future self will thank you.

...

## 12 CTEs with WITH

### 12.1 A Simple CTE

```
1 WITH large_counties (county_name, state_name, pop_est_2019) AS (  
2     SELECT county_name, state_name, pop_est_2019  
3     FROM us_counties_pop_est_2019  
4     WHERE pop_est_2019 >= 100000  
5 )  
6 SELECT state_name, count(*)  
7 FROM large_counties  
8 GROUP BY state_name  
9 ORDER BY count(*) DESC;
```

The CTE filters to counties with 100k+ population. The outer query counts them by state. Texas, Florida, and California top the list.

...

The column list after the CTE name is optional (it inherits from the subquery), but including it makes your intent explicit. Explicit is better than implicit.

### 12.2 CTEs for Joining Aggregations

Remember the derived table join? Here it is rewritten with CTEs:

```
1 WITH  
2     counties (st, pop_est_2018) AS (  
3         SELECT state_name, sum(pop_est_2018)  
4         FROM us_counties_pop_est_2019  
5         GROUP BY state_name  
6     ),  
7     establishments (st, establishment_count) AS (  
8         SELECT st, sum(establishments) AS establishment_count  
9         FROM cbp_naics_72_establishments  
10        GROUP BY st  
11    )  
12 SELECT counties.st,  
13        pop_est_2018,  
14        establishment_count,  
15        round((establishments.establishment_count /  
16            counties.pop_est_2018::numeric(10,1)) * 1000, 1)  
17        AS estabs_per_thousand
```

```

18 FROM counties JOIN establishments
19 ON counties.st = establishments.st
20 ORDER BY estabs_per_thousand DESC;

```

**Compare this to the derived table version.** Same result, dramatically more readable. Each CTE has a name and a clear purpose. In code review, this is the version that gets approved.

## 12.3 CTEs to Eliminate Redundancy

Remember the repeated subquery for median comparison? CTEs fix that:

```

1 WITH us_median AS (
2     SELECT percentile_cont(.5)
3         WITHIN GROUP (ORDER BY pop_est_2019) AS us_median_pop
4     FROM us_counties_pop_est_2019
5 )
6 SELECT county_name,
7     state_name AS st,
8     pop_est_2019,
9     us_median_pop,
10    pop_est_2019 - us_median_pop AS diff_from_median
11 FROM us_counties_pop_est_2019 CROSS JOIN us_median
12 WHERE (pop_est_2019 - us_median_pop) BETWEEN -1000 AND 1000;

```

The median is computed **once** in the CTE, then **CROSS JOIN**ed to every row. No repeated subqueries. Want to switch from median to 90th percentile? Change one number in one place.

