

# Mining Text Data

## DATA 351: Data Management with SQL

Lucas P. Cordova, Ph.D.

2026-04-06

This lecture covers string functions, POSIX regular expressions in PostgreSQL, extracting structured fields from messy narrative text, and full-text search with tsvector, tsquery, GIN indexes, ts\_headline, and ts\_rank. We use crime report narratives and State of the Union speeches. Based on Chapter 14 of Practical SQL, 2nd Edition.

### Table of contents

<b>1</b>	<b>What You'll Learn</b>	<b>1</b>
<b>2</b>	<b>Part 1: Environment</b>	<b>2</b>
<b>3</b>	<b>Part 2: String Functions</b>	<b>3</b>
<b>4</b>	<b>Part 3: Regular Expressions</b>	<b>6</b>
<b>5</b>	<b>Part 4: From Blob to Columns</b>	<b>9</b>
<b>6</b>	<b>Part 5: Full-Text Search</b>	<b>17</b>
<b>7</b>	<b>Part 6: Practice and Pitfalls</b>	<b>24</b>
<b>8</b>	<b>Summary</b>	<b>24</b>

## 1 What You'll Learn

### 1.1 Learning Objectives

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

1. Apply common string functions for trimming, case conversion, measuring length, and slicing substrings
2. Read basic POSIX regex notation (`\d`, `{n,m}`, groups, alternation) and use `~`, `~*`, `!~`, and `substring(... from ...)`
3. Use `regexp_match`, `regexp_matches`, `regexp_replace`, and `regexp_split_to_table` to transform text
4. Extract structured columns from semi-structured blobs using capture groups and array indexing
5. Build full-text search queries with `to_tsvector`, `to_tsquery`, the `@@` match operator, and GIN indexes
6. Present and rank search hits with `ts_headline` and `ts_rank` (including length normalization)

...

Text is not a luxury column type. It is where half your “CSV from a government website” problems live. Today we give SQL something sharper than hope.

## 2 Part 1: Environment

### 2.1 Create the Database and Load Data

```
1 createdb mining_text;
```

### 2.2 Create the Database and Load Data · Steps

Create a dedicated database for this chapter using whatever method you prefer or the command above in a terminal.

Download and run the setup script [mining\\_text.sql](#) ↓ while connected to `mining_text`. The file embeds all rows with `INSERT` statements (dollar-quoted text), so you can run it from pgAdmin or Beekeeper without configuring CSV paths.

Optional: [mining\\_text\\_csvs.zip](#) ↓ still has the raw extracts if you want to practice `\COPY` separately.

### 2.3 Tables You Will See Today

Table	Source	Role
county_regex_demo	bundled in mining_text.sql	Regex filtering practice (~, ~*)
crime_reports	embedded in mining_text.sql	Parse semi-structured police narratives
president_speeches	embedded in mining_text.sql	Full-text search on State of the Union addresses

## 2.4 Sanity Checks

```

1 SELECT count(*) FROM crime_reports;
2 SELECT count(*) FROM president_speeches;
3 SELECT count(*) FROM county_regex_demo;

```

## 2.5 Sanity Checks · Output

Run these after the setup script finishes to verify the counts of the tables.

```

1 SELECT count(*) FROM crime_reports;
2 -- 5
3 SELECT count(*) FROM president_speeches;
4 -- 79
5 SELECT count(*) FROM county_regex_demo;
6 -- 8

```

# 3 Part 2: String Functions

## 3.1 Why Strings Still Matter

JSON gets the hype. Arrays get the conference talks. Strings pay the rent: addresses, notes, log lines, pasted emails, and whatever your stakeholder calls “structured” because it has a colon somewhere.

PostgreSQL’s string toolkit is documented in the manual under [String Functions](#). We start with the ones you will actually type.

## 3.2 Casing and Cleanup

```
1 SELECT upper('Neal7');
2 SELECT lower('Randy');
3 SELECT initcap('at the end of the day');
4 SELECT initcap('Practical SQL'); -- acronyms: not perfect
```

## 3.3 Casing and Cleanup · Output

Run these in order and watch how forgiving real-world text is not.

```
1 SELECT upper('Neal7');
2 -- NEAL7
3 SELECT lower('Randy');
4 -- randy
5 SELECT initcap('at the end of the day');
6 -- At The End Of The Day
7 SELECT initcap('Practical SQL'); -- acronyms: not perfect
8 -- Practical Sql
```

`initcap` capitalizes word starts. It does not read minds. Your acronym table lives in application code, not here.

