

---

# **Introduction to PostGIS**

***Release 1.0***

**Mayra Zurbarán, Mark Leslie, Paul Ramsey**

**October 07, 2019**



## CONTENTS

<b>1</b>	<b>Welcome</b>	<b>1</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
<b>3</b>	<b>Setup</b>	<b>11</b>
<b>4</b>	<b>Basic</b>	<b>39</b>
<b>5</b>	<b>Advanced</b>	<b>103</b>
<b>6</b>	<b>Maintenance</b>	<b>155</b>
<b>7</b>	<b>Appendix A: PostGIS Functions</b>	<b>177</b>
<b>8</b>	<b>Appendix B: Glossary</b>	<b>179</b>
<b>9</b>	<b>Appendix C: License</b>	<b>181</b>
<b>Index</b>		<b>183</b>



---

CHAPTER  
ONE

---

WELCOME

## 1.1 Workshop Conventions

These sections conform to a number of conventions to make it easier to follow the conversation. This section gives a brief overview of what to expect in the way of typographic conventions, as well as a short overview of the structure of each workbook.

### 1.1.1 Directions

Directions for you, the workshop attendee, will be noted by **bold** font.

For example:

Click **Next** to continue.

### 1.1.2 Code

SQL query examples will be displayed in an offset box

```
SELECT postgis_full_version();
```

These examples can be entered into the query window or command line interface.

### 1.1.3 Notes

Notes are used to provide information that is useful but not critical to the overall understanding of the topic.

---

**Note:** If you haven't eaten an apple today, the doctor may be on the way.

---

### 1.1.4 Functions

Where function names are defined in the text, they will be rendered in a **bold** font.

For example:

**ST\_Touches(geometry A, geometry B)** returns TRUE if either of the geometries' boundaries intersect

### 1.1.5 Files, Tables and Column Names

File names, paths, table names and column names will be shown in fixed-width font.

For example:

Select the `name` column in the `nyc_streets` table.

### 1.1.6 Menus and Form elements

Menus/submenus and form elements such as fields or check boxes and other on-screen artifacts are displayed in *italics*.

For example:

Click on the *File > New* menu. Check the box that says *Confirm*.

### 1.1.7 Workflow

Sections are designed to be progressive. Each section will start with the assumption that you have completed and understood the previous section in the series and will build on that knowledge. A single section will progress through a handful of ideas and provide working examples wherever possible. At the end of a section, where appropriate, we have included a handful of exercises to allow you to try out the ideas we've presented. In some cases the section will include "Things To Try". These tasks contain more complex problems than the exercises and is designed to challenge participants with advanced knowledge.

## INTRODUCTION

### 2.1 What is a Spatial Database?

PostGIS is a spatial database. Oracle Spatial and SQL Server (2008 and later) are also spatial databases. But what does that mean; what is it that makes an ordinary database a spatial database?

The short answer, is...

**Spatial databases store and manipulate spatial objects like any other object in the database.**

The following briefly covers the evolution of spatial databases, and then reviews three aspects that associate *spatial* data with a database – data types, indexes, and functions.

1. **Spatial data types** refer to shapes such as point, line, and polygon;
2. Multi-dimensional **spatial indexing** is used for efficient processing of spatial operations;
3. **Spatial functions**, posed in [SQL](#), are for querying of spatial properties and relationships.

Combined, spatial data types, indexes, and functions provide a flexible structure for optimized performance and analysis.

#### 2.1.1 In the Beginning

In legacy first-generation [GIS](#) implementations, all spatial data is stored in flat files and special [GIS](#) software is required to interpret and manipulate the data. These first-generation management systems are designed to meet the needs of users where all required data is within the user's organizational domain. They are proprietary, self-contained systems specifically built for handling spatial data.

Second-generation spatial systems store some data in relational databases (usually the “attribute” or non-spatial parts) but still lack the flexibility afforded with direct integration.

**True spatial databases were born when people started to treat spatial features as first class database objects.**

Spatial databases fully integrate spatial data with an object relational database. The orientation changes from GIS-centric to database-centric.

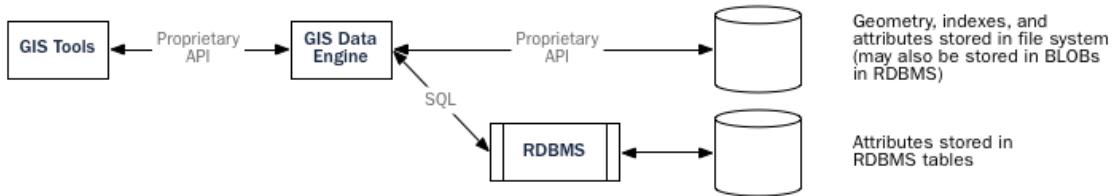
### Evolution of GIS Architectures

#### First-Generation GIS:



Geometry, indexes, and attributes stored in file system

#### Second-Generation GIS:



Geometry, indexes, and attributes stored in file system  
(may also be stored in BLOBs in RDBMS)

Attributes stored in RDBMS tables

#### Third-Generation GIS:



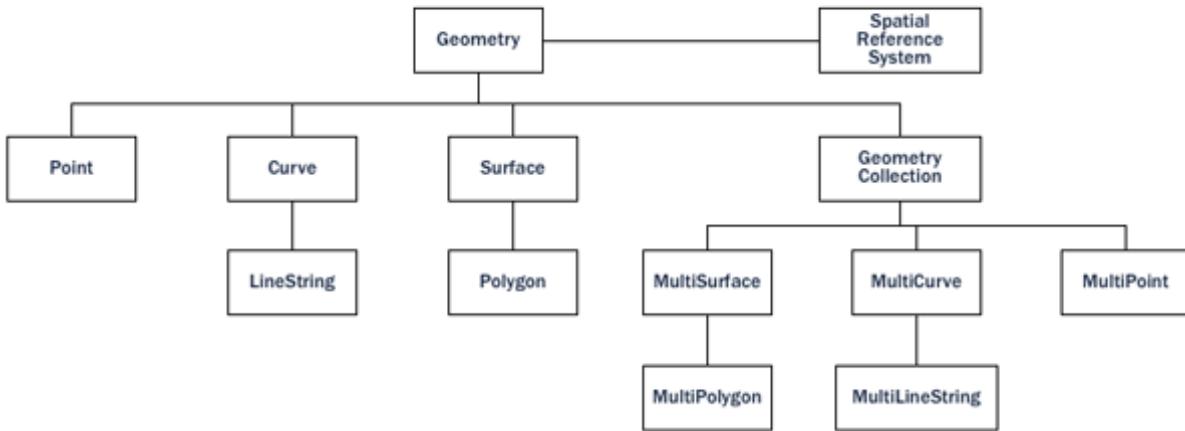
Geometry stored in ADTs in RDBMS tables with related business data

**Note:** A spatial database management system may be used in applications besides the geographic world. Spatial databases are used to manage data related to the anatomy of the human body, large-scale integrated circuits, molecular structures, and electro-magnetic fields, among others.

### 2.1.2 Spatial Data Types

An ordinary database has strings, numbers, and dates. A spatial database adds additional (spatial) types for representing **geographic features**. These spatial data types abstract and encapsulate spatial structures such as boundary and dimension. In many respects, spatial data types can be understood simply as shapes.

## Geometry Hierarchy



Spatial data types are organized in a type hierarchy. Each sub-type inherits the structure (attributes) and the behavior (methods or functions) of its super-type.

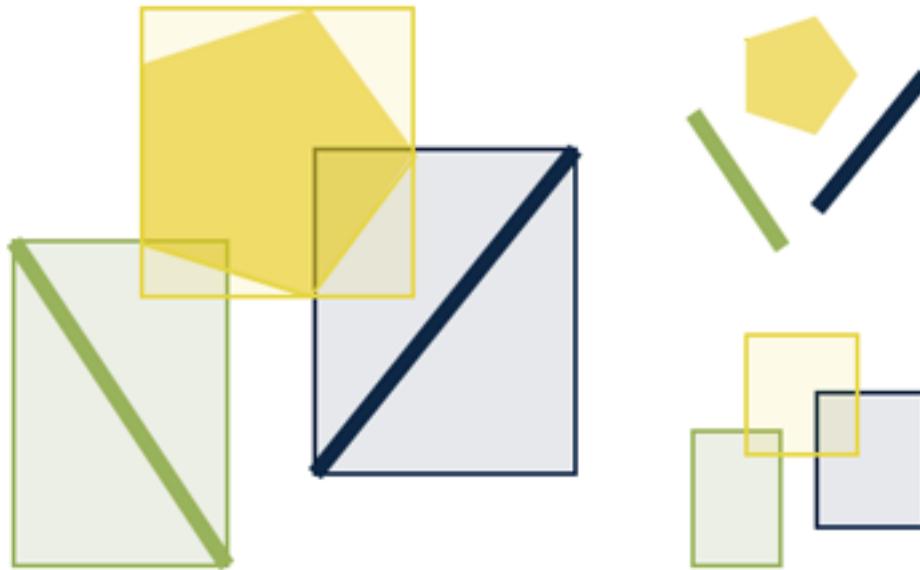
### 2.1.3 Spatial Indexes and Bounding Boxes

An ordinary database provides “access methods” – commonly known as **indexes** – to allow for fast and random access to subsets of data. Indexing for standard types (numbers, strings, dates) is usually done with **B-tree** indexes. A B-tree partitions the data using the natural sort order to put the data into a hierarchical tree.

The natural sort order of numbers, strings, and dates is simple to determine – every value is less than, greater than or equal to every other value. But because polygons can overlap, can be contained in one another, and are arrayed in a two-dimensional (or more) space, a B-tree cannot be used to efficiently index them. Real spatial databases provide a “spatial index” that instead answers the question “which objects are within this particular bounding box?”.

A **bounding box** is the smallest rectangle – parallel to the coordinate axes – capable of containing a given feature.

## Bounding Boxes



Bounding boxes are used because answering the question “is A inside B?” is very computationally intensive for polygons but very fast in the case of rectangles. Even the most complex polygons and linestrings can be represented by a simple bounding box.

Indexes have to perform quickly in order to be useful. So instead of providing exact results, as B-trees do, spatial indexes provide approximate results. The question “what lines are inside this polygon?” will be instead interpreted by a spatial index as “what lines have bounding boxes that are contained inside this polygon’s bounding box?”

The actual spatial indexes implemented by various databases vary widely. The most common implementation is the [R-tree](#) (used in PostGIS), but there are also [Quadtrees](#), and [grid-based indexes](#) implemented in other spatial databases.

### 2.1.4 Spatial Functions

For manipulating data during a query, an ordinary database provides **functions** such as concatenating strings, performing hash operations on strings, doing mathematics on numbers, and extracting information from dates. A spatial database provides a complete set of functions for analyzing geometric components, determining spatial relationships, and manipulating geometries. These spatial functions serve as the building block for any spatial project.

The majority of all spatial functions can be grouped into one of the following five categories:

1. **Conversion:** Functions that *convert* between geometries and external data formats.
2. **Management:** Functions that *manage* information about spatial tables and PostGIS administration.
3. **Retrieval:** Functions that *retrieve* properties and measurements of a Geometry.
4. **Comparison:** Functions that *compare* two geometries with respect to their spatial relation.
5. **Generation:** Functions that *generate* new geometries from others.

The list of possible functions is very large, but a common set of functions is defined by the *OGC SFSQL* and implemented (along with additional useful functions) by PostGIS.

## 2.2 What is PostGIS?

PostGIS turns the [PostgreSQL](#) Database Management System into a spatial database by adding support for the three features: spatial types, indexes, and functions. Because it is built on PostgreSQL, PostGIS automatically inherits important “enterprise” features as well as open standards for implementation

### 2.2.1 But what is PostgreSQL?

PostgreSQL is a powerful, object-relational database management system (ORDBMS). It is released under a BSD-style license and is thus free and open source software. As with many other open source programs, PostgreSQL is not controlled by any single company, but has a global community of developers and companies to develop it.

PostgreSQL was designed from the very start with type extension in mind – the ability to add new data types, functions and access methods at run-time. Because of this, the PostGIS extension can be developed by a separate development team, yet still integrate very tightly into the core PostgreSQL database.

### Why choose PostgreSQL?

A common question from people familiar with open source databases is, “Why wasn’t PostGIS built on MySQL?”.

PostgreSQL has:

- Proven reliability and transactional integrity by default (ACID)
- Careful support for SQL standards (full SQL92)
- Pluggable type extension and function extension
- Community-oriented development model
- No limit on column sizes (“TOAST”able tuples) to support big GIS objects
- Generic index structure (GiST) to allow R-Tree index
- Easy to add custom functions

Combined, PostgreSQL provides a very easy development path to add new spatial types. In the proprietary world, only Illustra (now [Informix](#) Universal Server) allows such easy extension. This is no coincidence; Illustra is a proprietary re-working of the original PostgreSQL code base from the 1980’s.

Because the development path for adding types to PostgreSQL was so straightforward, it made sense to start there. When MySQL released basic spatial types in version 4.1, the PostGIS team took a look at their code, and the exercise reinforced the original decision to use PostgreSQL. Because MySQL spatial objects had to be hacked on top of the string type as a special case, the MySQL code was spread over the entire code base. Development of PostGIS 0.1 took under a month. Doing a “MyGIS” 0.1 would have taken a lot longer, and as such, might never have seen the light of day.

### 2.2.2 Why not Shapefiles?

The [shapefile](#) (and other file formats) have been the standard way of storing and interacting with spatial data since GIS software was first written. However, these “flat” files have the following disadvantages:

- **Files require special software to read and write.** SQL is an abstraction for random data access and analysis. Without that abstraction, you will need to write all the access and analysis code yourself.
- **Concurrent users can cause corruption.** While it's possible to write extra code to ensure that multiple writes to the same file do not corrupt the data, by the time you have solved the problem and also solved the associated performance problem, you will have written the better part of a database system. Why not just use a standard database?
- **Complicated questions require complicated software to answer.** Complicated and interesting questions (spatial joins, aggregations, etc) that are expressible in one line of SQL in the database take hundreds of lines of specialized code to answer when programming against files.

Most users of PostGIS are setting up systems where multiple applications will be expected to access the data, so having a standard SQL access method simplifies deployment and development. Some users are working with large data sets; with files, they might be segmented into multiple files, but in a database they can be stored as a single large table.

In summation, the combination of support for multiple users, complex ad hoc queries, and performance on large data sets are what sets spatial databases apart from file-based systems.

### 2.2.3 A brief history of PostGIS

In the May of 2001, Refractions Research released the first version of PostGIS. PostGIS 0.1 had objects, indexes and a handful of functions. The result was a database suitable for storage and retrieval, but not analysis.

As the number of functions increased, the need for an organizing principle became clear. The “Simple Features for SQL” ([SFSQ](#)L) specification from the Open Geospatial Consortium provided such structure with guidelines for function naming and requirements.

With PostGIS support for simple analysis and spatial joins, [Mapserver](#) became the first external application to provide visualization of data in the database.

Over the next several years the number of PostGIS functions grew, but its power remained limited. Many of the most interesting functions (e.g., ST\_Intersects(), ST\_Buffer(), ST\_Union()) were very difficult to code. Writing them from scratch promised years of work.

Fortunately a second project, the “Geometry Engine, Open Source” or [GEOS](#), came along. The GEOS library provides the necessary algorithms for implementing the [SFSQ](#)SQL specification. By linking in GEOS, PostGIS provided complete support for [SFSQ](#)SQL by version 0.8.

As PostGIS data capacity grew, another issue surfaced: the representation used to store geometry proved relatively inefficient. For small objects like points and short lines, the metadata in the representation had as much as a 300% overhead. For performance reasons, it was necessary to put the representation on a diet. By shrinking the metadata header and required dimensions, overhead greatly reduced. In PostGIS 1.0, this new, faster, lightweight representation became the default.

Recent updates of PostGIS have worked on expanding standards compliance, adding support for curve-based geometries and function signatures specified in the ISO [SQL/MM](#) standard. Through a continued focus on performance, PostGIS 1.4 significantly improved the speed of geometry testing routines.

### 2.2.4 Who uses PostGIS?

For a complete list of case studies, see the [PostGIS case studies](#) page.

## Institut Géographique National, France

IGN is the national mapping agency of France, and uses PostGIS to store the high resolution topographic map of the country, “BDUni”. BDUni has more than 100 million features, and is maintained by a staff of over 100 field staff who verify observations and add new mapping to the database daily. The IGN installation uses the database transactional system to ensure consistency during update processes, and a [warm standby system](#) to maintain uptime in the event of a system failure.

## GlobeXplorer

GlobeXplorer is a web-based service providing online access to petabytes of global satellite and aerial imagery. GlobeXplorer uses PostGIS to manage the metadata associated with the imagery catalogue, so queries for imagery first search the PostGIS catalogue to find the location of the relevant images, then pull the images from storage and return them to the client. In building their system, GlobeXplorer tried other spatial databases but eventually settled on PostGIS because of the great combination of price and performance it offers.

### 2.2.5 What applications support PostGIS?

PostGIS has become a widely used spatial database, and the number of third-party programs that support storing and retrieving data using it has increased as well. The [programs that support PostGIS](#) include both open source and proprietary software on both server and desktop systems.

The following table shows a list of some of the software that leverages PostGIS:

Open/Free	Closed/Proprietary
<ul style="list-style-type: none"> <li>• Loading/Extracting           <ul style="list-style-type: none"> <li>– Shp2Pgsql</li> <li>– ogr2ogr</li> <li>– Dxf2PostGIS</li> </ul> </li> <li>• Web-Based           <ul style="list-style-type: none"> <li>– Mapserver</li> <li>– GeoServer (Java-based WFS / WMS -server )</li> <li>– SharpMap SDK - for ASP.NET 2.0</li> <li>– MapGuide Open Source (using FDO)</li> </ul> </li> <li>• Desktop           <ul style="list-style-type: none"> <li>– uDig</li> <li>– QGIS</li> <li>– mezoGIS</li> <li>– OpenJUMP</li> <li>– OpenEV</li> <li>– SharpMap SDK for Microsoft.NET 2.0</li> <li>– <b>ZigGIS</b> for ArcGIS/ objects.NET</li> <li>– GvSIG</li> <li>– GRASS</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Loading/Extracting           <ul style="list-style-type: none"> <li>– Safe FME Desktop Translator/Converter</li> </ul> </li> <li>• Web-Based           <ul style="list-style-type: none"> <li>– Ionic Red Spider (now ERDAS)</li> <li>– Cadcorp GeognoSIS</li> <li>– Iwan Mapserver</li> <li>– MapDotNet Server</li> <li>– MapGuide Enterprise (using FDO)</li> <li>– ESRI ArcGIS Server 9.3+</li> </ul> </li> <li>• Desktop           <ul style="list-style-type: none"> <li>– Cadcorp SIS</li> <li>– Microimages TNTmips GIS</li> <li>– ESRI ArcGIS 9.3+</li> <li>– Manifold</li> <li>– GeoConcept</li> <li>– MapInfo (v10)</li> <li>– AutoCAD Map 3D (using FDO)</li> </ul> </li> </ul>



## SETUP

### 3.1 Installation

PostGIS is an extension for the PostgreSQL database to deal with spatial data, so in order to use it we first need an installation of the database engine. The following instructions will guide you using the PostGIS binary installers available from the official project site: <https://postgis.net/install/>. There are options for multiple platforms.