## 12.4 When to Choose What (textualized)

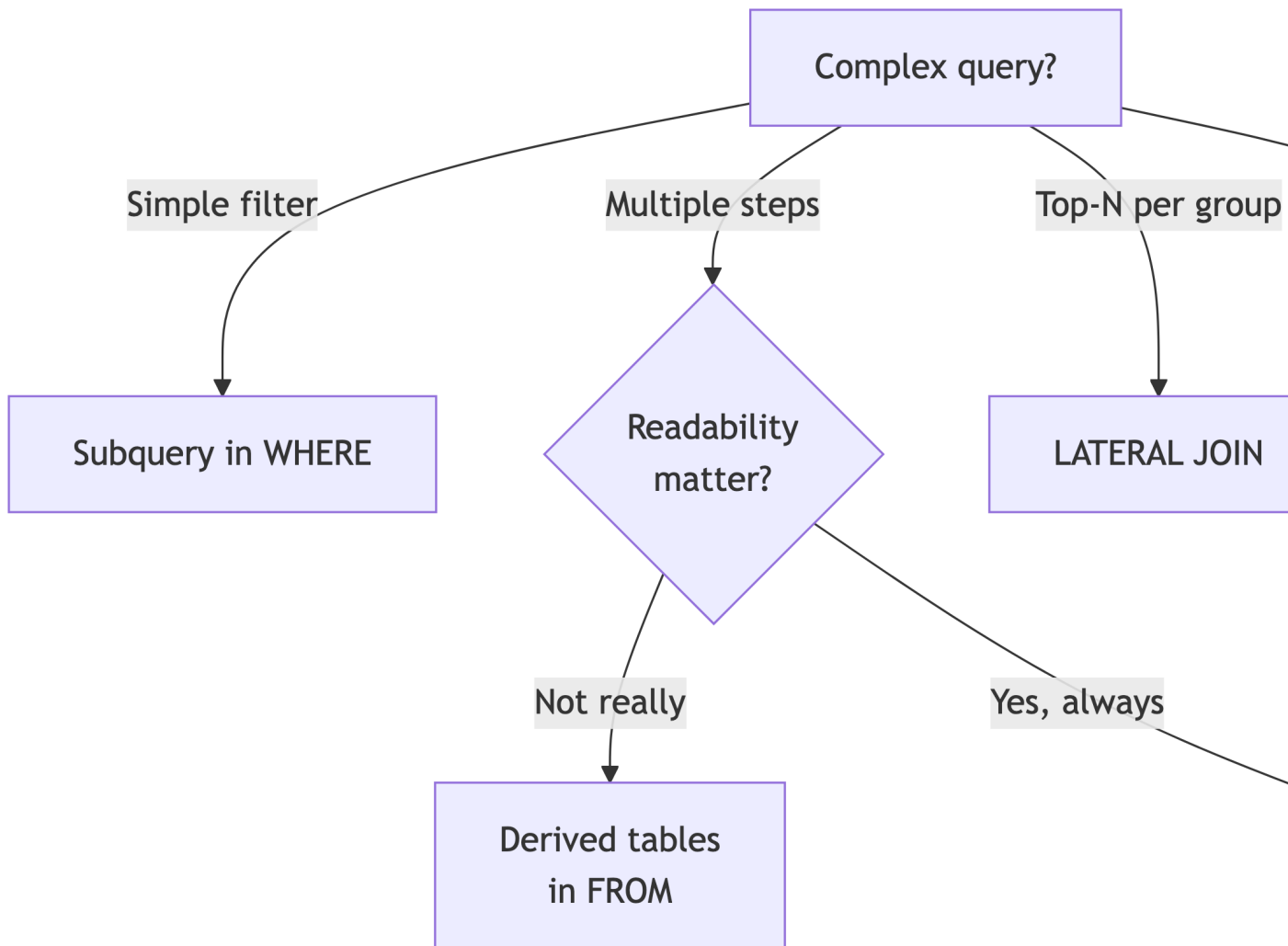
**CTEs win when:**

- You have multiple preprocessing steps
- The same computation is used more than once
- Someone else (or future you) needs to read the query
- You're joining aggregated datasets

**Derived tables win when:**

- It's a one-off simple aggregation
- You don't need to reference it more than once

## 12.5 When to Choose What (visualized)



## 12.6 TL;DR;DS: When to Choose What

My honest take: CTEs > nested subqueries for readability. Always.

## 12.7 Quick Check: CTE Syntax

Which keyword starts a Common Table Expression?

- A. CREATE TEMP TABLE

- B. WITH
- C. AS TABLE
- D. DEFINE

...

B. The WITH keyword defines one or more CTEs. You can chain multiple CTEs separated by commas before the main SELECT. They're sometimes called "WITH queries" for this reason.

## 12.8 Try It: State-Level Summary CTE

Using two CTEs, compute:

1. CTE 1: the total population per state
2. CTE 2: the number of counties per state

Join them and add a column for average population per county. Show state, total pop, county count, and avg pop per county. Order by average descending.

```

1  -- Your query here
...
1  WITH state_pop AS (
2      SELECT state_name, sum(pop_est_2019) AS total_pop
3      FROM us_counties_pop_est_2019
4      GROUP BY state_name
5  ),
6  state_counties AS (
7      SELECT state_name, count(*) AS num_counties
8      FROM us_counties_pop_est_2019
9      GROUP BY state_name
10 )
11 SELECT sp.state_name,
12        sp.total_pop,
13        sc.num_counties,
14        round(sp.total_pop::numeric / sc.num_counties, 0) AS avg_pop_per_county
15 FROM state_pop sp
16 JOIN state_counties sc ON sp.state_name = sc.state_name
17 ORDER BY avg_pop_per_county DESC;
```

Two clean, named steps, then a simple join. This is the CTE pattern at its best.

## 13 Part 6: Cross Tabulations

...

Pivot tables in SQL. Every analyst's favorite party trick.

...

## 14 crosstab()

### 14.1 Enable the Extension

```
1 CREATE EXTENSION tablefunc;
```

This loads the `crosstab()` function from the `tablefunc` module. You only need to do this once per database. Standard ANSI SQL doesn't have a `crosstab` function, but PostgreSQL provides one through its extension system.