## 3.4 Length, Position, Trim

```
1 SELECT char_length(' Pat ');
2 SELECT length(' Pat '); -- same for varchar; byte length differs on some types
3 SELECT position(', ' in 'Tan, Bella');
4 SELECT trim('s' from 'socks');
5 SELECT trim(trailing 's' from 'socks');
6 SELECT trim(' Pat ');
7 SELECT char_length(trim(' Pat '));
8 SELECT ltrim('socks', 's');
9 SELECT rtrim('socks', 's');
```

## 3.5 Length, Position, Trim · Output

Run these in order.

```

1 SELECT char_length(' Pat ');
2 -- 5
3 SELECT length(' Pat ');          -- same for varchar; byte length differs on some types
4 -- 5
5 SELECT position(', ' in 'Tan, Bella');
6 -- 4
7 SELECT trim('s' from 'socks');
8 -- ock
9 SELECT trim(trailing 's' from 'socks');
10 -- sock
11 SELECT trim(' Pat ');
12 -- Pat
13 SELECT char_length(trim(' Pat '));
14 -- 3
15 SELECT ltrim('socks', 's');
16 -- ocks
17 SELECT rtrim('socks', 's');
18 -- sock

```

If you are counting characters for validation rules, `char_length` after `trim` is the boring, correct path.

### 3.6 Slicing and Replacing

```

1 SELECT left('703-555-1212', 3);
2 SELECT right('703-555-1212', 8);
3 SELECT replace('bat', 'b', 'c');
4 SELECT replace('aaa', 'aa', 'b');

```

### 3.7 Slicing and Replacing · Output

Run these (including the `replace('aaa', ...)` check).

```

1 SELECT left('703-555-1212', 3);
2 -- 703
3 SELECT right('703-555-1212', 8);
4 -- 555-1212
5 SELECT replace('bat', 'b', 'c');
6 -- cat
7 SELECT replace('aaa', 'aa', 'b');
8 -- ba

```

**Quick check:** `replace('aaa', 'aa', 'b')` returns `ba` (left-to-right, non-overlapping substitution). String replacement is greedy in a way that will eventually humble everyone.

## 4 Part 3: Regular Expressions

### 4.1 What is a regular expression?

A **regular expression** (regex) is a small pattern language for describing text shape. You hand PostgreSQL a pattern string; the engine can test a match, return the matched substring, split on delimiters, or pull out captured pieces.

PostgreSQL uses **POSIX** regular expressions (documented under [Pattern Matching](#)). In SQL the pattern usually lives in a single-quoted string. Backslashes are easy to get wrong: the string parser and the regex engine both care about `\`, so you often write `'\\d'` (two characters in source, one backslash in the pattern) or use dollar-quoting when patterns get busy.

### 4.2 Notation

Piece	Meaning
<code>.</code>	one character (with caveats: not always newline)
<code>\d</code>	a digit 0–9
<code>\w</code>	a word character (letters, digits, underscore in the POSIX class)
<code>[0–9]</code>	any single digit from the list
<code>{n}</code>	repeat the previous “atom” exactly <code>n</code> times
<code>{n,m}</code>	repeat at least <code>n</code> and at most <code>m</code> times, so <code>{1,2}</code> means one or two
<code>+</code>	one or more; <code>*</code> zero or more; <code>?</code> optional (zero or one)
<code>^</code> <code>\$</code>	start and end of the string (for these lectures, think whole-field matching unless noted)
<code>\ </code>	alternation: <code>A\ B</code> is A or B
<code>()</code>	grouping; also <b>capture</b> for <code>regexp_match</code>
<code>(?: ... )</code>	group <b>without</b> capturing (still useful for precedence)

### 4.3 Example

`\d{1,2}/\d{1,2}/\d{2}` describes a date fragment like 4/16/17: one or two digits, slash, one or two digits, slash, two digits. Slashes are literal here; we sometimes escape them as `\/` in patterns for clarity or habit from other tools.

### 4.4 Operators You Must Know

Operator	Meaning
<code>~</code>	matches, case sensitive
<code>~*</code>	matches, case insensitive
<code>!~</code>	does not match, case sensitive
<code>!~*</code>	does not match, case insensitive

In `WHERE` clauses, `regex` turns a pile of `LIKE` patterns into one expressive filter. It also turns a simple bug into a performance mystery. Use indexes thoughtfully (later: GIN for FTS; for `regex` alone, sometimes a functional index, sometimes a redesign).