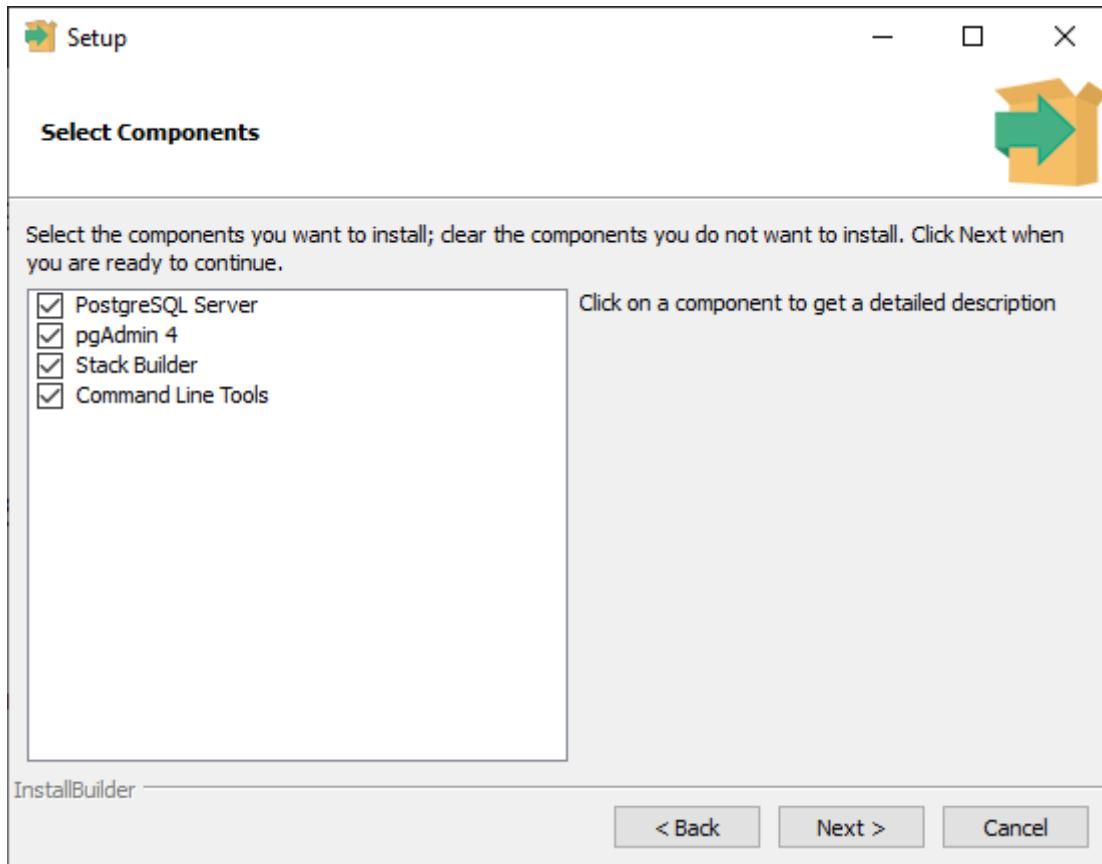
#### 3.1.1 Windows

For Windows, there is the EnterpriseDB PostgreSQL distributions: <https://www.postgresql.org/download/windows/>, an interactive installer that will provide the option for installing PostGIS in an installation dialog with “Stack Builder”:

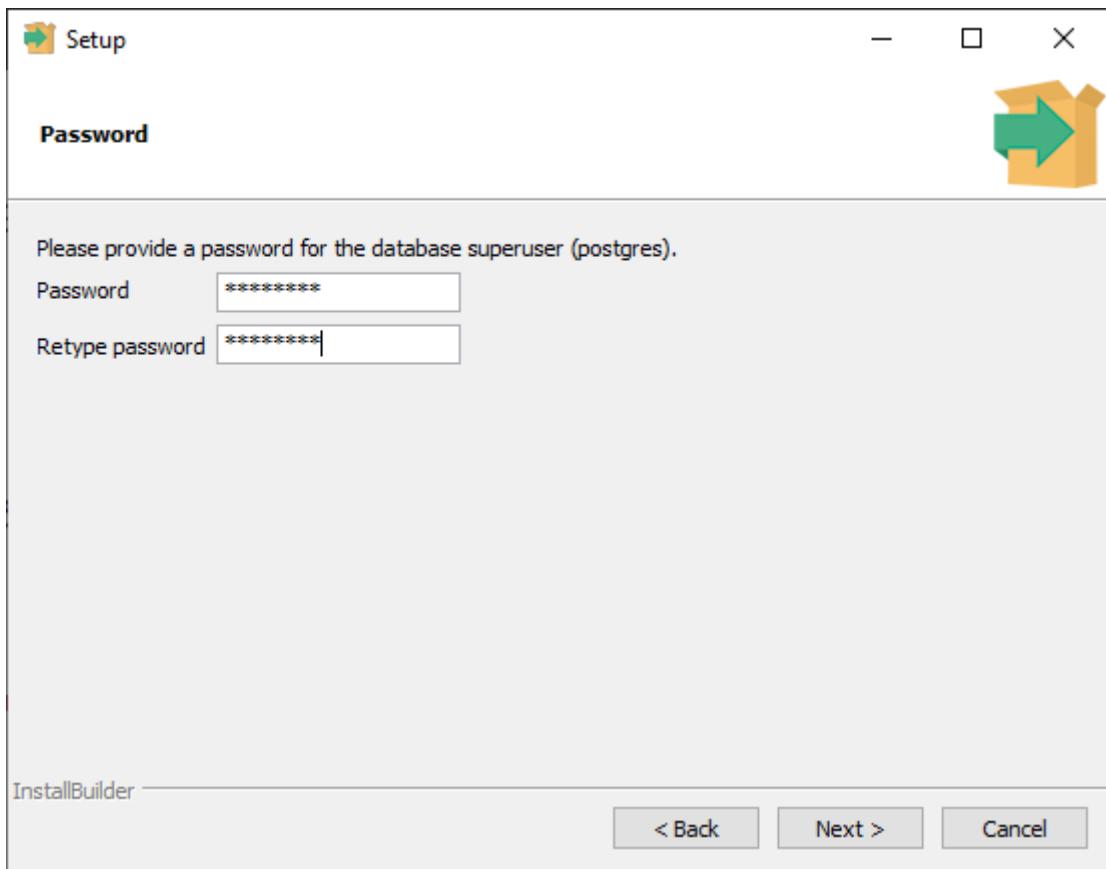
1. First install PostgreSQL with the EDB installer:



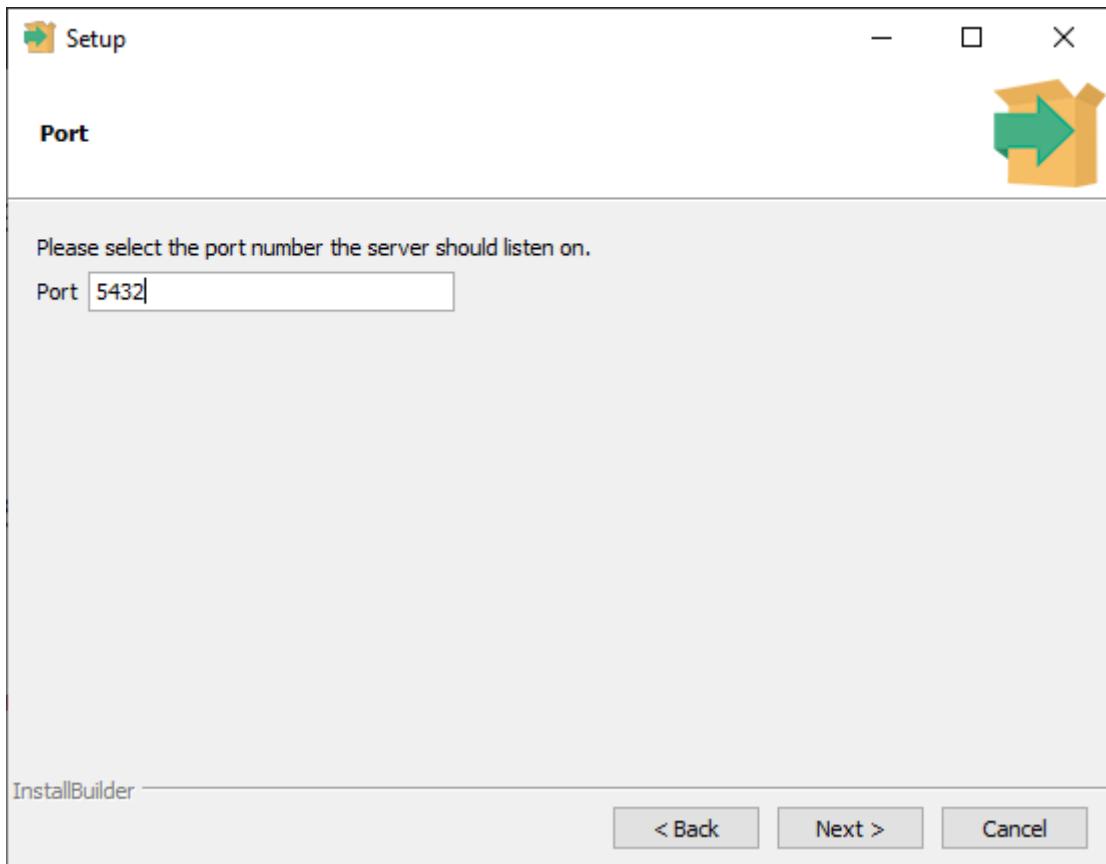
2. Follow the installation process leaving the directories as they are by default and make sure all components (postgreSQL Server, pgAdmin 4, Stack Builder and Command Line Tools) are selected to install.



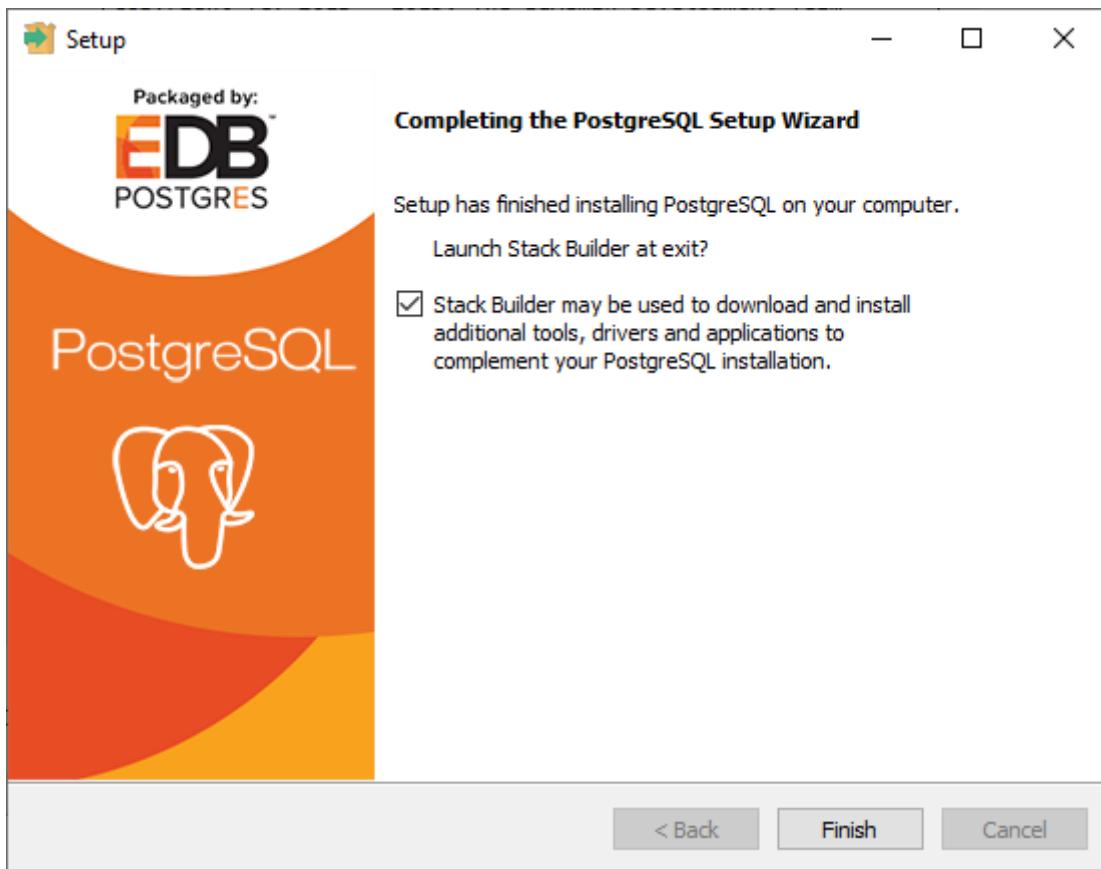
3. Set a password that you will use later for login. You may use `postgres`.



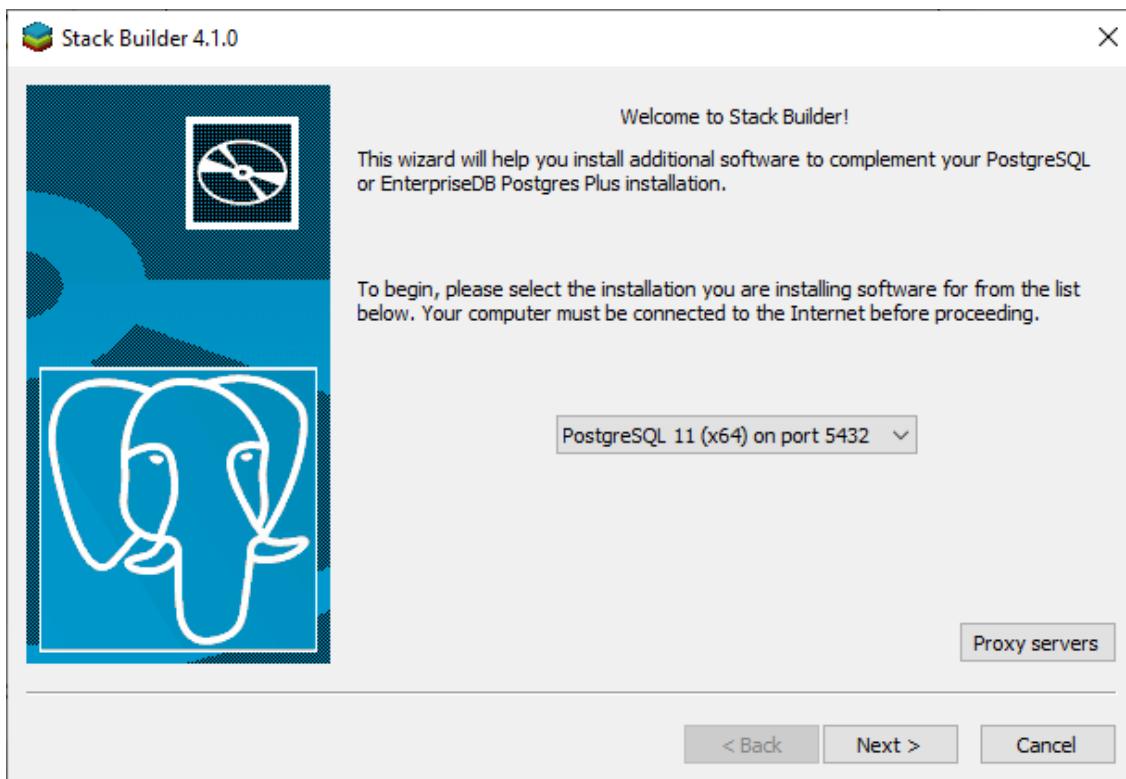
4. Leave the default port to 5432:



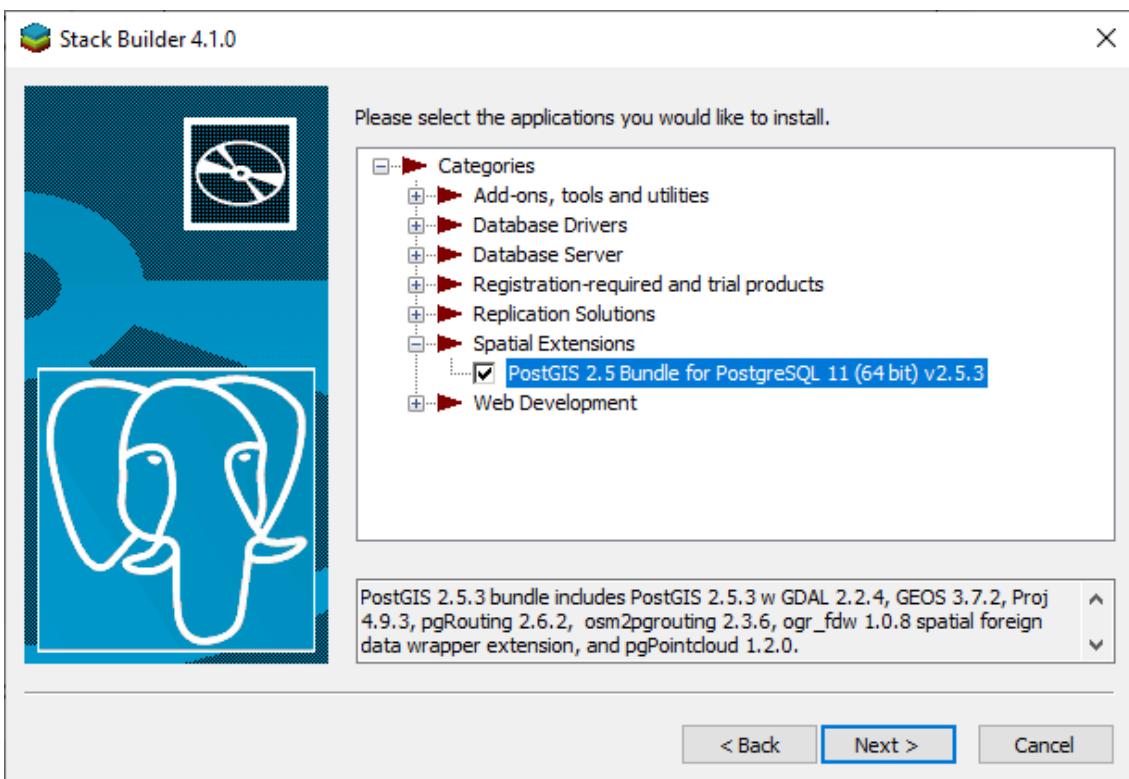
5. At the end of the PostgreSQL installation check the option to launch Stack Builder. This will continue to install PostGIS.