### 14.2 The Ice Cream Survey

```
1 SELECT *
2 FROM ice_cream_survey
3 ORDER BY response_id
4 LIMIT 10;
```

200 employees across three offices, each picking a flavor. We want a pivot: offices as rows, flavors as columns, counts as values.

...

The raw data is “long format”: one row per response. We need “wide format”: one row per office with flavor counts as columns. That transformation is what `crosstab()` does.

### 14.3 Generating the Crosstab

```
1 SELECT *
2 FROM crosstab(
3     'SELECT office, flavor, count(*)
4     FROM ice_cream_survey
5     GROUP BY office, flavor
6     ORDER BY office',
7     'SELECT flavor
```

```

8     FROM ice_cream_survey
9     GROUP BY flavor
10    ORDER BY flavor'
11 )
12 AS (office text,
13     chocolate bigint,
14     strawberry bigint,
15     vanilla bigint);

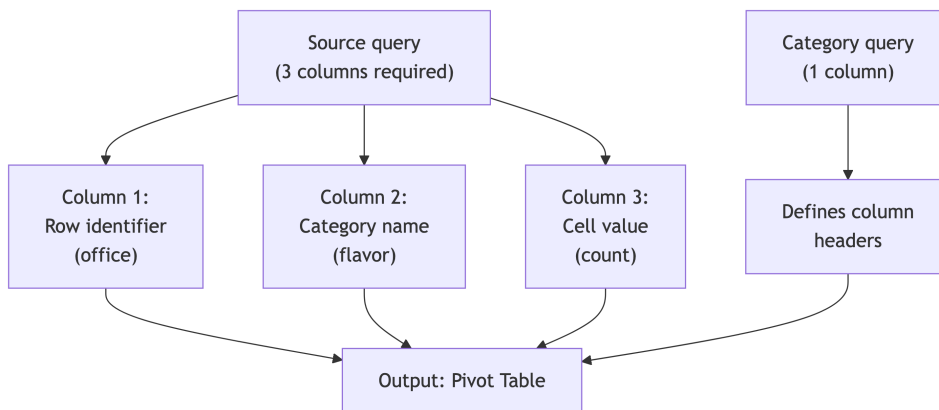
```

**Run this.** The first argument is the source query (must have exactly **3 columns**: row name, category, value). The second argument defines the category labels. The **AS** clause names the output columns.

...

**Takeaway:** Midtown loves chocolate and has zero interest in strawberry (NULL). The Uptown office is more evenly split. Decisions made, ice cream ordered.

## 14.4 How crosstab() Works



### Rules:

- Source query: exactly 3 columns (row, category, value)
- Category query: exactly 1 column (unique category values)
- AS clause: must list output columns in the order the category query returns them
- Both queries are passed as **strings** (single-quoted)
- Escaped quotes inside: use `'` (two single quotes)

## 14.5 Temperature Crosstab

A more complex example: median max temperature by station and month **Run this**. What are your observations?

```
1 SELECT *
2 FROM crosstab(
3     'SELECT station_name,
4         date_part('month', observation_date),
5         percentile_cont(.5)
6         WITHIN GROUP (ORDER BY max_temp)
7     FROM temperature_readings
8     GROUP BY station_name,
9         date_part('month', observation_date)
10    ORDER BY station_name',
11    'SELECT month FROM generate_series(1,12) month'
12 )
13 AS (station text,
14     jan numeric(3,0), feb numeric(3,0), mar numeric(3,0),
15     apr numeric(3,0), may numeric(3,0), jun numeric(3,0),
16     jul numeric(3,0), aug numeric(3,0), sep numeric(3,0),
17     oct numeric(3,0), nov numeric(3,0), dec numeric(3,0));
```

## 14.6 Key Points on Previous Query

---

station	jan	feb	mar	apr	may	jun	jul	aug	sep	oct	nov	dec
CHICAGO NORTHERLY	34	36	46	50	66	77	81	80	77	65	57	35
ISLAND IL US												
SEATTLE BOEING FIELD	50	54	56	64	66	71	76	77	69	62	55	42
WA US												
WAIKIKI 717.2 HI US	83	84	84	86	87	87	88	87	87	86	84	82

---

...

**Takeaway:** Three stations, 12 months, median temperatures. `generate_series(1,12)` creates the month column headers. Waikiki is consistently balmy. Chicago goes from freezing to pleasant. Seattle splits the difference.