### 4.5 Substring with a Pattern

```
1 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '.+');
2 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '\d{1,2} (?:a.m.|p.m.)');
3 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '^w+');
4 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from 'w+.$');
5 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from 'May|June');
6 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '\d{4}');
7 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from 'May \d, \d{4}');
```

### 4.6 Substring with a Pattern · Output

`substring(string from pattern)` returns the first match. **Run them** one at a time.

```
1 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '.+');
2 -- The game starts at 7 p.m. on May 2, 2024.
3 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '\d{1,2} (?:a.m.|p.m.)');
4 -- 7 p.m.
5 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '^w+');
6 -- The
7 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from 'w+.$');
```

```

8 -- 2024.
9 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from 'May|June');
10 -- May
11 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from '\d{4}');
12 -- 2024
13 SELECT substring('The game starts at 7 p.m. on May 2, 2024.' from 'May \d, \d{4}');
14 -- May 2, 2024

```

Notice how greedy `.+` is on the first pattern. Regex is a contract between you and the parser. Read the contract.

## 4.7 Filtering Rows

```

1 SELECT county_name
2 FROM county_regex_demo
3 WHERE county_name ~* '(lade|lare)'
4 ORDER BY county_name;
5
6 SELECT county_name
7 FROM county_regex_demo
8 WHERE county_name ~* 'ash' AND county_name !~ 'Wash'
9 ORDER BY county_name;

```

## 4.8 Filtering Rows · Output

We use `county_regex_demo` instead of the full Census extract:

```

1 SELECT county_name
2 FROM county_regex_demo
3 WHERE county_name ~* '(lade|lare)'
4 ORDER BY county_name;
5 -- Clare County
6 SELECT county_name
7 FROM county_regex_demo
8 WHERE county_name ~* 'ash' AND county_name !~ 'Wash'
9 ORDER BY county_name;
10 -- Ashland County
11 -- Nash County

```

**Run both queries.** Advance once per query, then per result line.

The second query is the classic “I want Ashland but not Washington” maneuver. If your data contains NULL county names, remember that `WHERE col ~ 'pattern'` does not match NULL. Neither does your excuse.

## 4.9 Replacing and Splitting

```
1 SELECT regexp_replace('05/12/2024', '\d{4}', '2023');
2 SELECT regexp_split_to_table('Four,score,and,seven,years,ago', ',');
3 SELECT regexp_split_to_array('Phil Mike Tony Steve', ' ');
4 SELECT array_length(regexp_split_to_array('Phil Mike Tony Steve', ' '), 1);
```

## 4.10 Replacing and Splitting · Output

Run these.

```
1 SELECT regexp_replace('05/12/2024', '\d{4}', '2023');
2 -- 05/12/2023
3 SELECT regexp_split_to_table('Four,score,and,seven,years,ago', ',');
4 -- Four
5 -- score
6 -- and
7 -- seven
8 -- years
9 -- ago
10 SELECT regexp_split_to_array('Phil Mike Tony Steve', ' ');
11 -- {Phil, Mike, Tony, Steve}
12 SELECT array_length(regexp_split_to_array('Phil Mike Tony Steve', ' '), 1);
13 -- 4
```

`regexp_split_to_table` is a set-returning function: it produces one row per piece. That makes it handy for normalizing multi-valued cells someone swore were “atomic.”

# 5 Part 4: From Blob to Columns

## 5.1 The Crime Narratives

```
1 SELECT crime_id, original_text
2 FROM crime_reports
3 ORDER BY crime_id
4 LIMIT 1;
```

## 5.2 The Crime Narratives · Output

Our `crime_reports` table loads only `original_text`. Everything else (dates, address, offense, case number) is hiding inside the blob.

Run this to open one row:

```
1 SELECT original_text
2 FROM crime_reports
3 ORDER BY crime_id
4 LIMIT 1;
5 -- 4/16/17-4/17/17
6 -- 2100-0900 hrs.
7 -- 46000 Block Ashmere Sq.
8 -- Sterling
9 -- Larceny: The victim reported that a
10 -- bicycle was stolen from their opened
11 -- garage door during the overnight hours.
12 -- C0170006614
```

Beautiful. Readable. Terrible for `GROUP BY`. We will fix that with patterns, not with interns.

## 5.3 First Date Match

Run this.

```
1 SELECT crime_id,
2         regexp_match(original_text, '\d{1,2}\/\d{1,2}\/\d{2}')
3 FROM crime_reports
4 ORDER BY crime_id;
```

## 5.4 First Date Match · Output

`regexp_match` returns the first match as a text array (or `NULL`).

```
1 SELECT crime_id,
2         regexp_match(original_text, '\d{1,2}\/\d{1,2}\/\d{2}')
3 FROM crime_reports
4 ORDER BY crime_id;
5 -- 1 | {4/16/17}
6 -- 2 | {4/8/17}
7 -- 3 | {4/4/17}
```

```
8 -- 4 | {04/10/17}
9 -- 5 | {04/09/17}
```