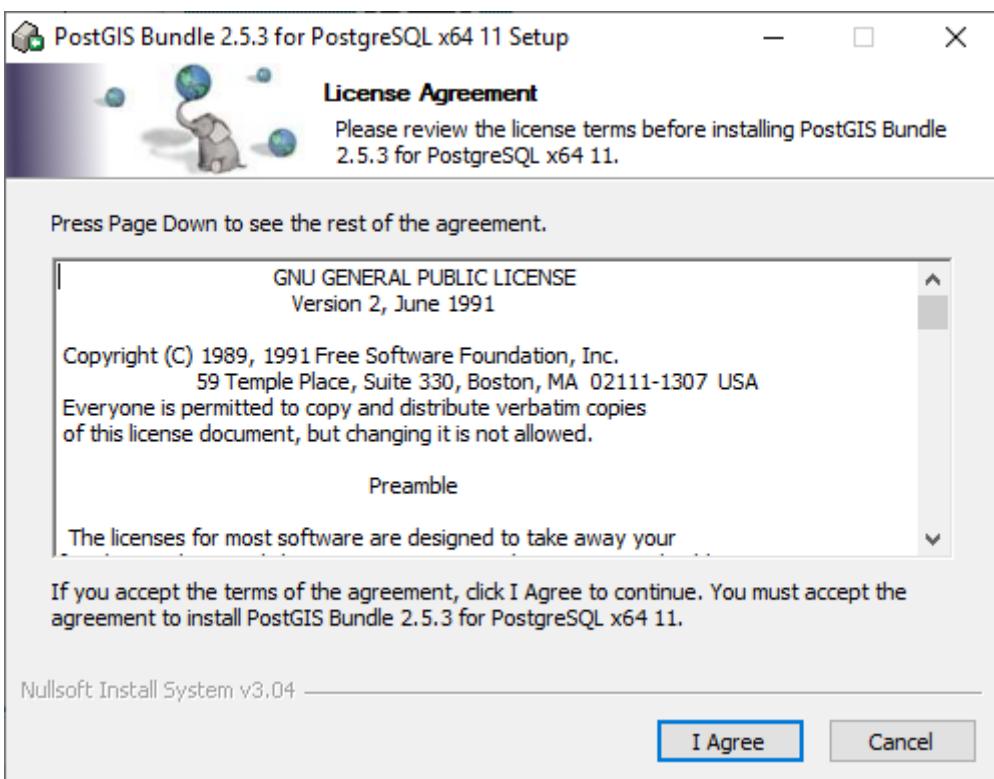
6. Select the recent PostgreSQL installation:



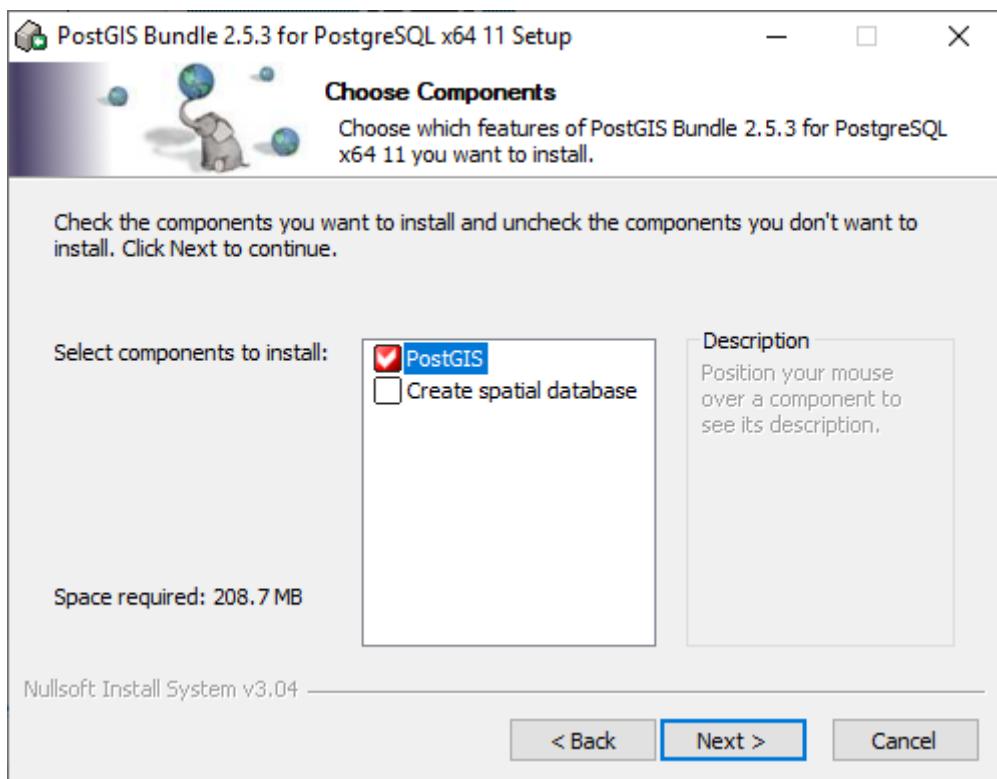
7. Search for PostGIS under **Spatial Extensions** and check it.



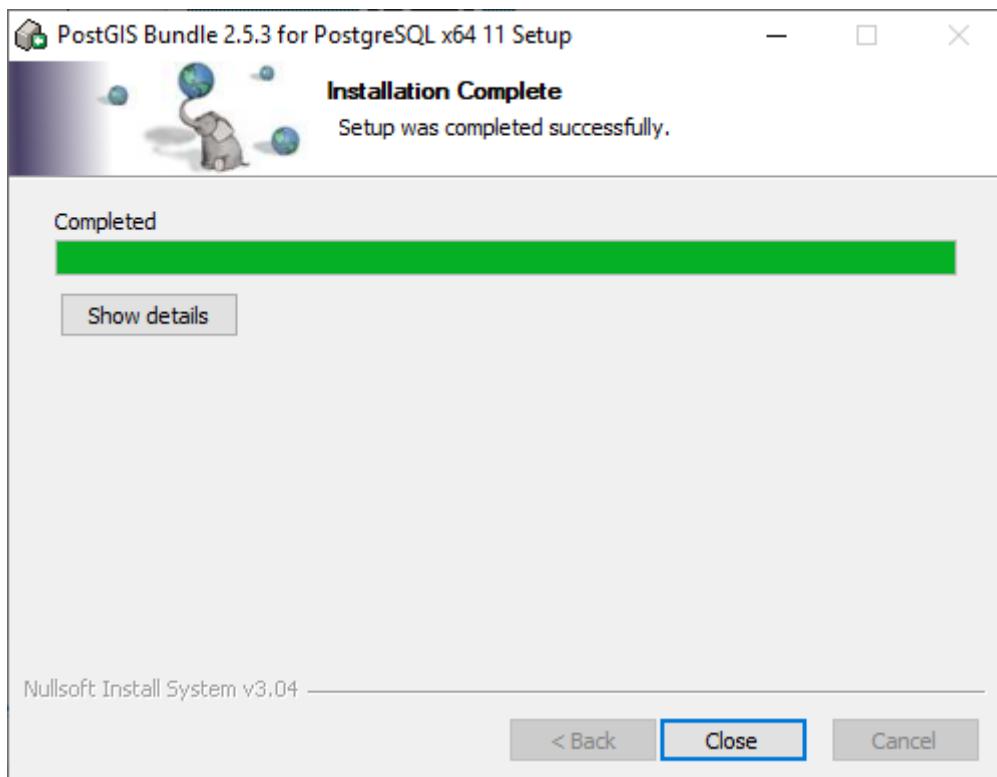
8. Follow the installation process and click on **I Agree** in the License Agreement step:



9. Choose the PostGIS component:



10. Leave the default directories and complete the installation:



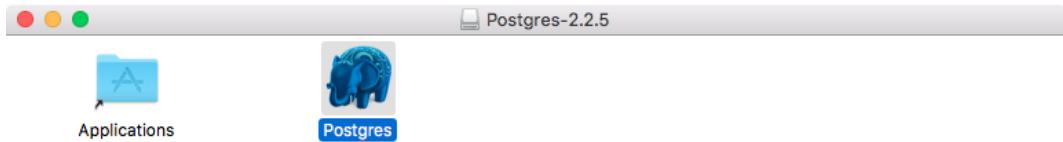
11. Go to **Start** and look for the PostgreSQL folder where pgAdmin 4 will be located. Open it and set a password to login everytime.

### 3.1.2 Mac OS X

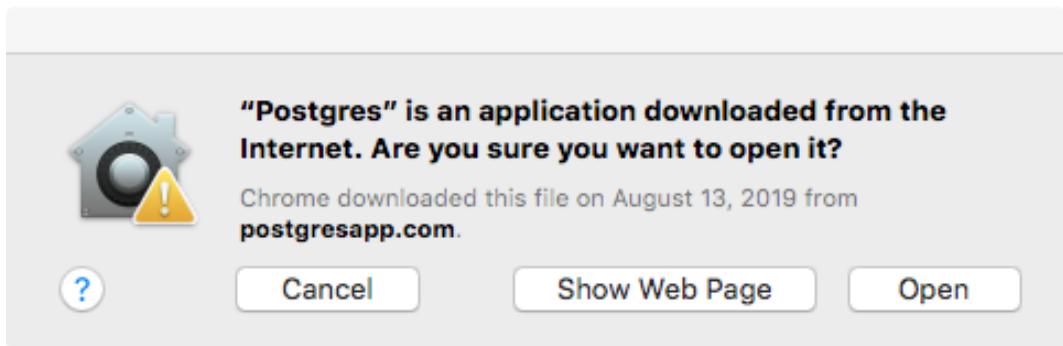
For OSX users it is possible to install PostGIS along with PostgreSQL using the Postgres app: <http://postgresapp.com/>. This is the easiest way to do the installation but you can also use the EnterpriseDb installer: <https://www.enterprisedb.com/downloads/postgres-postgresql-downloads> or the homebrew installer if you are familiar with it by running the command: brew install postgis.

The following steps will guide you in the installation using Postgres.app:

1. Download the dmg file: <https://postgresapp.com/downloads.html>.
2. Drag the Postgres icon into your applications folder.



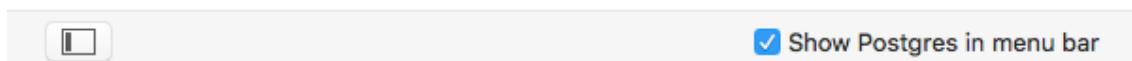
1. Double click the application and allow it to be opened:



1. Click initialize to start the process.

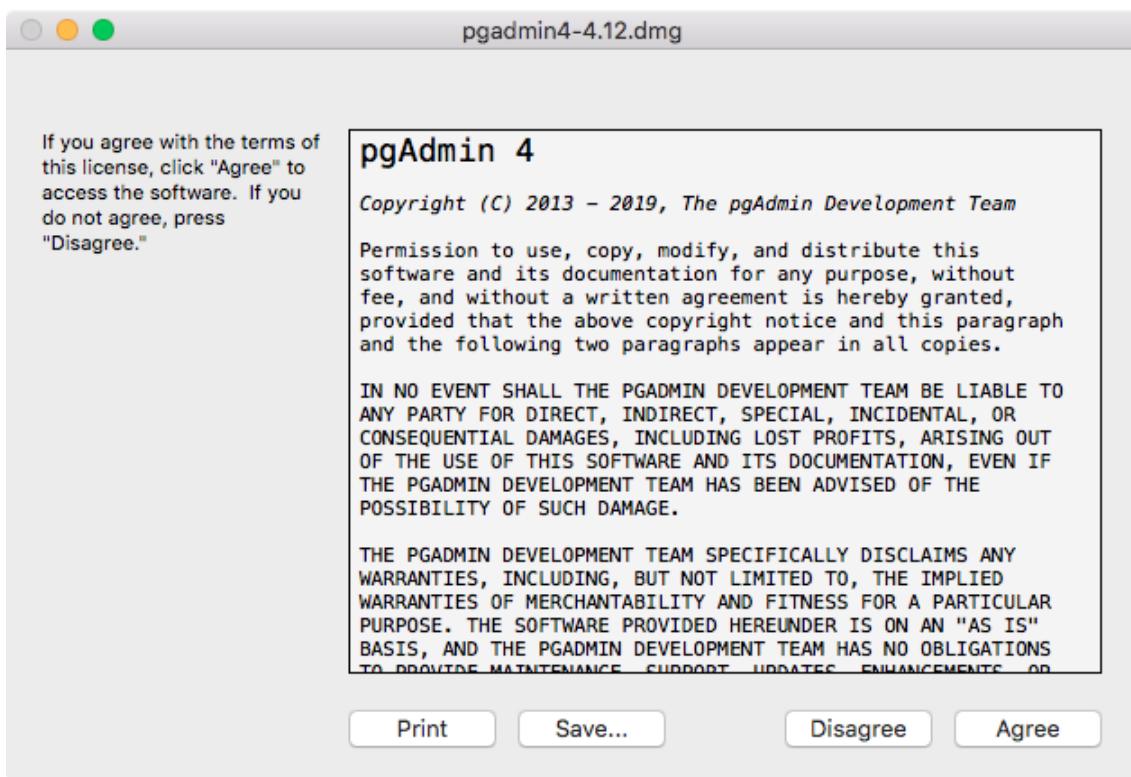


2. If successful, you will see the databases available to connect to, this means that the database engine is up and running. By clicking the icons of the databases you will be prompted to a command line but for this workshop, the user interface pgAdmin 4 will be preferred.

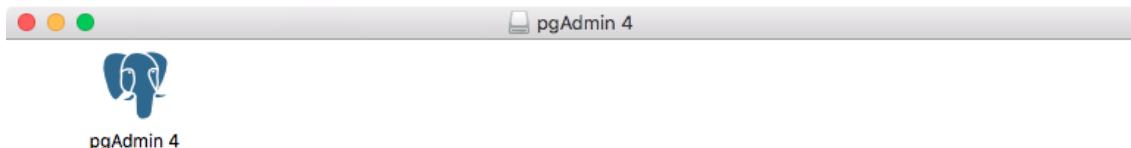


### Install pgAdmin 4 (Mac OS X)

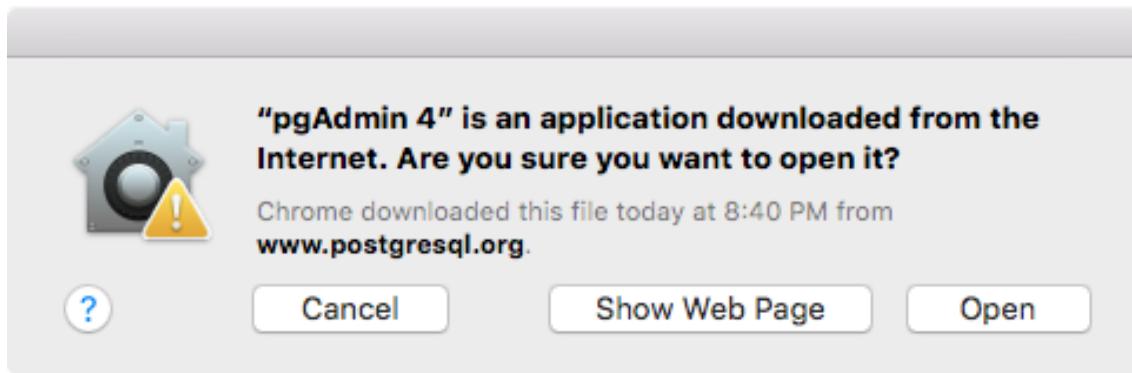
1. Go to: <https://www.pgadmin.org/download/> to get pgAdmin 4 for **Mac OS X**. Double click the installer and agree to the terms by clicking **Agree**.



2. For Mac OS X drag the pgAdmin 4 icon to your applications folder.



3. Open the installed application and allow it to run:



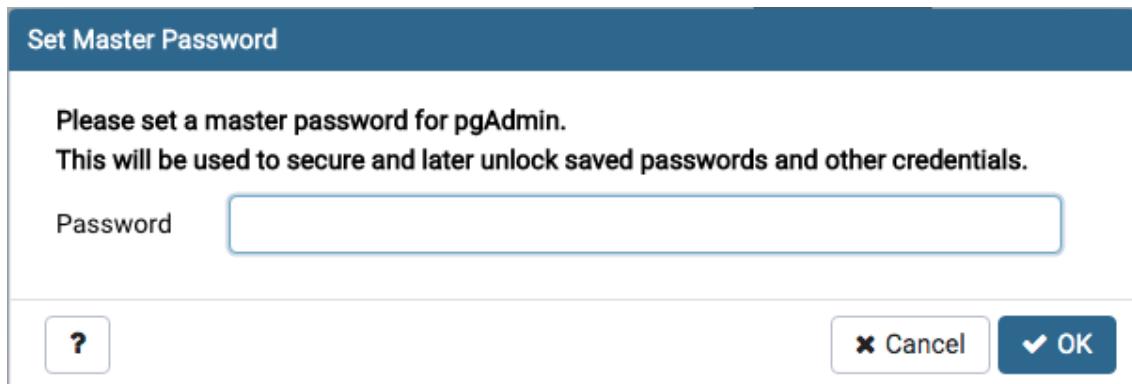
4. pgAdmin 4 is web-based so a tab will open in your browser window. The first time, it will prompt for a master password to use, set one, then enter it to see the servers.

## 3.2 Creating a Spatial Database

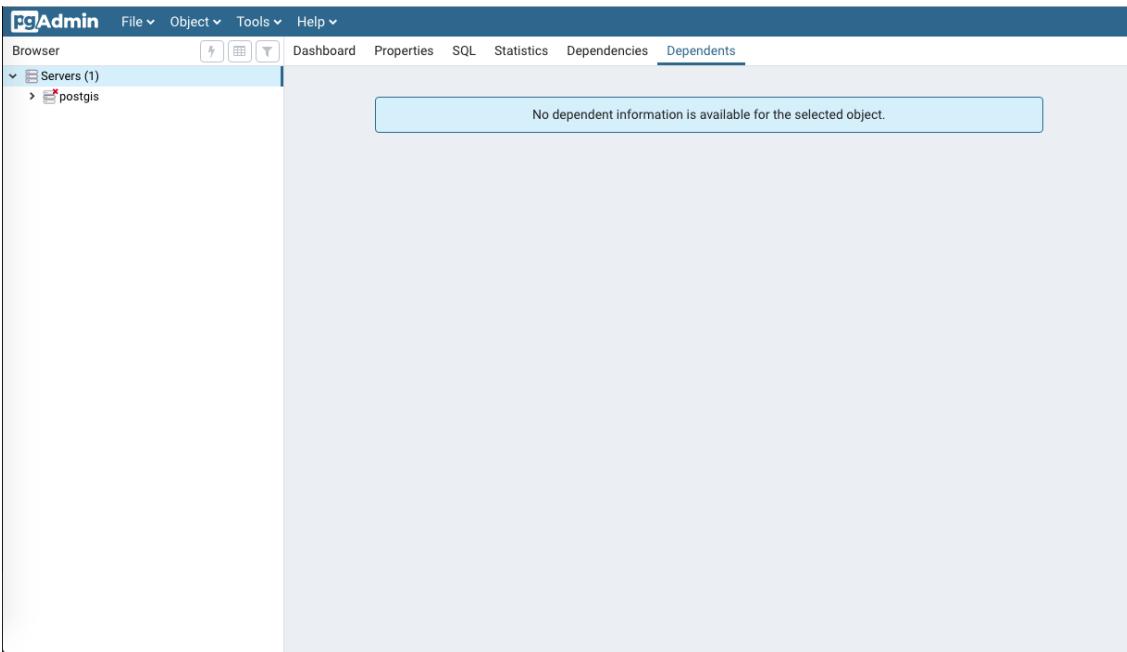
### 3.2.1 PgAdmin

PostgreSQL has a number of administrative front-ends. The primary one is `psql` a command-line tool for entering SQL queries. Another popular PostgreSQL front-end is the free and open source graphical tool `pgAdmin`. All queries done in pgAdmin can also be done on the command line with `psql`.

1. Find pgAdmin and start it up. If it is the first time you are running it, it will ask you to set up a master password, this will allow you to log in to pgAdmin everytime and connect to servers.



2. Afterwards it will show you a list of servers. Click on **Servers** to view the active servers.

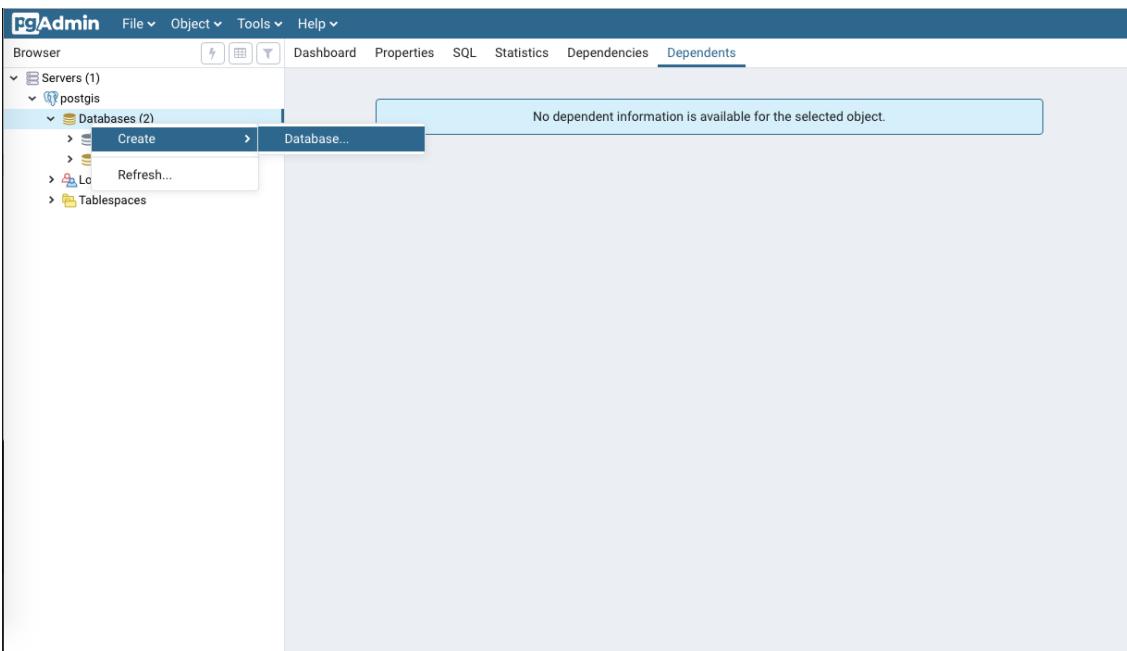


3. If this is the first time you have run pgAdmin, you should have a server entry for **PostGIS (localhost:5432)** already configured in pgAdmin. Click to display the databases in it.