## 14.7 Quick Check: Crosstab Requirements

The source query for `crosstab()` must return exactly how many columns?

- A. 2
- B. 3
- C. 4
- D. It depends on the output

...

B. Always three: (1) row identifier, (2) category, (3) value. The row identifier groups the pivot rows, the category determines which output column gets the value. No more, no less.

## 15 Part 7: CASE Expressions

...

Conditional logic inside SQL. Your if/else for data transformation.

...

## 16 Reclassifying Values with CASE

### 16.1 The CASE Pattern

```
1 CASE
2     WHEN condition THEN result
3     WHEN another_condition THEN result
4     ELSE fallback_result
5 END
```

- Evaluates conditions **top-to-bottom**
- Returns the **first match** and stops
- **ELSE** catches anything that fell through
- Without **ELSE**, unmatched rows get **NULL**
- Works in **SELECT**, **WHERE**, **ORDER BY**, and aggregations

...

Think of it as a series of if/else-if/else blocks. Order matters: put your most specific conditions first.

## 16.2 Reclassifying Temperature Data

```
1 SELECT max_temp,
2         CASE WHEN max_temp >= 90 THEN 'Hot'
3              WHEN max_temp >= 70 AND max_temp < 90 THEN 'Warm'
4              WHEN max_temp >= 50 AND max_temp < 70 THEN 'Pleasant'
5              WHEN max_temp >= 33 AND max_temp < 50 THEN 'Cold'
6              WHEN max_temp >= 20 AND max_temp < 33 THEN 'Frigid'
7              WHEN max_temp < 20 THEN 'Inhumane'
8              ELSE 'No reading'
9         END AS temperature_group
10 FROM temperature_readings
11 ORDER BY station_name, observation_date;
```

Each row gets classified into a group. The ranges cover all possible values with no gaps. The ELSE clause is our safety net (and yes, “Inhumane” is an editorial choice, not a meteorological term).

...

May Dr. Gore grant me absolution.

## 16.3 CASE in a Common Table Expression (Steps)

Combine CASE with a CTE to classify, then aggregate. The two-step pattern:

- **Step 1 (CTE):** classify every reading.
- **Step 2 (outer query):** count days per station per group.

## 16.4 Temperatures Collapsed

**Run this.** What are your observations?

```
1 WITH temps_collapsed (station_name, max_temperature_group) AS (
2     SELECT station_name,
3            CASE WHEN max_temp >= 90 THEN 'Hot'
4                 WHEN max_temp >= 70 AND max_temp < 90 THEN 'Warm'
5                 WHEN max_temp >= 50 AND max_temp < 70 THEN 'Pleasant'
6                 WHEN max_temp >= 33 AND max_temp < 50 THEN 'Cold'
7                 WHEN max_temp >= 20 AND max_temp < 33 THEN 'Frigid'
8                 WHEN max_temp < 20 THEN 'Inhumane'
9                 ELSE 'No reading'
10            END
```

```

11     FROM temperature_readings
12 )
13 SELECT station_name, max_temperature_group, count(*)
14 FROM temps_collapsed
15 GROUP BY station_name, max_temperature_group
16 ORDER BY station_name, count(*) DESC;

```

## 16.5 Dataset Output from Previous Query

station_name	max_temperature_group	count
CHICAGO NORTHERLY ISLAND IL US	Warm	133
CHICAGO NORTHERLY ISLAND IL US	Cold	92
CHICAGO NORTHERLY ISLAND IL US	Pleasant	91
CHICAGO NORTHERLY ISLAND IL US	Frigid	30
CHICAGO NORTHERLY ISLAND IL US	Inhumane	8
CHICAGO NORTHERLY ISLAND IL US	Hot	8
SEATTLE BOEING FIELD WA US	Pleasant	198
SEATTLE BOEING FIELD WA US	Warm	98
SEATTLE BOEING FIELD WA US	Cold	50
SEATTLE BOEING FIELD WA US	Hot	3
WAIKIKI 717.2 HI US	Warm	361
WAIKIKI 717.2 HI US	Hot	5

## 16.6 My Takeaway

Waikiki has 361 Warm days. Chicago has 8 Inhumane days. I know where I'd rather be.

...

## 16.7 Quick Check: CASE Evaluation

If a value of 95 is evaluated by a CASE statement where the first WHEN is `max_temp >= 70` and the second is `max_temp >= 90`, what happens?