## 5.5 All Date Matches

Row 1 has two dates in the narrative.

```
1 4/16/17-4/17/17
2 2100-0900 hrs.
3 46000 Block Ashmere Sq.
4 Sterling
. . .
```

**Run this** and compare to the previous slide.

```
1 SELECT crime_id,
2        regexp_matches(original_text, '\d{1,2}\/\d{1,2}\/\d{2}', 'g')
3 FROM crime_reports
4 ORDER BY crime_id;
```

## 5.6 All Date Matches · Output

`regexp_matches` with the `g` flag finds every non-overlapping match.

```
1 SELECT crime_id,
2        regexp_matches(original_text, '\d{1,2}\/\d{1,2}\/\d{2}', 'g')
3 FROM crime_reports
4 ORDER BY crime_id;
5 -- 1 | {4/16/17}
6 -- 1 | {4/17/17}
7 -- 2 | {4/8/17}
8 -- 3 | {4/4/17}
9 -- 4 | {04/10/17}
10 -- 5 | {04/09/17}
```

**Compare** the two outputs for rows that mention more than one calendar date.

## 5.7 Capture Groups Clean Up Noise

Run these:

```
1 SELECT crime_id,
2     regexp_match(original_text, '-\d{1,2}\/\d{1,2}\/\d{2}')
3 FROM crime_reports
4 ORDER BY crime_id;
5 SELECT crime_id,
6     regexp_match(original_text, '-(\d{1,2}\/\d{1,2}\/\d{2})')
7 FROM crime_reports
8 ORDER BY crime_id;
```

## 5.8 Capture Groups · Output

Parentheses in the pattern define capture groups. `regexp_match` still returns the whole match as element [1] when you wrap only part of the pattern. **Run both.**

```
1 SELECT crime_id,
2     regexp_match(original_text, '-\d{1,2}\/\d{1,2}\/\d{2}')
3 FROM crime_reports
4 ORDER BY crime_id;
5 -- 1 | {-4/17/17}
6 -- 2 |
7 -- 3 |
8 -- 4 |
9 -- 5 |
10 SELECT crime_id,
11     regexp_match(original_text, '-(\d{1,2}\/\d{1,2}\/\d{2})')
12 FROM crime_reports
13 ORDER BY crime_id;
14 -- 1 | {4/17/17}
15 -- 2 |
16 -- 3 |
17 -- 4 |
18 -- 5 |
```

The second version drops the leading hyphen by capturing only the date digits.

## 5.9 Pull Several Fields at Once

Run this:

```

1 SELECT
2     crime_id,
3     regexp_match(original_text, '(?:C0|S0)[0-9]+') AS case_number,
4     regexp_match(original_text, '\d{1,2}\/\d{1,2}\/\d{2}') AS date_1,
5     regexp_match(original_text, '\n(?:\w+ \w+|\w+)\n(.*):') AS crime_type,
6     regexp_match(original_text, '(?:Sq.|Plz.|Dr.|Ter.|Rd.)\n(\w+ \w+|\w+)\n')
7         AS city
8 FROM crime_reports
9 ORDER BY crime_id;

```

## 5.10 Pull Several Fields · Output

This is not pretty. It is honest.

```

1 SELECT
2     crime_id,
3     regexp_match(original_text, '(?:C0|S0)[0-9]+') AS case_number,
4     regexp_match(original_text, '\d{1,2}\/\d{1,2}\/\d{2}') AS date_1,
5     regexp_match(original_text, '\n(?:\w+ \w+|\w+)\n(.*):') AS crime_type,
6     regexp_match(original_text, '(?:Sq.|Plz.|Dr.|Ter.|Rd.)\n(\w+ \w+|\w+)\n')
7         AS city
8 FROM crime_reports
9 ORDER BY crime_id;
10 -- 1 | {C0170006614} | {4/16/17} | {Larceny} | {Sterling}
11 -- 2 | {C0170006162} | {4/8/17} | {Destruction of Property} | {Sterling}
12 -- 3 | {C0170006079} | {4/4/17} | {Larceny} | {Sterling}
13 -- 4 | {S0170006250} | {04/10/17} | {Larceny} | {Middleburg}
14 -- 5 | {S0170006211} | {04/09/17} | {Destruction of Property} | {Sterling}

```

Patterns like this are why we version-control SQL. The day the sheriff’s office changes a label from `hrs.` to `hours`, you will discover the true meaning of “breaking change.”