The PostGIS database has been installed with unrestricted access for local users (users connecting from the same machine as the database is running). That means that it will accept *any* password you provide. If you need to connect from a remote computer, the password for the `postgres` user has been set to `postgres`.

### 3.2.2 Creating a Database

1. Open the Databases tree item and have a look at the available databases. The `postgres` database is the user database for the default `postgres` user and is not too interesting to [yellow speech bubble].
2. Right-click on the Databases item and select New Database.



3. Fill in the New Database form as shown below and click OK.

Name	nyc
Owner	postgres

**Create - Database**

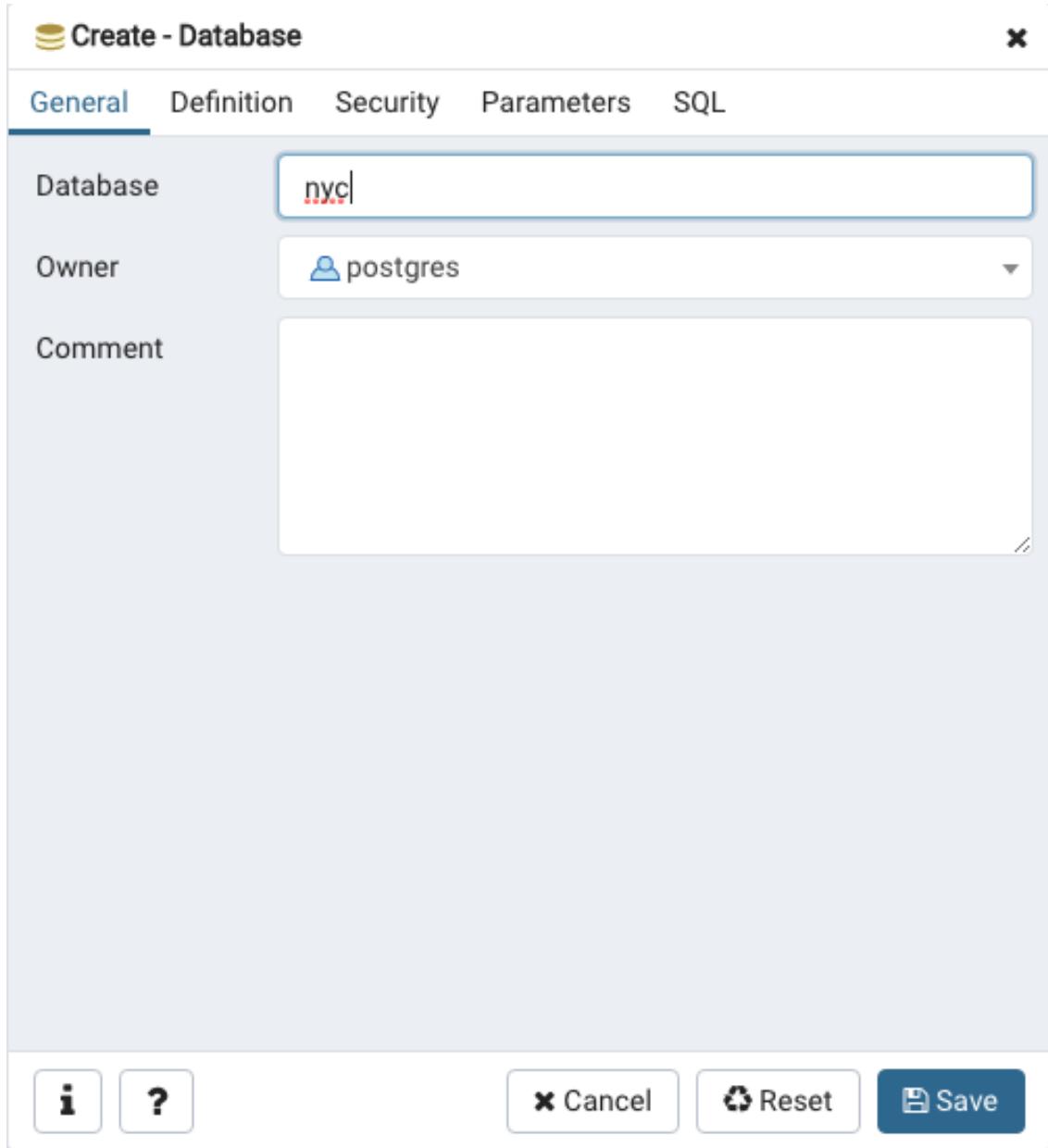
**General**    Definition    Security    Parameters    SQL

Database

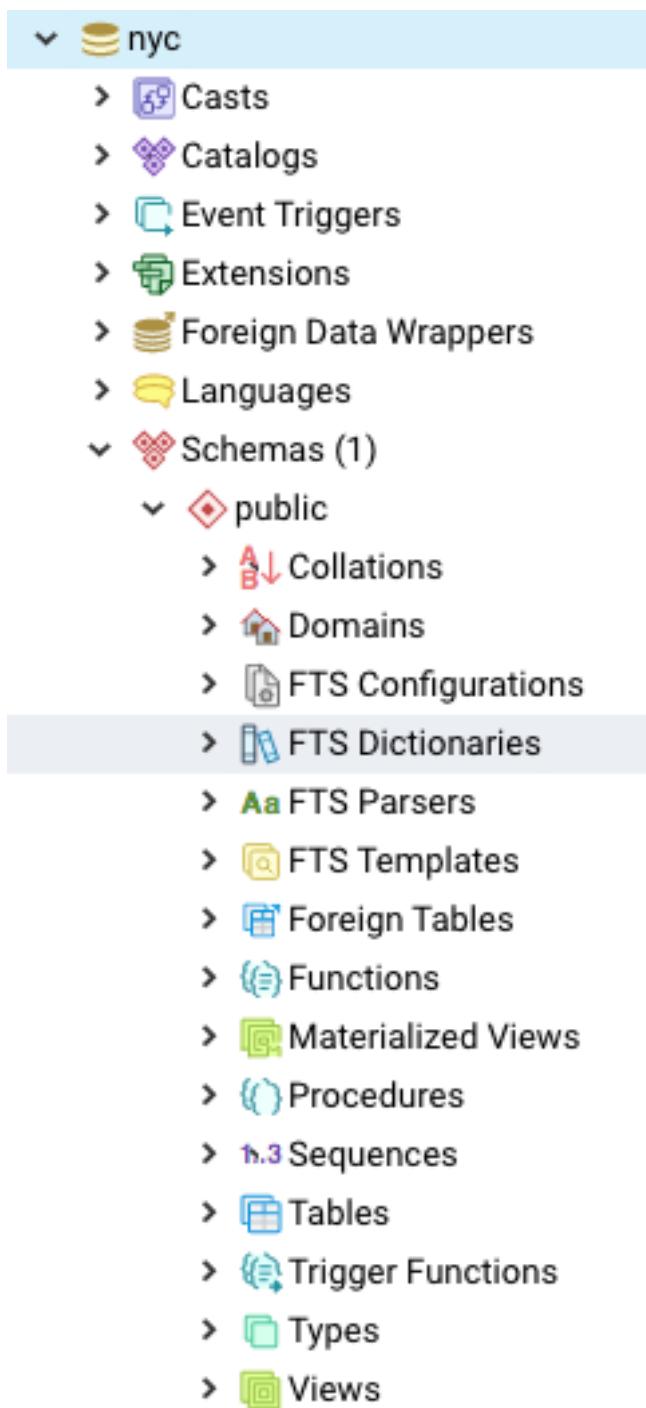
Owner

Comment

**i** **?**    **Cancel** **Reset** **Save**



4. Select the new nyc database and open it up to display the tree of objects. You will be the public schema.



5. Click on the SQL query button indicated below (or go to *Tools > Query Tool*).



6. Enter the following query into the query text field to load the PostGIS spatial extension:

```
CREATE EXTENSION postgis;
```

7. Click the **Execute** button in the toolbar (or press **F5**) to “Execute the query.”

8. Now confirm that PostGIS is installed by running a PostGIS function:

```
SELECT postgis_full_version();
```

You have successfully created a PostGIS spatial database!!

### 3.2.3 Function List

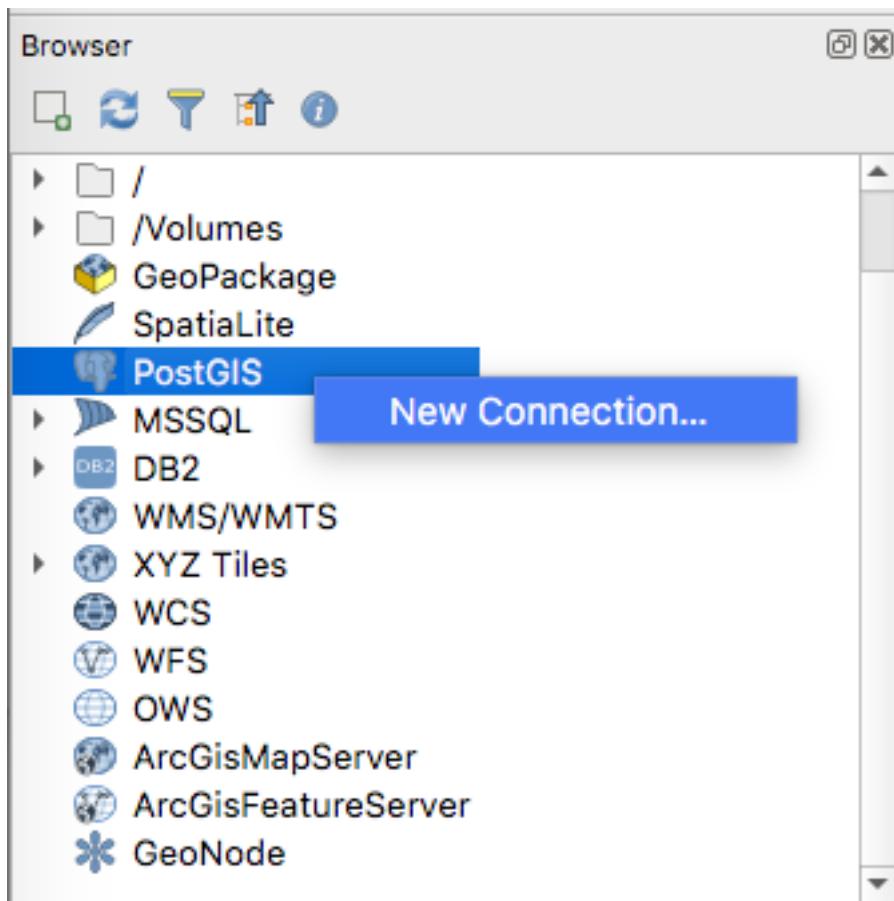
`PostGIS_Full_Version`: Reports full PostGIS version and build configuration info.

## 3.3 Loading spatial data

Supported by a wide variety of libraries and applications, PostGIS provides many options for loading data. This section will focus on the basics – loading shapefiles using QGIS DB Manager.

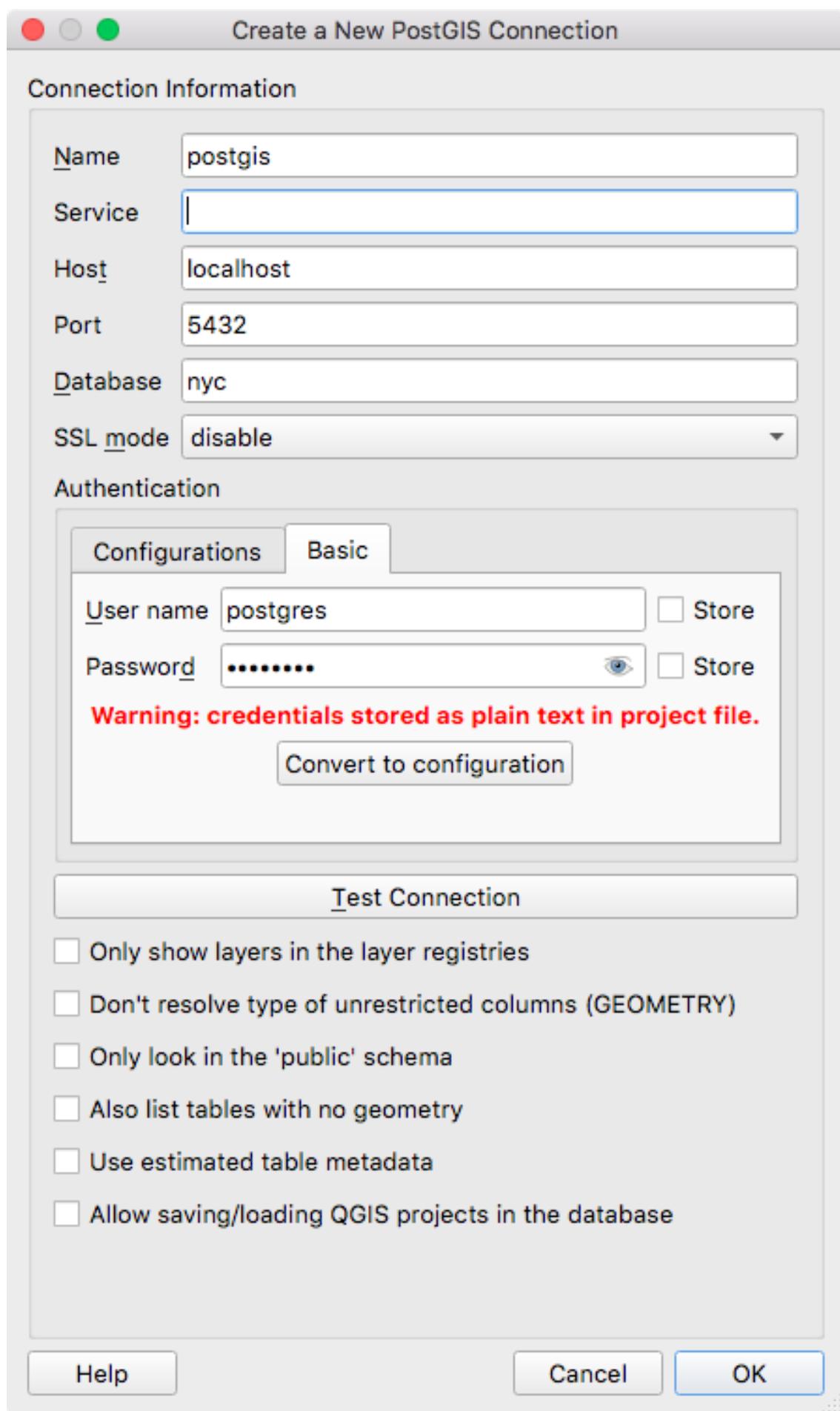
[QGIS](#) is a free and open source geographic information system to create, edit, visualize, analyse and publish geospatial information. It's available for multiple platforms: Windows, Mac, Linux, BSD. To get it, download the long term release (LTR) from the project website: <https://qgis.org/en/site/> forusers/download.html and follow the installation process. This tutorial was made using the version 3.4 (Madeira).

1. First, launch the QGIS software and navigate to the browser panel where you will find the PostGIS icon. Right click on it and select **New Connection...**

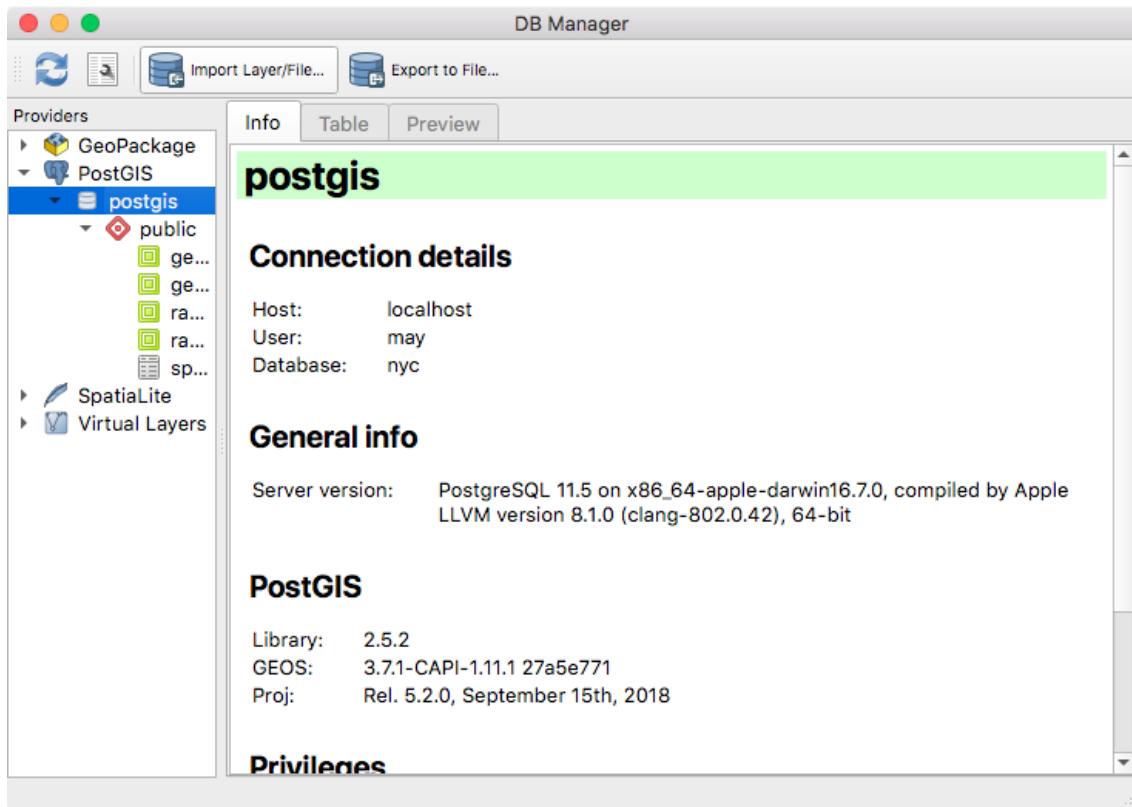


2. Fill in the connection details for the **Authentication** section and click on the **Test Connection** button. This will test the connection and report back if it works or not. If it is working, click the **Ok** button to create the connection.