- **A.** It matches the `>= 90` condition
- **B.** It matches the `>= 70` condition (first match wins)
- **C.** It matches both conditions
- **D.** It returns NULL

...

**B.** CASE evaluates top-to-bottom and returns the **first** match. `95 >= 70` is true, so it returns that result and stops. It never checks `>= 90`. This is why you put the most restrictive conditions first (`>= 90` before `>= 70`).

## 16.8 Try It: Population Tiers

Using a CTE with CASE on the `us_counties_pop_est_2019` table, classify each county as:

- 'Metro' (pop `>= 500,000`)
- 'Urban' (pop 100,000 to 499,999)
- 'Suburban' (pop 50,000 to 99,999)
- 'Rural' (pop `< 50,000`)

Then count the number of counties in each tier and compute the total population per tier. Which tier has the most counties? Which holds the most population?

Columns: `tier`, `num_counties`, `total_pop`, `avg_county_pop`.

```
1 -- Your query here
```

## 16.9 Population Tiers Query

```
1 WITH county_tiers AS (  
2     SELECT county_name, state_name, pop_est_2019,  
3         CASE  
4             WHEN pop_est_2019 >= 500000 THEN 'Metro'  
5             WHEN pop_est_2019 >= 100000 THEN 'Urban'  
6             WHEN pop_est_2019 >= 50000 THEN 'Suburban'  
7             ELSE 'Rural'  
8         END AS tier  
9     FROM us_counties_pop_est_2019  
10 )  
11 SELECT tier,  
12     count(*) AS num_counties,  
13     sum(pop_est_2019) AS total_pop,  
14     round(avg(pop_est_2019), 0) AS avg_county_pop  
15 FROM county_tiers  
16 GROUP BY tier  
17 ORDER BY total_pop DESC;
```

**Takeaway:** Rural has the most counties by far (~2,500+). Metro has the most total population. The vast majority of U.S. counties are small, but the vast majority of Americans live in big ones.

## 17 Part 8: Summary

### 17.1 What We Covered (a lot of ground! )

### 17.2 Key Techniques

Technique	What It Does	Textbook Listing
Scalar subquery in WHERE	Compute a threshold, filter on it	13-1
Subquery with DELETE	Filter rows to remove	13-2
Derived table (subquery in FROM)	Create a temporary table inline	13-3, 13-4
Column subquery in SELECT	Add a computed column per row	13-5, 13-6
IN with subquery	Check if value is in a result set	13-8
EXISTS / NOT EXISTS	Check membership (correlated)	13-9, 13-10
LATERAL in FROM	Reuse calculations across subqueries	13-11
LATERAL with JOIN	Per-row subquery (top-N per group)	13-12
CTE (WITH)	Named temporary result sets	13-13 to 13-15
crosstab()	Pivot long data to wide	13-17, 13-19
CASE	Conditional classification	13-20, 13-21

### 17.3 The Big Ideas

1. **Subqueries go everywhere:** WHERE, FROM, SELECT, even inside other subqueries and DML statements.
2. **Know your correlation:** Uncorrelated = runs once. Correlated = runs per row. The performance difference can be ginormous.

3. **EXISTS short-circuits.** For large-table membership checks, prefer it over IN. And never use NOT IN with nullable columns.
4. **LATERAL is the top-N-per-group tool.** “For each X, get the last N Y’s.” Learn it, love it.
5. **CTEs > nested subqueries** for readability. Always. Yes, I will die on this hill.
6. **crosstab() = long to wide.** Three columns in, pivot table out. Great for reporting.
7. **CASE then aggregate.** Classify first in a CTE, aggregate second in the outer query.

## 17.4 Up Next: In-Class Activity

Time to practice! Open the **Advanced Queries** assignment in CodeGrade.

You’ll work through problems covering:

- Scalar subqueries and derived tables
- Correlated subqueries and EXISTS
- LATERAL joins
- CTEs
- Cross tabulations and CASE

Use today’s lecture as your reference. You have ~30 minutes.

## 17.5 References

1. DeBarros, A. (2022). *Practical SQL* (2nd ed.). No Starch Press. Chapter 13.
2. PostgreSQL: [WITH Queries \(CTEs\)](#)
3. PostgreSQL: [Subquery Expressions](#)
4. PostgreSQL: [LATERAL Joins](#)
5. PostgreSQL: [tablefunc \(crosstab\)](#)