## 5.11 Arrays in PostgreSQL: [1]

Run this:

```

1 SELECT
2     crime_id,
3     (regexp_match(original_text, '(?:C0|S0)[0-9]+'))[1] AS case_number
4 FROM crime_reports
5 ORDER BY crime_id;

```

## 5.12 Arrays in PostgreSQL · Output

`regexp_match` returns `text[]`. Grab the first element explicitly.

```
1 SELECT
2     crime_id,
3     (regexp_match(original_text, '(?:C0|S0)[0-9]+'))[1] AS case_number
4 FROM crime_reports
5 ORDER BY crime_id;
6 -- 1 | C0170006614
7 -- 2 | C0170006162
8 -- 3 | C0170006079
9 -- 4 | S0170006250
10 -- 5 | S0170006211
```

If the match fails, `regexp_match` is `NULL` and `[1]` is `NULL`. Plan for that when you cast to timestamps.

## 5.13 Optional: Populate Columns

```
1 -- Pattern sketch (abbreviated): build a string, cast to timestampz
2 -- ((regexp_match(...))[1] || ' ' || (regexp_match(...))[1])::timestampz
```

## 5.14 Optional: Populate Columns · Notes

The textbook walks through `UPDATE` statements that cast extracted dates to `timestampz` and fill `street`, `city`, `crime_type`, and `description`. The expressions are long but repetitive: once you trust one `regexp_match`, the rest are the same idea with different patterns.

**Question for you:** Why might `EXISTS (SELECT regexp_matches(...))` appear in a `CASE` branch?

## 5.15 Answer

Because sometimes you need to know whether a second date exists before you try to parse it.

## 5.16 Walkthrough: Build date\_1 Once

Don't run this (unless you are using a scratch database or copy of the crime\_reports table):

```
1 CREATE temp TABLE crime_reports_copy AS SELECT * FROM crime_reports;
2
3 UPDATE crime_reports
4 SET date_1 =
5 (
6     (regexp_match(original_text, '\d{1,2}\.\d{1,2}\.\d{2}'))[1]
7     || ' ' ||
8     (regexp_match(original_text, '\.\d{2}\n(\d{4})'))[1]
9     || ' US/Eastern'
10 )::timestampz
11 RETURNING crime_id, date_1, original_text;
```

## 5.17 Walkthrough: Build date\_1 · Output

The book example concatenates the first date in MM/DD/YY form with the first time field (four digits) and a time zone. **Run once** (it updates every row). Use a scratch database or ROLLBACK if you are experimenting.

```
1 UPDATE crime_reports_copy
2 SET date_1 =
3 (
4     (regexp_match(original_text, '\d{1,2}\.\d{1,2}\.\d{2}'))[1]
5     || ' ' ||
6     (regexp_match(original_text, '\.\d{2}\n(\d{4})'))[1]
7     || ' US/Eastern'
8 )::timestampz
9 RETURNING crime_id, date_1, original_text;
10 -- 1 | 2017-04-17 01:00:00+00 | ...
11 -- 2 | 2017-04-08 20:00:00+00 | ...
12 -- 3 | 2017-04-04 18:00:00+00 | ...
13 -- 4 | 2017-04-10 20:05:00+00 | ...
14 -- 5 | 2017-04-09 16:00:00+00 | ...
```

## 5.18 Why RETURNING?

RETURNING shows you what you broke before you commit. This is not Tinder. Do not swipe right on untested regex.

If a row lacks one of the pieces, the corresponding `regexp_match` is `NULL` and the concatenation collapses in ways PostgreSQL will explain loudly. That is good. Silent failure is how you become a historian instead of a data engineer.

## 5.19 View the Cleaned Rows

Run this:

```
1 SELECT date_1,
2         street,
3         city,
4         crime_type
5 FROM crime_reports_copy
6 ORDER BY crime_id;
```