<b>Username</b>	postgres
<b>Password</b>	postgres
<b>Server Host</b>	localhost 5432
<b>Database</b>	nyc

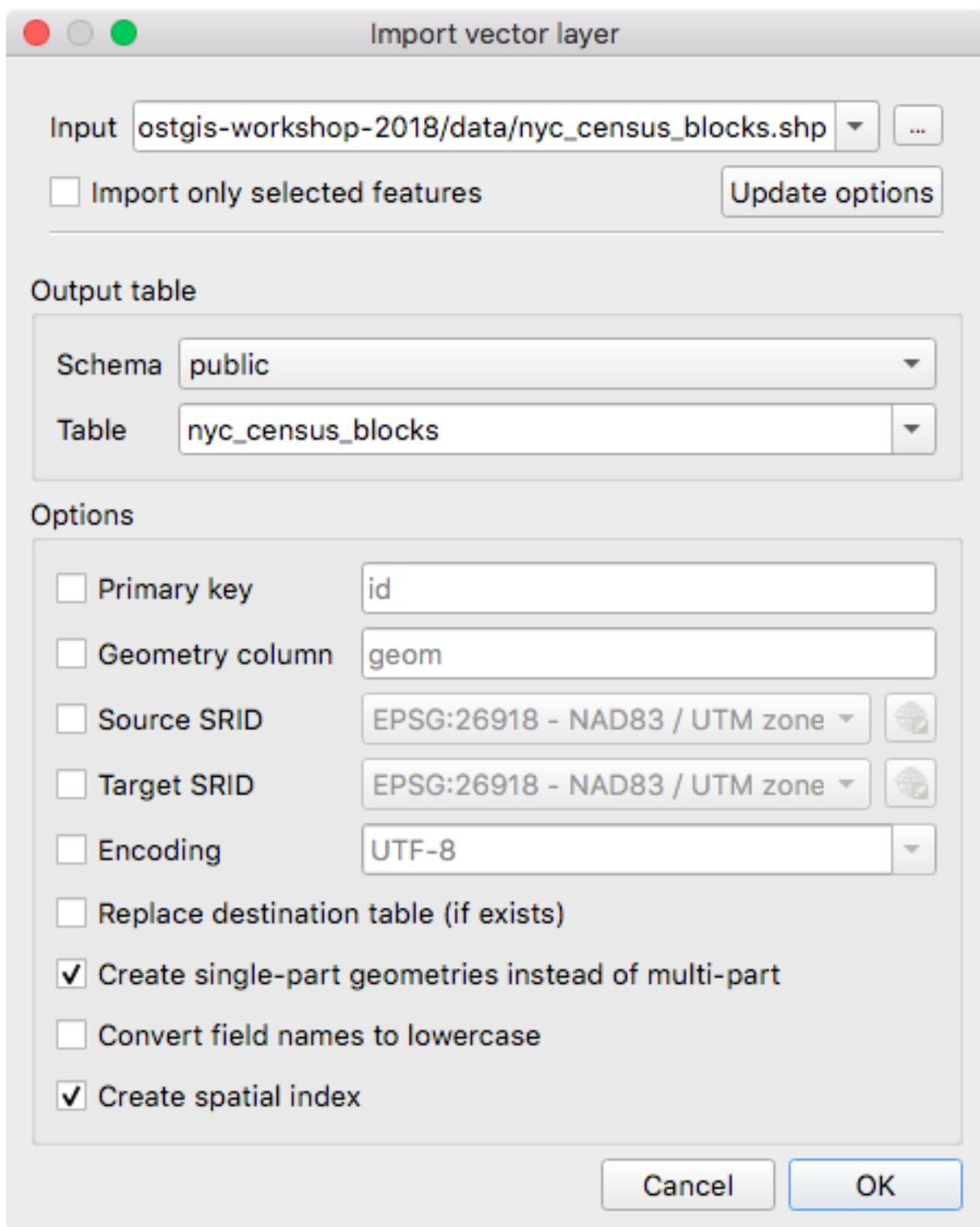


3. Next, open the **DB Manager** browser from the menu **Database > DB Manager...** and select the postgis connection added in the previous step.



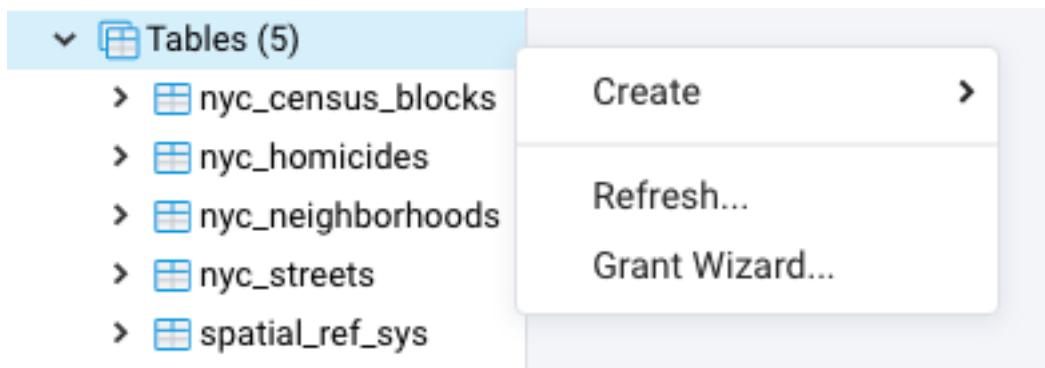
4. Then **Click on Import Layer/File** and navigate to the folder from the data bundle:`\postgis-workshop-2018\data`. Select the `nyc_census_blocks.shp` file and tick the checkboxes next to **Create single-part geometries instead of multi-part** and **Create spatial index**. Make sure the SRID value for the file is set to **26918** and the schema is set to public. Note that the primary key, geometry column, source SRID, target SRID and encoding are already filled in using the shapefile, but you can optionally change them (**Don't!** There are steps later in the workshop that expect the default names.)

<b>Destination Schema</b>	public
<b>SRID</b>	26918
<b>Destination Table</b>	nyc_census_blocks
<b>Geometry Column</b>	geom
<b>Create single-part geometries instead of multi-part</b>	true
<b>Create spatial index</b>	true



5. Click the **Ok** button to te the table. It may take a few minutes to load, but this is the largest file in our **test** set.
6. Repeat the import process for the remaining shapefiles in the data directory:
  - nyc\_streets.shp
  - nyc\_neighborhoods.shp
  - nyc\_subway\_stations.shp
  - nyc\_homicides.shp
7. When all the files are loaded, click the “Refresh” button in pgAdmin to update the tree view. You

should see your four tables show up in the **Databases > nyc > Schemas > public > Tables** section of the tree.



### 3.3.1 Shapefiles? What's that?

You may be asking yourself – “What’s this shapefile thing?” A “shapefile” commonly refers to a collection of files with .shp, .shx, .dbf, and other extensions on a common prefix name (e.g., nyc\_census\_blocks). The actual shapefile relates specifically to files with the .shp extension. However, the .shp file alone is incomplete for distribution without the required supporting files.

Mandatory files:

- .shp—shape format; the feature geometry itself
- .shx—shape index format; a positional index of the feature geometry
- .dbf —attribute format; columnar attributes for each shape, in dBase III

Optional files include:

- .prj—projection format; the coordinate system and projection information, a plain text file describing the projection using well-known text format

The Db Manager importer makes shape data usable in PostGIS by converting it from binary data into a series of SQL commands that are then run in the database to load the data.

### 3.3.2 SRID 26918? What's with that?

Most of the import process is self-explanatory, but even experienced GIS professionals can trip over an **SRID**.

An “SRID” stands for “Spatial Reference IDentifier.” It defines all the parameters of our data’s geographic coordinate system and projection. An SRID is convenient because it packs all the information about a map projection (which can be quite complex) into a single number.

You can see the definition of our workshop map projection by looking it up either in an online database,

- <http://spatialreference.org/ref/epsg/26918/>

or directly inside PostGIS with a query to the `spatial_ref_sys` table.

```
SELECT srtext FROM spatial_ref_sys WHERE srid = 26918;
```

---

**Note:** The PostGIS `spatial_ref_sys` table is an *OGC*-standard table that defines all the spatial reference systems known to the database. The data shipped with PostGIS, lists over 3000 known spatial

reference systems and details needed to transform/re-project between them.

---

In both cases, you see a textual representation of the **26918** spatial reference system (pretty-printed here for clarity):

```
PROJCS["NAD83 / UTM zone 18N",
    GEOGCS["NAD83",
        DATUM["North_American_Datum_1983",
            SPHEROID["GRS 1980", 6378137, 298.257222101, AUTHORITY["EPSG", "7019"]],
            AUTHORITY["EPSG", "6269"]],
        PRIMEM["Greenwich", 0, AUTHORITY["EPSG", "8901"]],
        UNIT["degree", 0.01745329251994328, AUTHORITY["EPSG", "9122"]],
        AUTHORITY["EPSG", "4269"]],
    UNIT["metre", 1, AUTHORITY["EPSG", "9001"]],
    PROJECTION["Transverse_Mercator"],
    PARAMETER["latitude_of_origin", 0],
    PARAMETER["central_meridian", -75],
    PARAMETER["scale_factor", 0.9996],
    PARAMETER["false_easting", 500000],
    PARAMETER["false_northing", 0],
    AUTHORITY["EPSG", "26918"],
    AXIS["Easting", EAST],
    AXIS["Northing", NORTH]]
```

If you open up the `nyc_neighborhoods.prj` file from the data directory, you'll see the same projection definition.

A common problem for people getting started with PostGIS is figuring out what SRID number to use for their data. All they have is a `.prj` file. But how do humans translate a `.prj` file into the correct SRID number?

The easy answer is to use a computer. Plug the contents of the `.prj` file into <http://prj2epsg.org>. This will give you the number (or a list of numbers) that most closely match your projection definition. There aren't numbers for *every* map projection in the world, but most common ones are contained within the `prj2epsg` database of standard numbers.

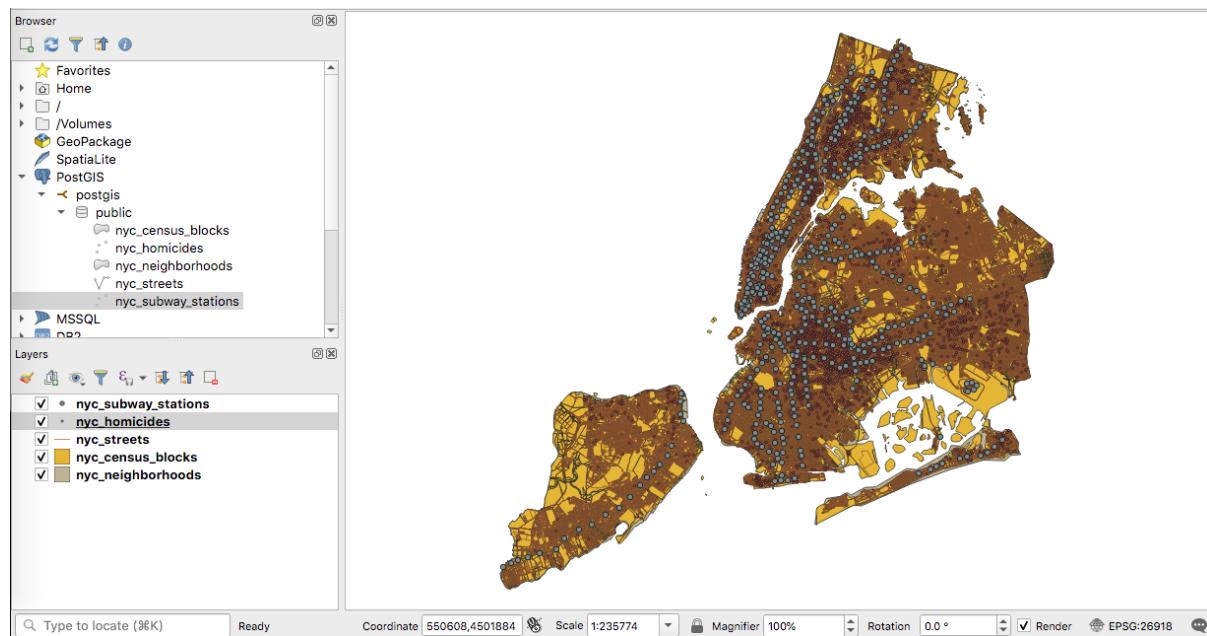
The screenshot shows the Prj2EPSG service interface. At the top, it says "Prj2EPSG" and "Boundless". Below that, a message states: "Prj2EPSG is a simple service for converting well-known text projection information from .prj files into standard EPSG codes." A text input field contains a WKT definition for NAD83 / UTM zone 18N. Below the input field is a "Choose File" button with "No file chosen" next to it. A "Convert" button is present. The "Results" section shows a single exact match: "26918 - NAD83 / UTM zone 18N".

Data you receive from local agencies—such as New York City—will usually be in a local projection noted by “state plane” or “UTM”. Our projection is “Universal Transverse Mercator (UTM) Zone 18 North” or EPSG:26918.

### 3.3.3 Things to Try: View data using QGIS

You can use QGIS for quickly looking at data too. You can view a number of data formats including flat shapefiles and a PostGIS database. Its graphical interface allows for easy exploration of your data, as well as simple testing and fast styling.

To view the imported data, click on the connection in the left panel and **Double-click** on the tables created to add them to the map.



### 3.4 About our data

The data for this workshop is four shapefiles for New York City, and one attribute table of sociodemographic variables. We've loaded our shapefiles as PostGIS tables and will add sociodemographic data later in the workshop.

The following describes the number of records and table attributes for each of our datasets. These attribute values and relationships are fundamental to our future analysis.

To explore the nature of your tables in pgAdmin, right-click a highlighted table and select **Properties**. You will find a summary of table properties, including a list of table attributes within the **Columns** tab.

#### 3.4.1 nyc\_census\_blocks

A census block is the smallest geography for which census data is reported. All higher level census geographies (block groups, tracts, metro areas, counties, etc) can be built from unions of census blocks. We have attached some demographic data to our collection of blocks.

Number of records: 36592

<b>blkid</b>	A 15-digit code that uniquely identifies every census <b>block</b> . Eg: 360050001009000
<b>popn_total</b>	Total number of people in the census block
<b>popn_white</b>	Number of people self-identifying as “White” in the block
<b>popn_black</b>	Number of people self-identifying as “Black” in the block
<b>popn_nativ</b>	Number of people self-identifying as “Native American” in the block
<b>popn_asian</b>	Number of people self-identifying as “Asian” in the block
<b>popn_other</b>	Number of people self-identifying with other categories in the block
<b>boroname</b>	Name of the New York borough. Manhattan, The Bronx, Brooklyn, Staten Island, Queens
<b>geom</b>	Polygon boundary of the block

---

**Note:** To get census data into GIS, you need to join two pieces of information: the actual data (text),

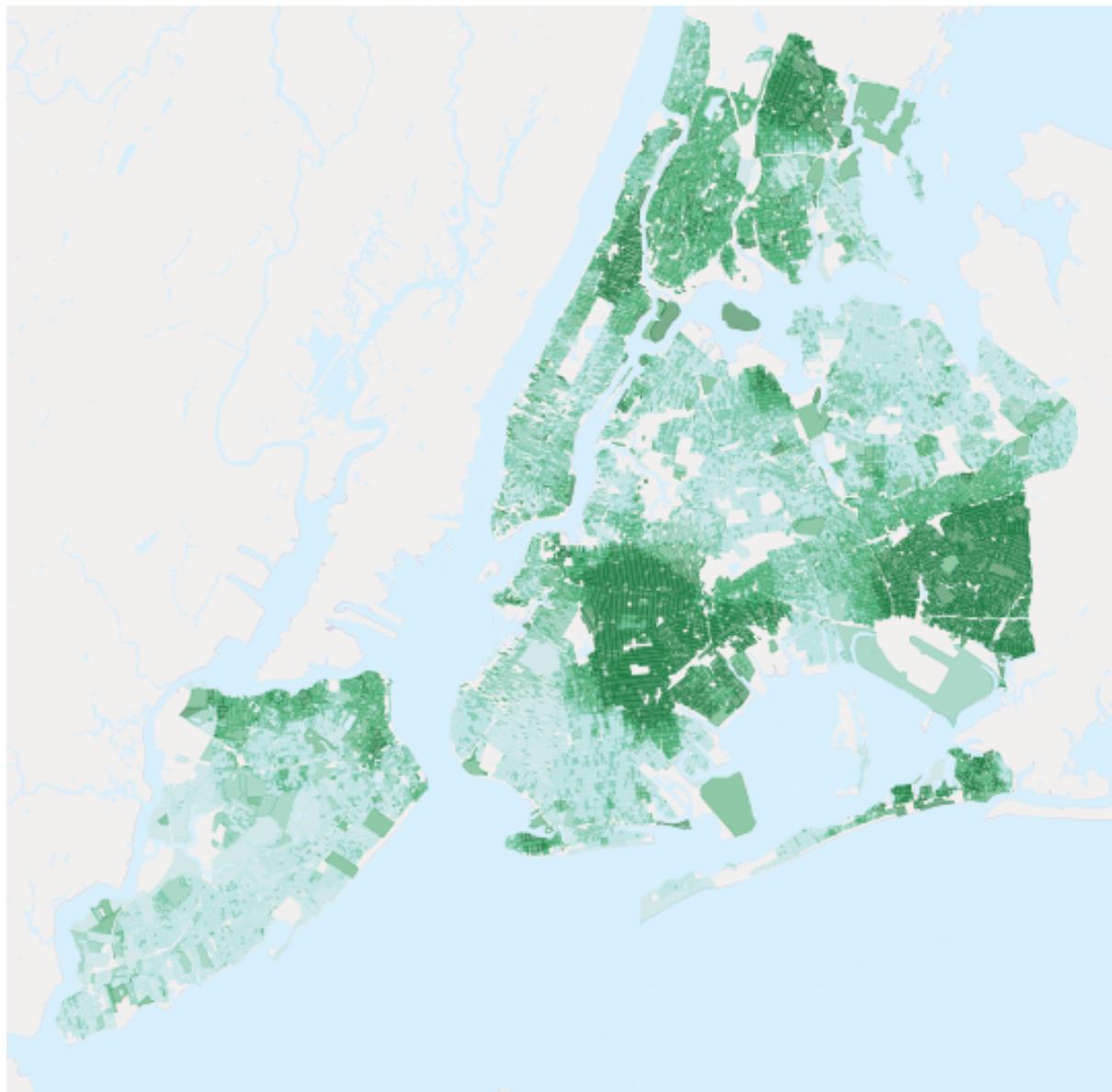


Fig. 3.1: *Black population as a percentage of Total Population*

and the boundary files (spatial). There are many options for getting the data, including downloading data and boundaries from the Census Bureau's [American FactFinder](#).

---

### 3.4.2 nyc\_neighborhoods

New York has a rich history of neighborhood names and extent. Neighborhoods are social constructs that do not follow lines laid down by the government. For example, the Brooklyn neighborhoods of Carroll Gardens, Red Hook, and Cobble Hill were once collectively known as "South Brooklyn." And now, depending on which real estate agent you talk to, the same four blocks in the-neighborhood-formerly-known-as-Red-Hook can be referred to as Columbia Heights, Carroll Gardens West, or Red Hook!

Number of records: 129

<b>name</b>	Name of the neighborhood
<b>boroname</b>	Name of the New York borough. Manhattan, The Bronx, Brooklyn, Staten Island, Queens
<b>geom</b>	Polygon boundary of the neighborhood

### 3.4.3 nyc\_streets

The street centerlines form the transportation network of the city. These streets have been flagged with types in order to distinguish between such thoroughfares as back alleys, arterial streets, freeways, and smaller streets. Desirable areas to live might be on residential streets rather than next to a freeway.

Number of records: 19091

<b>name</b>	Name of the street
<b>oneway</b>	Is the street one-way? "yes" = yes, "" = no
<b>type</b>	Road type (primary, secondary, residential, motorway)
<b>geom</b>	Linear centerline of the street

### 3.4.4 nyc\_subway\_stations

The subway stations link the upper world where people live to the invisible network of subways beneath. As portals to the public transportation system, station locations help determine how easy it is for different people to enter the subway system.

Number of records: 491

<b>name</b>	Name of the station
<b>borough</b>	Name of the New York borough. Manhattan, The Bronx, Brooklyn, Staten Island, Queens
<b>routes</b>	Subway lines that run through this station
<b>transfers</b>	Lines you can transfer to via this station
<b>express</b>	Stations where express trains stop, "express" = yes, "" = no
<b>geom</b>	Point location of the station

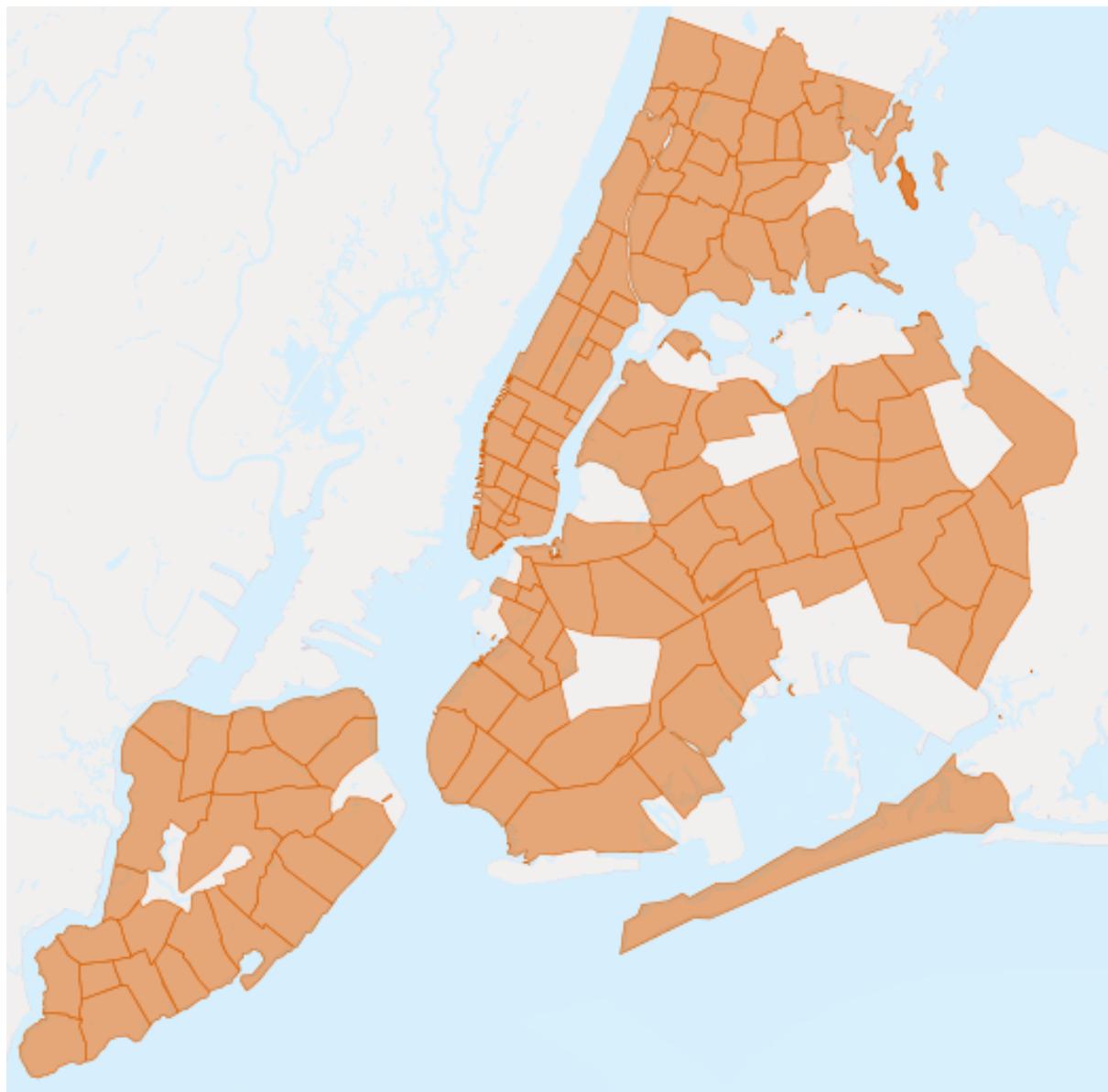


Fig. 3.2: *The neighborhoods of New York City*

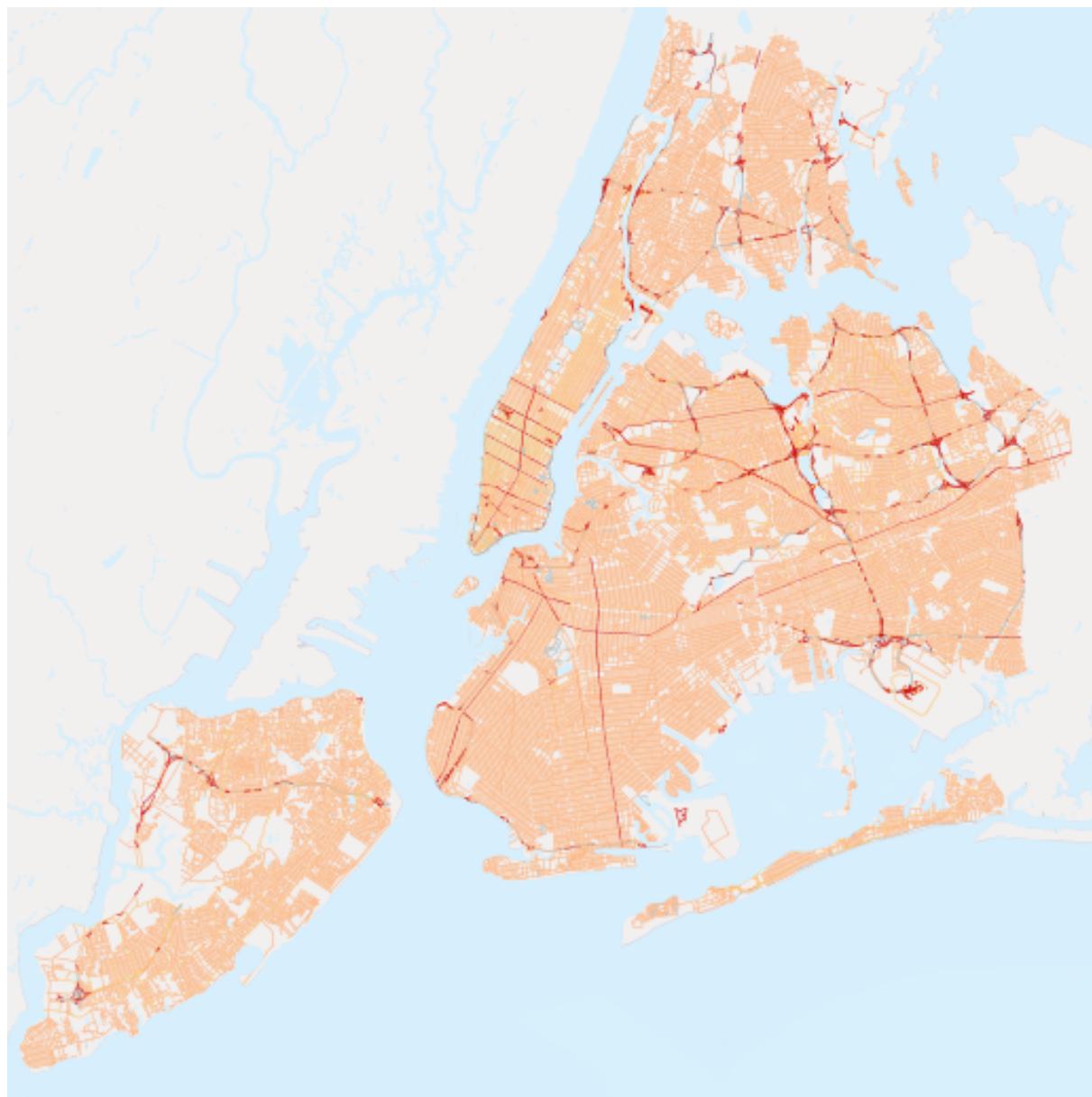


Fig. 3.3: *The streets of New York City. Major roads are in red.*



Fig. 3.4: *Point locations for New York City subway stations*

### 3.4.5 nyc\_census\_sociodata

There is a rich collection of social-economic data collected during the census process, but only at the larger geography level of census tract. Census blocks combine to form census tracts (and block groups). We have collected some social-economic at a census tract level to answer some of these more interesting questions about New York City.

---

**Note:** The `nyc_census_sociodata` is a data table. We will need to connect it to Census geographies before conducting any spatial analysis.

---

<b>tractid</b>	An 11-digit code that uniquely identifies every census <b>tract</b> . (“36005000100”)
<b>transit_total</b>	Number of workers in the tract
<b>transit_private</b>	Number of workers in the tract who use private automobiles / motorcycles
<b>transit_public</b>	Number of workers in the tract who take public transit
<b>transit_walk</b>	Number of workers in the tract who walk
<b>transit_other</b>	Number of workers in the tract who use other forms like walking / biking
<b>transit_none</b>	Number of workers in the tract who work from home
<b>transit_time_mins</b>	Total number of minutes spent in transit by all workers in the tract (minutes)
<b>family_count</b>	Number of families in the tract
<b>family_income_median</b>	Median family income in the tract (dollars)
<b>family_income_mean</b>	Average family income in the tract (dollars)
<b>family_income_aggregate</b>	Total income of all families in the tract (dollars)
<b>edu_total</b>	Number of people with educational history
<b>edu_no_highschool_dipl</b>	Number of people with no high school diploma
<b>edu_highschool_dipl</b>	Number of people with high school diploma and no further education
<b>edu_college_dipl</b>	Number of people with college diploma and no further education
<b>edu_graduate_dipl</b>	Number of people with graduate school diploma

---

**CHAPTER  
FOUR**

---

**BASIC**

## 4.1 Simple SQL

*SQL*, or “Structured Query Language”, is a means of asking questions of, and updating data in, relational databases. You have already seen SQL when we created our first database. Recall:

```
SELECT postgis_full_version();
```

But that was a question about the database. Now that we’ve loaded data into our database, let’s use SQL to ask questions of the data! For example,

“What are the names of all the neighborhoods in New York City?”

Open up the SQL query window in pgAdmin by clicking the SQL button



then enter the following query in to the query window

```
SELECT name FROM nyc_neighborhoods;
```

and click the **Execute Query** button (again the thunder).



The query will run for a few (milli)seconds and return the 129 results.

The screenshot shows a PostgreSQL query editor interface. At the top, there are tabs for "Query Editor" (which is selected) and "Query History". To the right of the editor area is a "Scratch Pad" section with a close button. Below the editor area, there are tabs for "Data Output" (selected), "Explain", "Messages", and "Notifications". The main content area displays the results of a query:

```
1  SELECT name FROM nyc_neighborhoods;
```

	name
1	Bensonhurst
2	East Village
3	West Village
4	Throggs Neck
5	Wakefield-Williamsbridge
6	Auburndale
7	Battery Park
8	Carnegie Hill

But what exactly happened here? To understand, let's begin with the four “verbs” of SQL,

- `SELECT`, returns rows in response to a query
- `INSERT`, adds new rows to a table
- `UPDATE`, alters existing rows in a table
- `DELETE`, removes rows from a table

We will be working almost exclusively with `SELECT` in order to ask questions of tables using spatial functions.

### 4.1.1 `SELECT` queries

A select query is generally of the form:

```
SELECT some_columns FROM some_data_source WHERE some_condition;
```

---

**Note:** For a synopsis of all `SELECT` parameters, see the [PostgresSQL documentation](#).

---

The `some_columns` are either column names or functions of column values. The `some_data_source` is either a single table, or a composite table created by joining two tables on a key or condition. The `some_condition` is a filter that restricts the number of rows to be returned.

“What are the names of all the neighborhoods in Brooklyn?”

We return to our `nyc_neighborhoods` table with a filter in hand. The table contains all the neighborhoods in New York, but we only want the ones in Brooklyn.

```
SELECT name
  FROM nyc_neighborhoods
 WHERE boroname = 'Brooklyn';
```

The query will run for even fewer (milli)seconds and return the 23 results.

Sometimes we will need to apply a function to the results of our query. For example,

“What is the number of letters in the names of all the neighborhoods in Brooklyn?”

Fortunately, PostgreSQL has a string length function, **char\_length(string)**.

```
SELECT char_length(name)
FROM nyc_neighborhoods
WHERE boroname = 'Brooklyn';
```

Often, we are less interested in the individual rows than in a statistic that applies to all of them. So knowing the lengths of the neighborhood names might be less interesting than knowing the average length of the names. Functions that take in multiple rows and return a single result are called “aggregate” functions.

PostgreSQL has a series of built-in aggregate functions, including the general purpose **avg()** for average values and **stddev()** for standard deviations.

“What is the average number of letters and standard deviation of number of letters in the names of all the neighborhoods in Brooklyn?”

```
SELECT avg(char_length(name)), stddev(char_length(name))
FROM nyc_neighborhoods
WHERE boroname = 'Brooklyn';
```

avg		stddev
11.7391304347826087		3.9105613559407395

The aggregate functions in our last example were applied to every row in the result set. What if we want the summaries to be carried out over smaller groups within the overall result set? For that we add a **GROUP BY** clause. Aggregate functions often need an added **GROUP BY** statement to group the result-set by one or more columns.

“What is the average number of letters in the names of all the neighborhoods in New York City, reported by borough?”

```
SELECT boroname, avg(char_length(name)), stddev(char_length(name))
FROM nyc_neighborhoods
GROUP BY boroname;
```

We include the **boroname** column in the output result so we can determine which statistic applies to which borough. In an aggregate query, you can only output columns that are either (a) members of the grouping clause or (b) aggregate functions.

boroname	avg	stddev
Brooklyn	11.7391304347826087	3.9105613559407395
Manhattan	11.8214285714285714	4.3123729948325257
The Bronx	12.0416666666666667	3.6651017740975152
Queens	11.6666666666666667	5.0057438272815975
Staten Island	12.2916666666666667	5.2043390480959474

### 4.1.2 Function List

`avg(expression)`: PostgreSQL aggregate function that returns the average value of a numeric column.

`char_length(string)`: PostgreSQL string function that returns the number of character in a string.

`stddev(expression)`: PostgreSQL aggregate function that returns the standard deviation of input values.

## 4.2 Simple SQL Exercises

Using the `nyc_census_blocks` table, answer the following questions (don't peak at the answers!).

Here is some helpful information to get started. Recall from the [About Our Data](#) section our `nyc_census_blocks` table definition.

<b>blkid</b>	A 15-digit code that uniquely identifies every census <b>block</b> . (“360050001009000”)
<b>popn_total</b>	Total number of people in the census block
<b>popn_white</b>	Number of people self-identifying as “white” in the block
<b>popn_black</b>	Number of people self-identifying as “black” in the block
<b>popn_nativ</b>	Number of people self-identifying as “native american” in the block
<b>popn_asian</b>	Number of people self-identifying as “asias” in the block
<b>popn_other</b>	Number of people self-identifying with other categories in the block
<b>hous_total</b>	Number of housing units in the block
<b>hous_own</b>	Number of owner-occupied housing units in the block
<b>hous_rent</b>	Number of renter-occupied housing units in the block
<b>boroname</b>	Name of the New York borough. Manhattan, The Bronx, Brooklyn, Staten Island, Queens
<b>geom</b>	Polygon boundary of the block

And, here are some common SQL aggregation functions you might find useful:

- `avg()` - the average (mean) of the values in a set of records
- `sum()` - the sum of the values in a set of records
- `count()` - the number of records in a set of records

Now the questions:

- **“What is the geometry value for the street named ‘Atlantic Commons’?”**

```
SELECT ST_AsText(geom)
  FROM nyc_streets
 WHERE name = 'Atlantic Commons';
```

```
MULTILINESTRING((586781.701577724 4504202.15314339, 586863.51964484,
→4504215.9881701))
```

- **“What neighborhood and borough is Atlantic Commons in?”**

```
SELECT name, boroname
  FROM nyc_neighborhoods
 WHERE ST_Intersects(
    geom,
```

```
ST_GeomFromText('LINESTRING(586782 4504202,586864 4504216)', 26918)
);
```

name	boroname
Fort Green	Brooklyn

**Note:** “Hey, why did you change from a ‘MULTILINESTRING’ to a ‘LINESTRING’?” Spatially they describe the same shape, so going from a single-item multi-geometry to a singleton saves a few keystrokes.

More importantly, we also rounded the coordinates to make them easier to read, which does actually change results: we couldn’t use the ST\_Touches() predicate to find out which roads join Atlantic Commons, because the coordinates are not exactly the same anymore.

- “What streets does Atlantic Commons join with?”

```
SELECT name
FROM nyc_streets
WHERE ST_DWithin(
    geom,
    ST_GeomFromText('LINESTRING(586782 4504202,586864 4504216)', 26918),
    0.1
);
```

name
Cumberland St
Atlantic Commons



- “Approximately how many people live on (within 50 meters of) Atlantic Commons?”

```
SELECT Sum(popn_total)
  FROM nyc_census_blocks
 WHERE ST_DWithin(
    geom,
    ST_GeomFromText('LINESTRING(586782 4504202, 586864 4504216)',_
→26918),
    50
);
```

```
1438
```

### 4.2.1 Function List

`avg(expression)`: PostgreSQL aggregate function that returns the average value of a numeric column.

`count(expression)`: PostgreSQL aggregate function that returns the number of records in a set of records.

`sum(expression)`: PostgreSQL aggregate function that returns the sum of records in a set of records.

## 4.3 Geometries

### 4.3.1 Introduction

In the previous *section*, we loaded a variety of data. Before we start playing with our data lets have a look at some simpler examples. In pgAdmin, once again select the `nyc` database and open the SQL query tool. Paste this example SQL code into the pgAdmin SQL Editor window (removing any text that may be there by default) and then execute.

```
CREATE TABLE geometries (name varchar, geom geometry);

INSERT INTO geometries VALUES
('Point', 'POINT(0 0)'),
('Linestring', 'LINESTRING(0 0, 1 1, 2 1, 2 2)'),
('Polygon', 'POLYGON((0 0, 1 0, 1 1, 0 1, 0 0))'),
('PolygonWithHole', 'POLYGON((0 0, 10 0, 10 10, 0 10, 0 0), (1 1, 1 2, 2
→2, 2 1, 1 1))'),
('Collection', 'GEOMETRYCOLLECTION(POINT(2 0), POLYGON((0 0, 1 0, 1 1, 0
→1, 0 0)))');

SELECT name, ST_AsText(geom) FROM geometries;
```

The screenshot shows a PostgreSQL/PostGIS query editor interface. At the top, there's a toolbar with various icons for file operations, search, and database management. Below the toolbar, the connection information is displayed: "nyc/postgres@postgis". The main area is divided into two panes: "Query Editor" and "Query History". The "Query Editor" pane contains the following SQL code:

```

1 CREATE TABLE geometries (name varchar, geom geometry);
2
3 INSERT INTO geometries VALUES
4     ('Point', 'POINT(0 0)'),
5     ('Linestring', 'LINESTRING(0 0, 1 1, 2 1, 2 2)'),
6     ('Polygon', 'POLYGON((0 0, 1 0, 1 1, 0 1, 0 0))'),
7     ('PolygonWithHole', 'POLYGON((0 0, 10 0, 10 10, 0 10, 0 0),(1 1, 1 2, 2 2 1 1))'),
8     ('Collection', 'GEOMETRYCOLLECTION(POINT(2 0),POLYGON((0 0, 1 0, 1 1, 0 1)))')
9
10    SELECT name, ST_AsText(geom) FROM geometries;

```

The "Data Output" pane below shows the results of the query:

	name	st_astext
1	Point	POINT(0 0)
2	Linestring	LINESTRING(0 0, 1 1, 2 1, 2 2)
3	Polygon	POLYGON((0 0, 1 0, 1 1, 0 1, 0 0))
4	PolygonWithHole	POLYGON((0 0, 10 0, 10 10, 0 10, 0 0),(1 1, 1 2, 2 2 1 1))
5	Collection	GEOMETRYCOLLECTION(POINT(2 0),POLYGON((0 0, 1 0, 1 1, 0 1)))

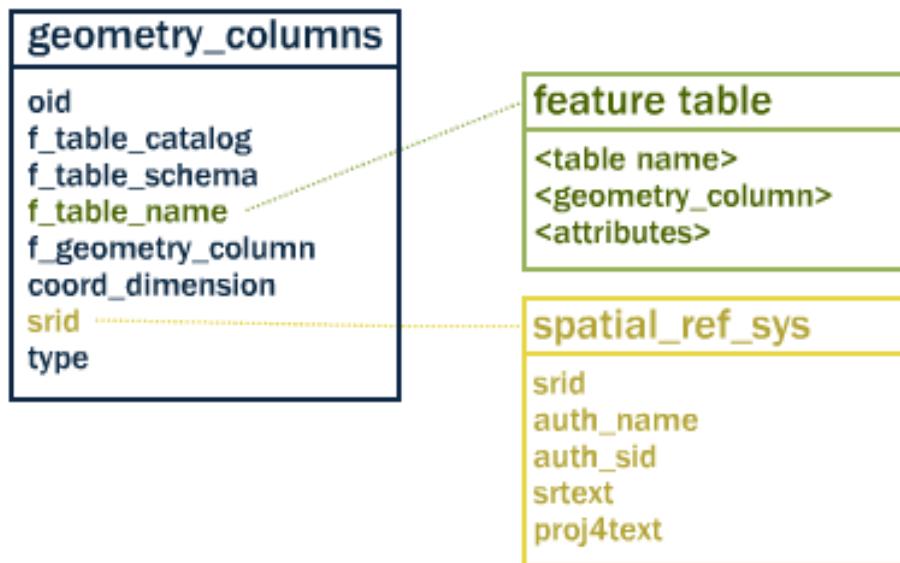
The above example CREATEs a table (**geometries**) then INSERTs five geometries: a point, a line, a polygon, a polygon with a hole, and a collection. Finally, the inserted rows are SELECTed and displayed in the Output pane.