## 5.20 View the Cleaned Rows · Output

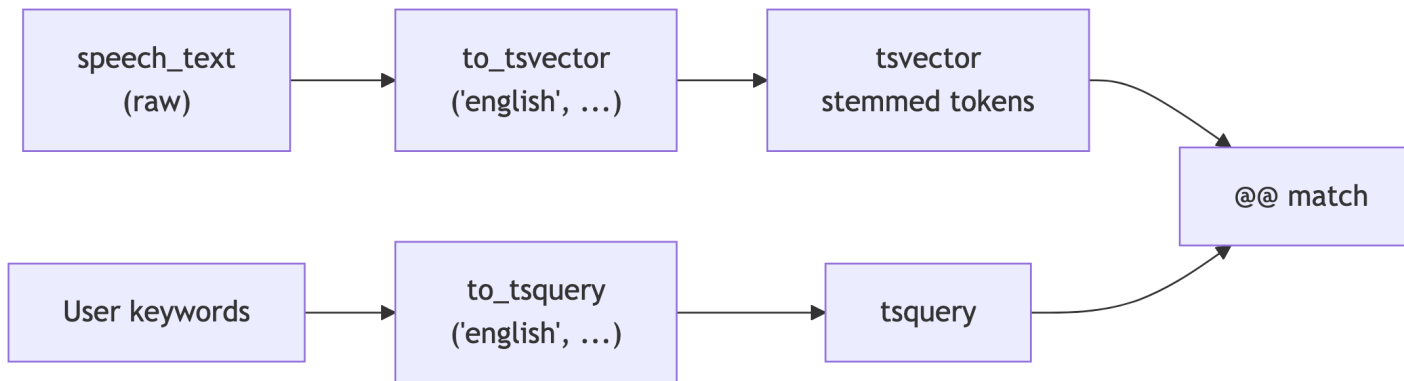
After a full `UPDATE` of all parsed columns (as in the book), a simple report query becomes possible. The `UPDATE` sets `date_1` only (other parsed columns stay `NULL` until you run the full `UPDATE` block that sets all parsed columns).

```
1 SELECT date_1,
2         street,
3         city,
4         crime_type
5 FROM crime_reports_copy
6 ORDER BY crime_id;
7 -- 2017-04-17 01:00:00+00 | | |
8 -- 2017-04-08 20:00:00+00 | | |
9 -- 2017-04-04 18:00:00+00 | | |
10 -- 2017-04-10 20:05:00+00 | | |
11 -- 2017-04-09 16:00:00+00 | | |
```

Compare mentally to the blob you started with. That difference is the whole lecture in one `SELECT`.

## 6 Part 5: Full-Text Search

### 6.1 Pipeline: From Speech to Search



Raw text enters, tokens come out, queries meet documents at @@. The GIN index sits on the `tsvector` column so the middle of this diagram does not require a sequential read of every presidency.

### 6.2 `pg_ts_config`

Run this:

```
1 SELECT cfgname FROM pg_ts_config;
```

### 6.3 `pg_ts_config` · Output

Full-text search is not `LIKE '%tax%'`. It tokenizes, stems, drops noise words, and gives you operators that behave well at scale.

Check which language configurations exist.

```
1 SELECT cfgname FROM pg_ts_config;
2 -- arabic
3 -- armenian
4 -- basque
5 -- catalan
6 -- danish
7 -- dutch
8 -- english
```

```
9 -- finnish
10 -- ... (more rows on your server)
```

## 6.4 to\_tsvector and @@

```
1 SELECT to_tsvector('english', 'I am walking across the sitting room to sit with you.');
```

```
2
```

```
3 SELECT to_tsquery('english', 'walking & sitting');
```

```
4
```

```
5 SELECT to_tsvector('english', 'I am walking across the sitting room')
6         @@ to_tsquery('english', 'walking & sitting');
```

```
7
```

```
8 SELECT to_tsvector('english', 'I am walking across the sitting room')
9         @@ to_tsquery('english', 'walking & running');
```

## 6.5 to\_tsvector and @@ · Output

Convert document and query:

**Run these.**

```
1 SELECT to_tsvector('english', 'I am walking across the sitting room to sit with you.');
```

```
2 -- 'across':4 'room':7 'sit':6,9 'walk':3
```

```
3
```

```
4 SELECT to_tsquery('english', 'walking & sitting');
```

```
5 -- 'walk' & 'sit'
```

```
6
```

```
7 SELECT to_tsvector('english', 'I am walking across the sitting room')
8         @@ to_tsquery('english', 'walking & sitting');
```

```
9 -- t
```

```
10
```

```
11 SELECT to_tsvector('english', 'I am walking across the sitting room')
12         @@ to_tsquery('english', 'walking & running');
```

```
13 -- f
```

@@ answers: does this document match this query?

## 6.6 Operators Inside to\_tsquery

- & AND
- | OR
- ! NOT

- <-> adjacent terms (phrase)
- <N> distance (terms within N positions)

We will lean on the State of the Union table. Our setup script already adds `search_speech_text` and a GIN index.

## 6.7 Index for Search

```
1 -- Already in mining_text.sql after INSERT + UPDATE:
2 -- CREATE INDEX search_idx ON president_speeches USING gin (search_speech_text);
```

## 6.8 Index for Search · Notes

GIN is the usual choice for `tsvector` columns. Without it, you are scanning speeches the way I scroll my inbox: slowly and with regret.