### 4.3.2 Metadata Tables

In conformance with the Simple Features for SQL ([SFSQ](#)) specification, PostGIS provides two tables to track and report on the geometry types available in a given database.

- The first table, `spatial_ref_sys`, defines all the spatial reference systems known to the database and will be described in greater detail later.
- The second table (actually, a view), `geometry_columns`, provides a listing of all “features” (defined as an object with geometric attributes), and the basic details of those features.

## Table Relationships



Let's have a look at the `geometry_columns` table in our database. Paste this command in the Query Tool as before:

```
SELECT * FROM geometry_columns;
```

Screenshot of a PostgreSQL query tool interface showing the results of the `SELECT * FROM geometry_columns;` query.

**Query Editor:**

```
1 SELECT * FROM geometry_columns;
```

**Results:**

	f_table_catalog	f_table_schema	f_table_name	f_geometry_column	coord_dimension	srid	type
1	nyc	public	nyc_census_blo...	geom		2	MULTIPOLYGON
2	nyc	public	nyc_homicides	geom		2	POINT
3	nyc	public	nyc_neighborhoods	geom		2	MULTIPOLYGON
4	nyc	public	nyc_streets	geom		2	MULTILINESTRING
5	nyc	public	nyc_subway_stations	geom		2	POINT
6	nyc	public	geometries	geom		0	GEOMETRY

- `f_table_catalog`, `f_table_schema`, and `f_table_name` provide the fully qualified name of the feature table containing a given `geometry`.
- `f_geometry_column` is the name of the column that contains a geometry – for feature tables with multiple geometry columns, there will be one record for each.
- `coord_dimension` and `srid` define the dimension of the geometry (2-, 3- or 4-dimensional) and the Spatial Reference system identifier that refers to the `spatial_ref_sys` table respectively.
- The `type` column defines the type of geometry as described below; we've seen Point and Linestring types so far.

By querying this table, GIS clients and libraries can determine what to expect when retrieving data and can perform any necessary projection, processing or rendering without needing to inspect each geometry.

---

**Note:** Do some or all of your `nyc` tables not have an `srid` of 26918? It's easy to fix by updating the table

```
SELECT UpdateGeometrySRID('nyc_neighborhoods', 'geom', 26918);
```

---

### 4.3.3 Representing Real World Objects

The Simple Features for SQL ([SFSQ](#)L) specification, the original guiding standard for PostGIS development, defines how a real world object is represented. By taking a continuous shape and digitizing it at a fixed resolution we achieve a passable representation of the object. SFSQ only handled 2-dimensional representations. PostGIS has extended that to include 3- and 4-dimensional representations; more recently the SQL-Multimedia Part 3 ([SQL/MM](#)) specification has officially defined their own representation.

Our example table contains a mixture of different geometry types. We can collect general information about each object using functions that read the geometry metadata.

- `ST_GeometryType (geometry)` returns the type of the geometry
- `ST_NDims (geometry)` returns the number of dimensions of the geometry
- `ST_SRID (geometry)` returns the spatial reference identifier number of the geometry

```
SELECT name, ST_GeometryType(geom), ST_NDims(geom), ST_SRID(geom)
  FROM geometries;
```

name	st_geometrype	st_ndims	st_srid
Point	ST_Point	2	0
Polygon	ST_Polygon	2	0
PolygonWithHole	ST_Polygon	2	0
Collection	ST_GeometryCollection	2	0
Linestring	ST_LineString	2	0

### Points



A spatial **point** represents a single location on the Earth. This point is represented by a single coordinate (including either 2-, 3- or 4-dimensions). Points are used to represent objects when the exact details, such as shape and size, are not important at the target scale. For example, cities on a map of the world can be described as points, while a map of a single state might represent cities as polygons.

```
SELECT ST_AsText(geom)  
FROM geometries  
WHERE name = 'Point';
```

```
POINT(0 0)
```

Some of the specific spatial functions for working with points are:

- **ST\_X(geometry)** returns the X ordinate
- **ST\_Y(geometry)** returns the Y ordinate

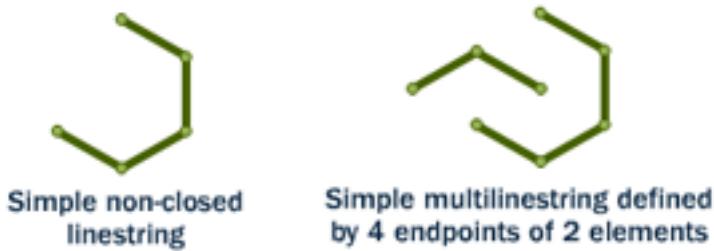
So, we can read the ordinates from a point like this:

```
SELECT ST_X(geom), ST_Y(geom)  
FROM geometries  
WHERE name = 'Point';
```

The New York City subway stations (`nyc_subway_stations`) table is a data set represented as points. The following SQL query will return the geometry associated with one point (in the **ST\_AsText** column).

```
SELECT name, ST_AsText(geom)  
FROM nyc_subway_stations  
LIMIT 1;
```

### Linestrings



A **linestring** is a path between locations. It takes the form of an ordered series of two or more points. Roads and rivers are typically represented as linestrings. A linestring is said to be **closed** if it starts and ends on the same point. It is said to be **simple** if it does not cross or touch itself (except at its endpoints if it is closed). A linestring can be both **closed** and **simple**.

The street network for New York (`nyc_streets`) was loaded earlier in the workshop. This dataset contains details such as name, and type. A single real world street may consist of many linestrings, each representing a segment of road with different attributes.

The following SQL query will return the geometry associated with one linestring (in the `ST_AsText` column).

```
SELECT ST_AsText(geom)
  FROM geometries
 WHERE name = 'Linestring';
```

```
LINESTRING(0 0, 1 1, 2 1, 2 2)
```

Some of the specific spatial functions for working with linestrings are:

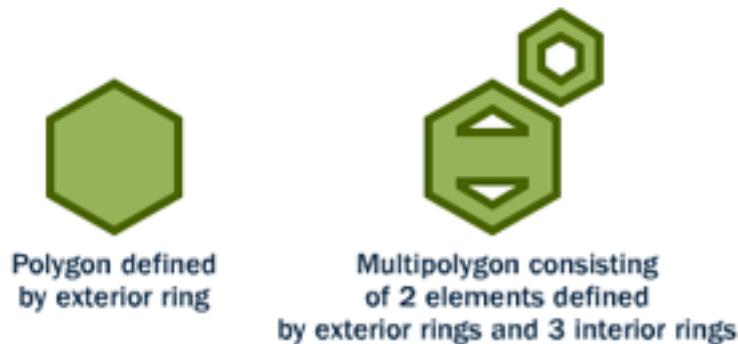
- `ST_Length(geometry)` returns the length of the linestring
- `ST_StartPoint(geometry)` returns the first coordinate as a point
- `ST_EndPoint(geometry)` returns the last coordinate as a point
- `ST_NPoints(geometry)` returns the number of coordinates in the linestring

So, the length of our linestring is:

```
SELECT ST_Length(geom)
  FROM geometries
 WHERE name = 'Linestring';
```

```
3.41421356237309
```

## Polygons



A polygon is a representation of an area. The outer boundary of the polygon is represented by a ring. This ring is a linestring that is both closed and simple as defined above. Holes within the polygon are also represented by rings.

Polygons are used to represent objects whose size and shape are important. City limits, parks, building footprints or bodies of water are all commonly represented as polygons when the scale is sufficiently high to see their area. Roads and rivers can sometimes be represented as polygons.

The following SQL query will return the geometry associated with one linestring (in the `ST_AsText` column).

```
SELECT ST_AsText(geom)
  FROM geometries
 WHERE name LIKE 'Polygon%';
```

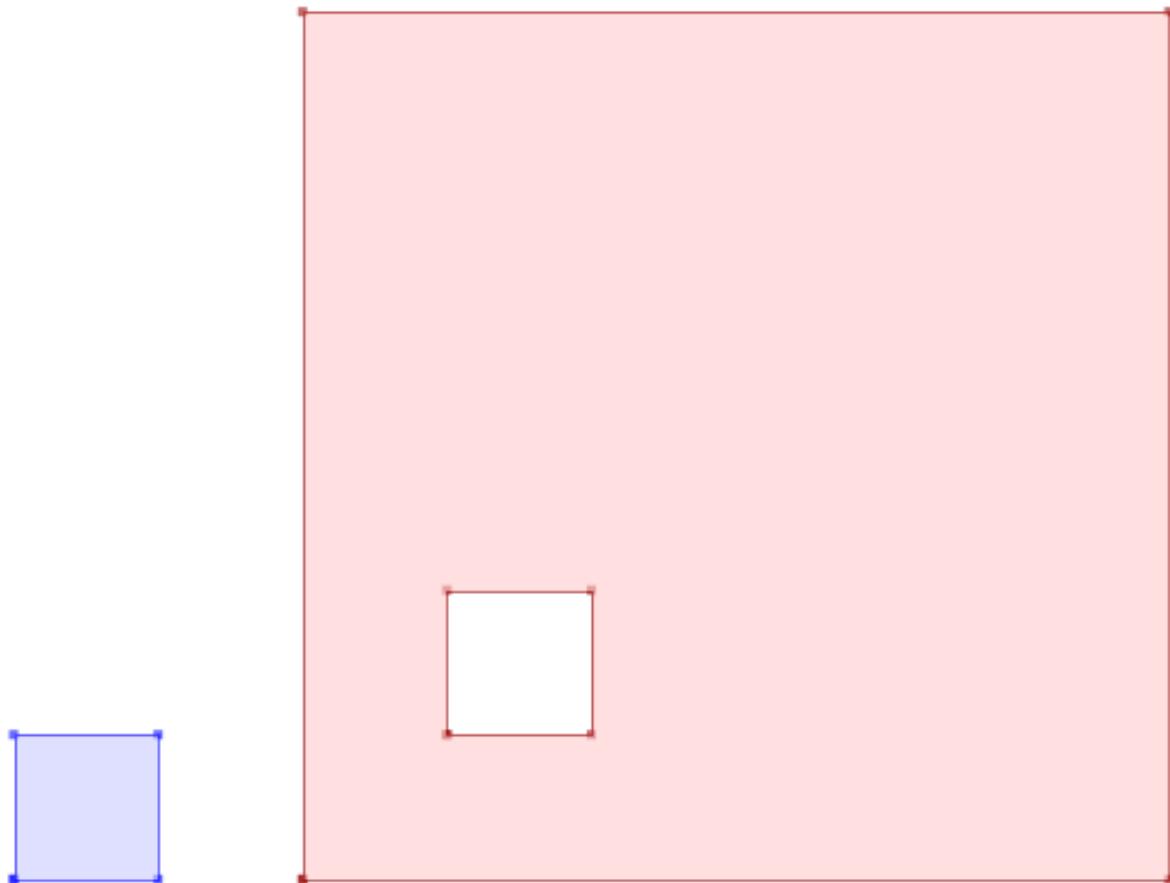
---

**Note:** Rather than using an = sign in our WHERE clause, we are using the LIKE operator to carry out a string matching operation. **You may be used to the “\*” symbol as a “glob” for pattern matching, but in SQL the “%” symbol is used**, along with the LIKE operator to tell the system to do globbing.

---

```
POLYGON((0 0, 1 0, 1 1, 0 1, 0 0))
POLYGON((0 0, 10 0, 10 10, 0 10, 0 0), (1 1, 1 2, 2 2, 2 1, 1 1))
```

The first polygon has only one ring. The second one has an interior “hole”. Most graphics systems include the concept of a “polygon”, but GIS systems are relatively unique in allowing polygons to explicitly have holes.



Some of the specific spatial functions for working with polygons are:

- **ST\_Area(geometry)** returns the area of the polygons
- **ST\_NRings(geometry)** returns the number of rings (usually 1, more of there are holes)
- **ST\_ExteriorRing(geometry)** returns the outer ring as a linestring
- **ST\_InteriorRingN(geometry, n)** returns a specified interior ring as a linestring
- **ST\_Perimeter(geometry)** returns the length of all the rings

We can calculate the area of our polygons using the area function:

```
SELECT name, ST_Area(geom)
  FROM geometries
 WHERE name LIKE 'Polygon%';
```

Polygon	1
PolygonWithHole	99

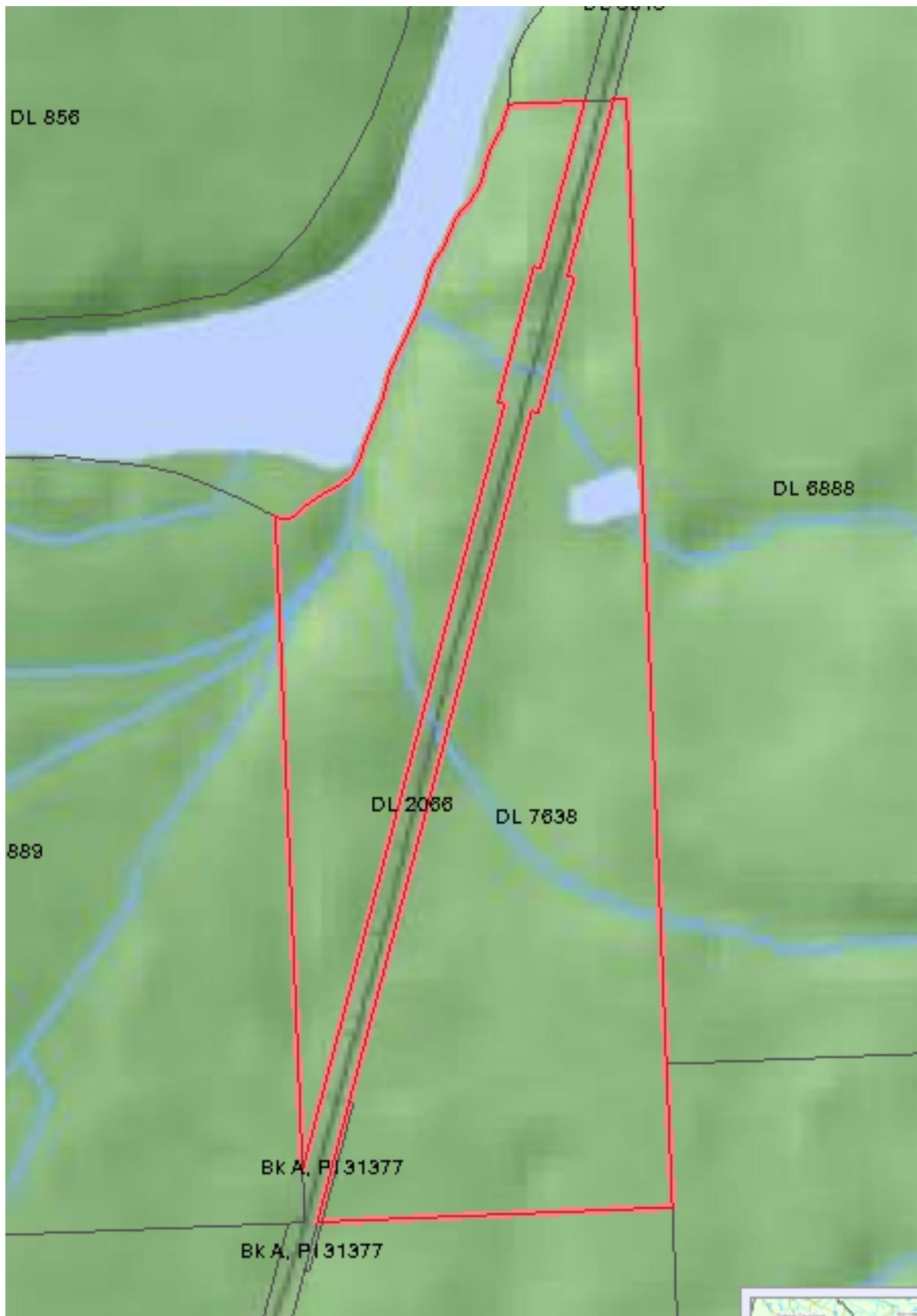
Note that the polygon with a hole has an area that is the area of the outer shell (a 10x10 square) minus the area of the hole (a 1x1 square).

## Collections

There are four collection types, which group multiple simple geometries into sets.

- **MultiPoint**, a collection of points
- **MultiLineString**, a collection of linestrings
- **MultiPolygon**, a collection of polygons
- **GeometryCollection**, a heterogeneous collection of any geometry (including other collections)

Collections are another concept that shows up in GIS software more than in generic graphics software. They are useful for directly modeling real world objects as spatial objects. For example, how to model a lot that is split by a right-of-way? As a **MultiPolygon**, with a part on either side of the right-of-way.



Our example collection contains a polygon and a point:

```
SELECT name, ST_AsText(geom)
  FROM geometries
 WHERE name = 'Collection';
```

```
GEOGRAPHYCOLLECTION(POINT(2 0),POLYGON((0 0, 1 0, 1 1, 0 1, 0 0)))
```



Some of the specific spatial functions for working with collections are:

- **ST\_NumGeometries(geometry)** returns the number of parts in the collection
- **ST\_GeometryN(geometry, n)** returns the specified part
- **ST\_Area(geometry)** returns the total area of all polygonal parts
- **ST\_Length(geometry)** returns the total length of all linear parts

#### 4.3.4 Geometry Input and Output

Within the database, geometries are stored on disk in a format only used by the PostGIS program. In order for external programs to insert and retrieve useful geometries, they need to be converted into a format that other applications can understand. Fortunately, PostGIS supports emitting and consuming geometries in a large number of formats:

- Well-known text ([WKT](#))
  - **ST\_GeomFromText(text, srid)** returns geometry
  - **ST\_AsText(geometry)** returns text
  - **ST\_AsEWKT(geometry)** returns text
- Well-known binary ([WKB](#))
  - **ST\_GeomFromWKB(bytea)** returns geometry
  - **ST\_AsBinary(geometry)** returns bytea
  - **ST\_AsEWKB(geometry)** returns bytea
- Geographic Mark-up Language ([GML](#))
  - **ST\_GeomFromGML(text)** returns geometry
  - **ST\_AsGML(geometry)** returns text
- Keyhole Mark-up Language ([KML](#))
  - **ST\_GeomFromKML(text)** returns geometry
  - **ST\_AsKML(geometry)** returns text
- [GeoJSON](#)
  - **ST\_AsGeoJSON(geometry)** returns text

- Scalable Vector Graphics (*SVG*)
    - **ST\_AsSVG(geometries)** returns text

The most common use of a constructor is to turn a text representation of a geometry into an internal representation:

Note that in addition to a text parameter with a geometry representation, we also have a numeric parameter providing the *SRID* of the geometry.

The following SQL query shows an example of [WKB](#) representation (the call to `encode()` is required to convert the binary output into an ASCII form for printing):

```
SELECT encode(  
    ST_AsBinary(ST_GeometryFromText('LINESTRING(0 0,1 0)'),  
    'hex');
```

For the purposes of this workshop we will continue to use WKT to ensure you can read and understand the geometries we're viewing. However, most actual processes, such as viewing data in a GIS application, transferring data to a web service, or processing data remotely, WKB is the format of choice.

Since WKT and WKB were defined in the [SFSQL](#) specification, they do not handle 3- or 4-dimensional geometries. For these cases PostGIS has defined the Extended Well Known Text (EWKT) and Extended Well Known Binary (EWKB) formats. These provide the same formatting capabilities of WKT and WKB with the added dimensionality.