## 6.9 Find Speeches

Run this:

```
1 SELECT president, speech_date
2 FROM president_speeches
3 WHERE search_speech_text @@ to_tsquery('english', 'Vietnam')
4 ORDER BY speech_date;
```

## 6.10 Find Speeches · Output

```
1 SELECT president, speech_date
2 FROM president_speeches
3 WHERE search_speech_text @@ to_tsquery('english', 'Vietnam')
4 ORDER BY speech_date;
5 -- John F. Kennedy | 1961-05-25
6 -- Lyndon B. Johnson | 1966-01-12
7 -- Lyndon B. Johnson | 1967-01-10
8 -- ... (16 more rows; 19 total in this bundle)
```

## 6.11 Snippets with ts\_headline

Run this:

```
1 SELECT president,
2     speech_date,
3     ts_headline(
4         speech_text,
5         to_tsquery('english', 'tax'),
6         'StartSel = <,
7         StopSel = >,
8         MinWords=5,
9         MaxWords=7,
10        MaxFragments=1'
11     )
12 FROM president_speeches
13 WHERE search_speech_text @@ to_tsquery('english', 'tax')
14 ORDER BY speech_date;
```

## 6.12 Snippets with ts\_headline · Output

```
1 SELECT president,
2     speech_date,
3     ts_headline(
4         speech_text,
5         to_tsquery('english', 'tax'),
6         'StartSel = <,
7         StopSel = >,
8         MinWords=5,
9         MaxWords=7,
10        MaxFragments=1'
11     )
12 FROM president_speeches
13 WHERE search_speech_text @@ to_tsquery('english', 'tax')
14 ORDER BY speech_date;
15 -- Harry S. Truman | 1946-01-21 | ... <tax> POLICY ...
16 -- Harry S. Truman | 1947-01-06 | excise <tax> rates which, under the present
17 -- Harry S. Truman | 1948-01-07 | point in our <tax> structure. ...
18 -- ... (many more rows mention tax)
```

That angle-bracket markup is for HTML consumers. In slides, it still reads clearly as “here is the hit in context.”

## 6.13 AND / NOT and Phrase Proximity

Run each query:

```
1 SELECT president,
2     speech_date,
3     ts_headline(
4         speech_text,
5         to_tsquery('english', 'transportation & !roads'),
6         'StartSel = <, StopSel = >, MinWords=5, MaxWords=7, MaxFragments=1'
7     )
8 FROM president_speeches
9 WHERE search_speech_text @@ to_tsquery('english', 'transportation & !roads')
10 ORDER BY speech_date;
11 SELECT president,
12     speech_date,
13     ts_headline(
14         speech_text,
15         to_tsquery('english', 'military <-> defense'),
16         'StartSel = <, StopSel = >, MinWords=5, MaxWords=7, MaxFragments=1'
17     )
18 FROM president_speeches
19 WHERE search_speech_text @@ to_tsquery('english', 'military <-> defense')
20 ORDER BY speech_date;
```

## 6.14 AND / NOT and Phrase Proximity · Output

```
1 SELECT president,
2     speech_date,
3     ts_headline(
4         speech_text,
5         to_tsquery('english', 'transportation & !roads'),
6         'StartSel = <, StopSel = >, MinWords=5, MaxWords=7, MaxFragments=1'
7     )
8 FROM president_speeches
9 WHERE search_speech_text @@ to_tsquery('english', 'transportation & !roads')
10 ORDER BY speech_date;
11 -- Harry S. Truman | 1947-01-06 | Mr. President, Mr. Speaker, Members
12 -- Harry S. Truman | 1949-01-05 | Mr. President, Mr. Speaker, Members
13 SELECT president,
14     speech_date,
15     ts_headline(
```

```

16         speech_text,
17         to_tsquery('english', 'military <-> defense'),
18         'StartSel = <, StopSel = >, MinWords=5, MaxWords=7, MaxFragments=1'
19     )
20 FROM president_speeches
21 WHERE search_speech_text @@ to_tsquery('english', 'military <-> defense')
22 ORDER BY speech_date;
23 -- Dwight D. Eisenhower | 1956-01-05 | To the Congress of the
24 -- Dwight D. Eisenhower | 1958-01-09 | Mr. President, Mr. Speaker, Members
25 -- Dwight D. Eisenhower | 1959-01-09 | Mr. President, Mr. Speaker, Members

```

Try <2> instead of <-> in a copy of the query to allow a small gap between words. Language is messy. PostgreSQL knows.

## 6.15 Ranking

Run both and compare ordering.

```

1  SELECT president,
2         speech_date,
3         ts_rank(
4             search_speech_text,
5             to_tsquery('english', 'war & security & threat & enemy')
6         ) AS score
7 FROM president_speeches
8 WHERE search_speech_text @@
9        to_tsquery('english', 'war & security & threat & enemy')
10 ORDER BY score DESC
11 LIMIT 5;
12 SELECT president,
13         speech_date,
14         ts_rank(
15             search_speech_text,
16             to_tsquery('english', 'war & security & threat & enemy'),
17             2
18         )::numeric AS score
19 FROM president_speeches
20 WHERE search_speech_text @@
21        to_tsquery('english', 'war & security & threat & enemy')
22 ORDER BY score DESC
23 LIMIT 5;

```

## 6.16 Ranking · Output

```
1 SELECT president,
2     speech_date,
3     ts_rank(
4         search_speech_text,
5         to_tsquery('english', 'war & security & threat & enemy')
6     ) AS score
7 FROM president_speeches
8 WHERE search_speech_text @@
9     to_tsquery('english', 'war & security & threat & enemy')
10 ORDER BY score DESC
11 LIMIT 5;
12 -- William J. Clinton | 1997-02-04 | 0.35810584
13 -- George W. Bush | 2004-01-20 | 0.29587495
14 -- George W. Bush | 2003-01-28 | 0.28381455
15 -- Harry S. Truman | 1946-01-21 | 0.25752166
16 -- William J. Clinton | 2000-01-27 | 0.22214262
17 SELECT president,
18     speech_date,
19     ts_rank(
20         search_speech_text,
21         to_tsquery('english', 'war & security & threat & enemy'),
22         2
23     )::numeric AS score
24 FROM president_speeches
25 WHERE search_speech_text @@
26     to_tsquery('english', 'war & security & threat & enemy')
27 ORDER BY score DESC
28 LIMIT 5;
29 -- George W. Bush | 2004-01-20 | 0.000103
30 -- William J. Clinton | 1997-02-04 | 0.000098
31 -- George W. Bush | 2003-01-28 | 0.000096
32 -- Jimmy Carter | 1979-01-23 | 0.000090
33 -- Lyndon B. Johnson | 1968-01-17 | 0.000073
```

Long speeches accumulate token hits. Normalization option 2 divides rank by document length.

Compare the ordering before and after normalization. This is the difference between “mentions the keywords a lot” and “mentions them a lot relative to length.”

The book points to a live demo at [anthonydebarros.com/sotu](http://anthonydebarros.com/sotu). Worth clicking after class.

## 7 Part 6: Practice and Pitfalls

### 7.1 Try It: Crime Case Numbers

Write a query that returns `crime_id` and the case number extracted from `original_text` using the same `(?:C0|S0)[0-9]+` pattern, but filter to rows where the case starts with `S0`.

**Stretch:** Add a column that counts how many distinct dates appear in the narrative using `regexp_matches` with the `g` flag and a subquery or lateral trick. (If you get stuck, move on. The point is to think in matches, not to win a Nobel prize in SQL golf.)

### 7.2 Try It: Speeches

Find speeches where the query `economy & jobs` matches, show `ts_headline`, and sort by normalized rank (`ts_rank` with the length divisor). Pick a president you have opinions about. Keep the opinions out of the query.

### 7.3 Common Footguns (Collected)

- Forgetting that `regexp_match` returns `NULL` when nothing matches, then indexing `[1]` in a cast chain
- Using `LIKE` patterns with regex operators by mistake (they are not interchangeable)
- Building huge `OR` chains when a single well-factored regex would do
- Skipping the GIN index on `tsvector` columns and wondering why search “works on my laptop”

## 8 Summary

### 8.1 What We Covered

---

Topic	Core functions and ideas
Strings	<code>upper</code> , <code>lower</code> , <code>initcap</code> , <code>trim</code> , <code>left/right</code> , <code>replace</code> , <code>char_length</code>
Regex filters	<code>~</code> , <code>~*</code> , <code>!~</code> , <code>substring</code> , <code>regexp_replace</code> , <code>splits</code>
Structured extraction	<code>regexp_match</code> , <code>regexp_matches</code> , capture groups, <code>[1]</code> indexing
Full-text search	<code>to_tsvector</code> , <code>to_tsquery</code> , <code>@@</code> , <code>ts_headline</code> , <code>ts_rank</code>
Performance	GIN index on <code>tsvector</code>

---

## 8.2 The One-Sentence Version

Treat text like data: measure it, constrain it, extract what you need, and index what you search. Everything else is a spreadsheet someone emailed you as prose.

## 8.3 References

1. DeBarros, A. (2022). *Practical SQL* (2nd ed.). No Starch Press. Chapter 14.
2. PostgreSQL: [Pattern Matching](#)
3. PostgreSQL: [String Functions](#)
4. PostgreSQL: [Regular Expressions](#)
5. PostgreSQL: [Full Text Search](#)
6. PostgreSQL: [ts\\_rank](#)