Here is an example of a 3D linestring in WKT:

```
SELECT ST_AsText(ST_GeometryFromText ('LINESTRING(0 0 0,1 0 0,1 1 2)'));
```

LINESTRING Z (0 0 0, 1 0 0, 1 1 2)

Note that the text representation changes! This is because the text input routine for PostGIS is liberal in what it consumes. It will consume

- hex-encoded EWKB,
  - extended well-known text, and
  - ISO standard well-known text.

On the output side, the **ST\_AsText** function is conservative, and only emits ISO standard well-known text.

In addition to the `ST_GeometryFromText` function, there are many other ways to create geometries from well-known text or similar formatted inputs:

```
-- Using ST_GeomFromText with the SRID parameter
SELECT ST_GeomFromText('POINT(2 2)', 4326);

-- Using ST_GeomFromText without the SRID parameter
SELECT ST_SetSRID(ST_GeomFromText('POINT(2 2)'), 4326);

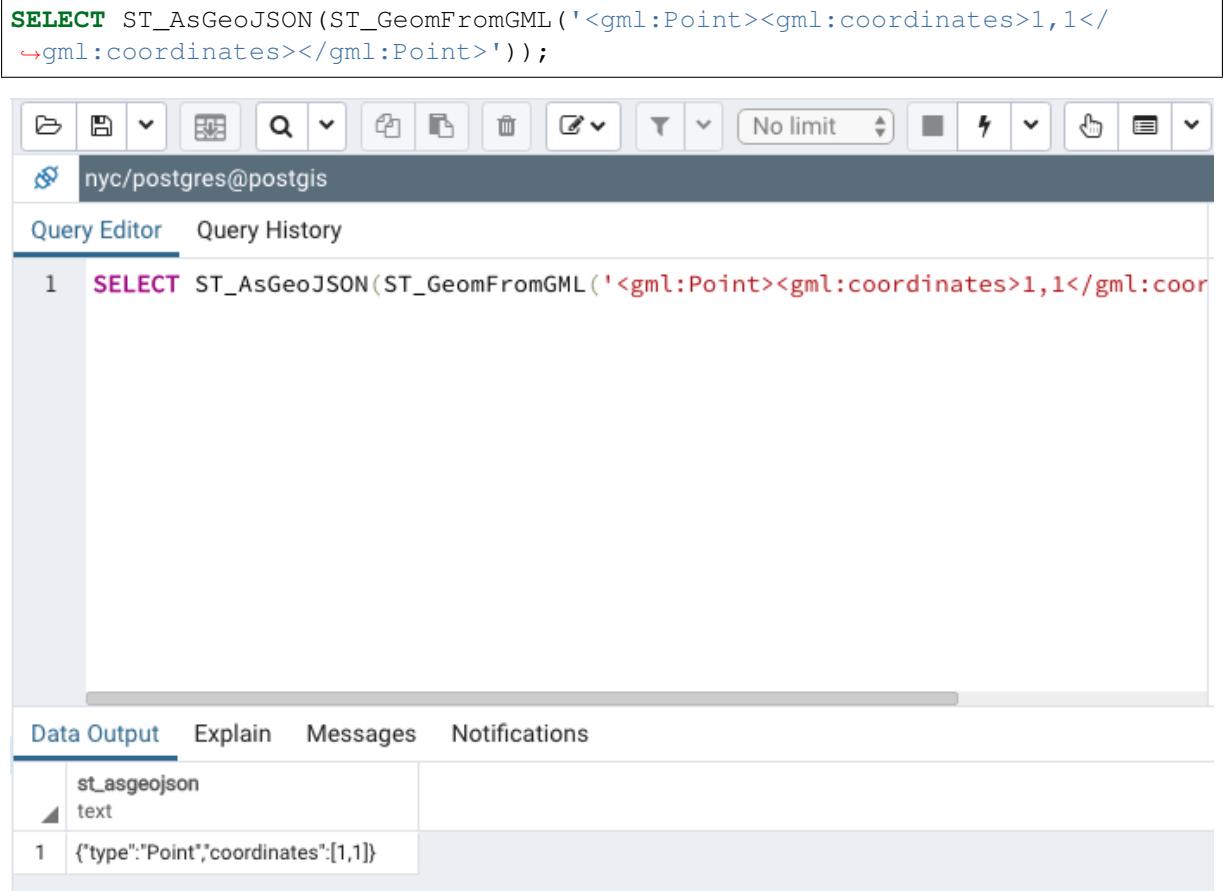
-- Using a ST_Make* function
SELECT ST_SetSRID(ST_MakePoint(2, 2), 4326);

-- Using PostgreSQL casting syntax and ISO WKT
```

```
SELECT ST_SetSRID('POINT(2 2)::geometry, 4326);

-- Using PostgreSQL casting syntax and extended WKT
SELECT 'SRID=4326;POINT(2 2)::geometry';
```

In addition to emitters for the various forms (WKT, WKB, GML, KML, JSON, SVG), PostGIS also has consumers for four (WKT, WKB, GML, KML). Most applications use the WKT or WKB geometry creation functions, but the others work too. Here's an example that consumes GML and output JSON:



The screenshot shows the pgAdmin III interface. The connection is set to 'nyc/postgres@postgis'. The 'Query Editor' tab is selected. A single-line query is entered: 'SELECT ST\_AsGeoJSON(ST\_GeomFromGML('<gml:Point><gml:coordinates>1,1</gml:coordinates></gml:Point>'))'. The results pane shows a table with one row labeled 'st\_asgeojson' containing the JSON output: '{ "type": "Point", "coordinates": [1,1] }'.

st_asgeojson
{ "type": "Point", "coordinates": [1,1] }

### 4.3.5 Casting from Text

The [WKT](#) strings we've seen so far have been of type 'text' and we have been converting them to type 'geometry' using PostGIS functions like `ST_GeomFromText()`.

PostgreSQL includes a short form syntax that allows data to be converted from one type to another, the casting syntax, *oldtype::newtype*. So for example, this SQL converts a double into a text string.

```
SELECT 0.9::text;
```

Less trivially, this SQL converts a [WKT](#) string into a geometry:

```
SELECT 'POINT(0 0)::geometry;
```

One thing to note about using casting to create geometries: unless you specify the SRID, you will get a geometry with an unknown SRID. You can specify the SRID using the “extended” well-known text form, which includes an SRID block at the front:

```
SELECT 'SRID=4326;POINT(0 0)' ::geometry;
```

It's very common to use the casting notation when working with *WKT*, as well as *geometry* and *geography* columns (see geography).

### 4.3.6 Function List

**ST\_Area:** Returns the area of the surface if it is a polygon or multi-polygon. For “geometry” type area is in SRID units. For “geography” area is in square meters.

**ST\_AsText:** Returns the Well-Known Text (WKT) representation of the geometry/geography without SRID metadata.

**ST\_AsBinary:** Returns the Well-Known Binary (WKB) representation of the geometry/geography without SRID meta data.

**ST\_EndPoint:** Returns the last point of a LINESTRING geometry as a POINT.

**ST\_AsEWKB:** Returns the Well-Known Binary (WKB) representation of the geometry with SRID meta data.

**ST\_AsEWKT:** Returns the Well-Known Text (WKT) representation of the geometry with SRID meta data.

**ST\_AsGeoJSON:** Returns the geometry as a GeoJSON element.

**ST\_AsGML:** Returns the geometry as a GML version 2 or 3 element.

**ST\_AsKML:** Returns the geometry as a KML element. Several variants. Default version=2, default precision=15.

**ST\_AsSVG:** Returns a Geometry in SVG path data given a geometry or geography object.

**ST\_ExteriorRing:** Returns a line string representing the exterior ring of the POLYGON geometry. Return NULL if the geometry is not a polygon.

**ST\_GeometryN:** Returns the 1-based Nth geometry if the geometry is a GEOMETRYCOLLECTION, (MULTI)POINT, (MULTI)LINESTRING, (MULTI)CURVE, (MULTI)POLYGON or POLYHEDRAL-SURFACE. Otherwise, return NULL.

**ST\_GeomFromGML:** Takes as input GML representation of geometry and outputs a PostGIS geometry object.

**ST\_GeomFromKML:** Takes as input KML representation of geometry and outputs a PostGIS geometry object

**ST\_GeomFromText:** Returns a specified ST\_Geometry value from Well-Known Text representation (WKT).

**ST\_GeomFromWKB:** Creates a geometry instance from a Well-Known Binary geometry representation (WKB) and optional SRID.

**ST\_GeometryType:** Returns the geometry type of the ST\_Geometry value.

**ST\_InteriorRingN:** Returns the Nth interior linestring ring of the polygon geometry. Return NULL if the geometry is not a polygon or the given N is out of range. index starts at 1.

**ST\_Length:** Returns the 2d length of the geometry if it is a linestring or multilinestring. geometry are in units of spatial reference and geography are in meters (default spheroid)

**ST\_NDims:** Returns coordinate dimension of the geometry. Values are: 2,3 or 4.

**ST\_NPoints:** Returns the number of points (vertexes) in a geometry.

**ST\_NRings:** If the geometry is a polygon or multi-polygon returns the number of rings. Unlike :command: *NumInteriorRings*, it counts the outer rings as well.

**ST\_NumGeometries:** If geometry is a GEOMETRYCOLLECTION (or MULTI\*) returns the number of geometries, for single geometries will return 1, otherwise return NULL.

**ST\_Perimeter:** Returns the length measurement of the boundary of an ST\_Surface or ST\_MultiSurface value. (Polygon, Multipolygon)

**ST\_SRID:** Returns the spatial reference identifier for the ST\_Geometry as defined in spatial\_ref\_sys table.

**ST\_StartPoint:** Returns the first point of a LINESTRING or CIRCULARLINESTRING geometry as a POINT.

**ST\_X:** Returns the X coordinate of the point, or NULL if not available. Input must be a point.

**ST\_Y:** Returns the Y coordinate of the point, or NULL if not available. Input must be a point.

## 4.4 Geometry Exercises

Here's a reminder of all the functions we have seen so far. They should be useful for the exercises!

- **sum(expression)** aggregate to return a sum for a set of records
- **count(expression)** aggregate to return the size of a set of records
- **ST\_GeometryType(geometry)** returns the type of the geometry
- **ST\_NDims(geometry)** returns the number of dimensions of the geometry
- **ST\_SRID(geometry)** returns the spatial reference identifier number of the geometry
- **ST\_X(point)** returns the X ordinate
- **ST\_Y(point)** returns the Y ordinate
- **ST\_Length(linestring)** returns the length of the linestring
- **ST\_StartPoint(geometry)** returns the first coordinate as a point
- **ST\_EndPoint(geometry)** returns the last coordinate as a point
- **ST\_NPoints(geometry)** returns the number of coordinates in the linestring
- **ST\_Area(geometry)** returns the area of the polygons
- **ST\_NRings(geometry)** returns the number of rings (usually 1, more if there are holes)
- **ST\_ExteriorRing(polygon)** returns the outer ring as a linestring
- **ST\_InteriorRingN(polygon, integer)** returns a specified interior ring as a linestring
- **ST\_Perimeter(geometry)** returns the length of all the rings
- **ST\_NumGeometries(multi/geomcollection)** returns the number of parts in the collection
- **ST\_GeometryN(geometry, integer)** returns the specified part of the collection
- **ST\_GeomFromText(text)** returns geometry

- **ST\_AsText (geometry)** returns WKT text
- **ST\_AsEWKT (geometry)** returns EWKT text
- **ST\_GeomFromWKB (bytea)** returns geometry
- **ST\_AsBinary (geometry)** returns WKB bytea
- **ST\_AsEWKB (geometry)** returns EWKB bytea
- **ST\_GeomFromGML (text)** returns geometry
- **ST\_AsGML (geometry)** returns GML text
- **ST\_GeomFromKML (text)** returns geometry
- **ST\_AsKML (geometry)** returns KML text
- **ST\_AsGeoJSON (geometry)** returns JSON text
- **ST\_AsSVG (geometry)** returns SVG text

Also remember the tables we have available:

- nyc\_census\_blocks
  - blkid, popn\_total, boroname, geom
- nyc\_streets
  - name, type, geom
- nyc\_subway\_stations
  - name, geom
- nyc\_neighborhoods
  - name, boroname, geom

### 4.4.1 Exercises

- “What is the area of the ‘West Village’ neighborhood?”

```
SELECT ST_Area(geom)
  FROM nyc_neighborhoods
 WHERE name = 'West Village';
```

```
1044614.5296486
```

---

**Note:** The area is given in square meters. To get an area in hectares, divide by 10000. To get an area in acres, divide by 4047.

---

- “What is the area of Manhattan in acres?” (Hint: both nyc\_census\_blocks and nyc\_neighborhoods have a boroname in them.)

```
SELECT Sum(ST_Area(geom)) / 4047
  FROM nyc_neighborhoods
 WHERE boroname = 'Manhattan';
```

```
13965.3201224118
```

or...

```
SELECT Sum(ST_Area(geom)) / 4047
  FROM nyc_census_blocks
 WHERE boroname = 'Manhattan';
```

```
14601.3987215548
```

- “How many census blocks in New York City have a hole in them?”

```
SELECT Count(*)
  FROM nyc_census_blocks
 WHERE ST_NumInteriorRings(ST_GeometryN(geom, 1)) > 0;
```

**Note:** The ST\_NRings() functions might be tempting, but it also counts the exterior rings of multi-polygons as well as interior rings. In order to run ST\_NumInteriorRings() we need to convert the MultiPolygon geometries of the blocks into simple polygons, so we extract the first polygon from each collection using ST\_GeometryN(). Yuck!

```
43
```

- “What is the total length of streets (in kilometers) in New York City?” (Hint: The units of measurement of the spatial data are meters, there are 1000 meters in a kilometer.)

```
SELECT Sum(ST_Length(geom)) / 1000
  FROM nyc_streets;
```

```
10418.9047172
```

- “How long is ‘Columbus Cir’ (Columbus Circle)?

```
SELECT ST_Length(geom)
  FROM nyc_streets
 WHERE name = 'Columbus Cir';
```

```
308.34199
```

- “What is the JSON representation of the boundary of the ‘West Village’?”

```
SELECT ST_AsGeoJSON(geom)
  FROM nyc_neighborhoods
 WHERE name = 'West Village';
```

```
{"type": "MultiPolygon", "coordinates": [
    [[[583263.2776595836, 4509242.6260239873],
      [583276.81990686338, 4509378.825446927], ...,
      [583263.2776595836, 4509242.6260239873]]]]}
```

The geometry type is “MultiPolygon”, interesting!

- “How many polygons are in the ‘West Village’ multipolygon?”

```
SELECT ST_NumGeometries(geom)
  FROM nyc_neighborhoods
 WHERE name = 'West Village';
```

```
1
```

---

**Note:** It is not uncommon to find single-element MultiPolygons in spatial tables. Using MultiPolygons allows a table with only one geometry type to store both single- and multi-geometries without using mixed types.

---

- “What is the length of streets in New York City, summarized by type?”

```
SELECT type, Sum(ST_Length(geom)) AS length
  FROM nyc_streets
 GROUP BY type
 ORDER BY length DESC;
```

type	length
residential	8629870.33786606
motorway	403622.478126363
tertiary	360394.879051303
motorway_link	294261.419479668
secondary	276264.303897926
unclassified	166936.371604458
primary	135034.233017947
footway	71798.4878378096
service	28337.635038596
trunk	20353.5819826076
cycleway	8863.75144825929
pedestrian	4867.05032825026
construction	4803.08162103562
residential; motorway_link	3661.57506293745
trunk_link	3202.18981240201
primary_link	2492.57457083536
living_street	1894.63905457332
primary; residential; motorway_link; residential	1367.76576941335
undefined	380.53861910346
steps	282.745221342127
motorway_link; residential	215.07778911517

---

**Note:** The ORDER BY length DESC clause sorts the result by length in descending order. The result is that most prevalent types are first in the list.

---

## 4.5 Spatial Relationships

So far we have only used spatial functions that measure (**ST\_Area**, **ST\_Length**), serialize (**ST\_AsGML**), and deserialize (**ST\_GeomFromText**) geometries. What these functions have in common is that they only work on one geometry at a time.

Spatial databases are powerful because they not only store geometry, they also have the ability to compare *relationships between geometries*.

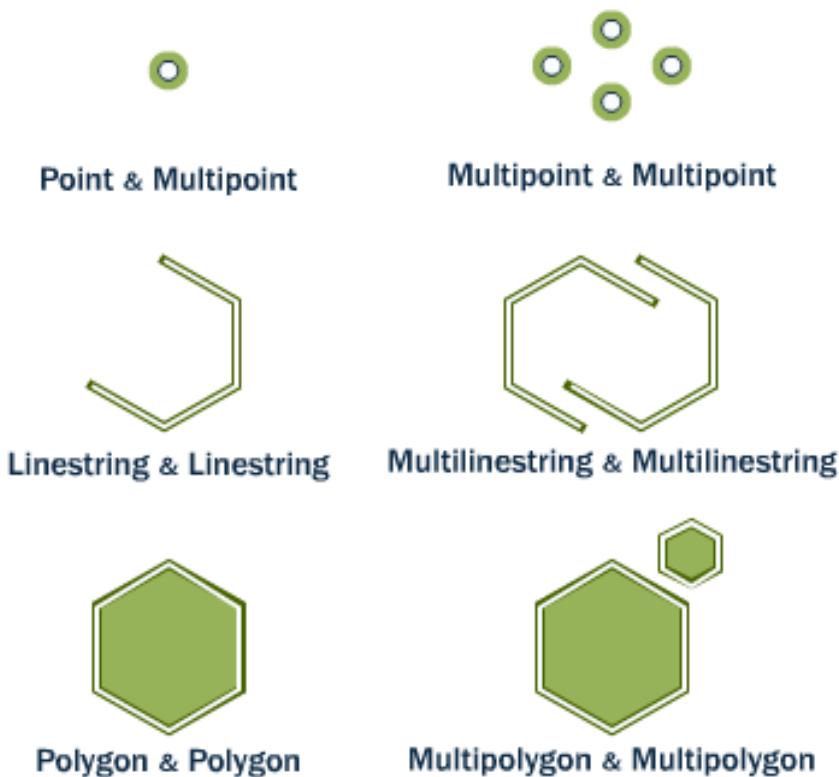
Questions like “Which are the closest bike racks to a park?” or “Where are the intersections of subway lines and streets?” can only be answered by comparing geometries representing the bike racks, streets, and subway lines.

The OGC standard defines the following set of methods to compare geometries.

#### 4.5.1 ST\_Equals

**ST\_Equals (geometry A, geometry B)** tests the spatial equality of two geometries.

## Equals



ST\_Equals returns TRUE if two geometries of the same type have identical x,y coordinate values, i.e. if the second shape is equal (identical) to the first shape.

First, let's retrieve a representation of a point from our `nyc_subway_stations` table. We'll take just the entry for ‘Broad St’.

```
SELECT name, geom, ST_AsText(geom)
FROM nyc_subway_stations
WHERE name = 'Broad St';
```

name	geom	st_
↳astext		
↳-----	-----+-----+	-----+
Broad St	0101000020266900000EEBD4CF27CF2141BC17D69516315141	↳
↳POINT(583571 4506714)		

Then, plug the geometry representation back into an **ST\_Equals** test:

```
SELECT name
FROM nyc_subway_stations
WHERE ST_Equals(geom, '0101000020266900000EEBD4CF27CF2141BC17D69516315141
↳');
```

Broad St
----------

---

**Note:** The representation of the point was not very human readable (0101000020266900000EEBD4CF27CF2141BC17D69516315141) but it was an exact representation of the coordinate values. For a test like equality, using the exact coordinates is necessary.

---

### 4.5.2 ST\_Intersects, ST\_Disjoint, ST\_Crosses and ST\_Overlaps

**ST\_Intersects**, **ST\_Crosses**, and **ST\_Overlaps** test whether the interiors of the geometries intersect.

**ST\_Intersects(geometry A, geometry B)** returns t (TRUE) if the two shapes have any space in common, i.e., if their boundaries or interiors intersect.

The opposite of **ST\_Intersects** is **ST\_Disjoint(geometry A, geometry B)**. If two geometries are disjoint, they do not intersect, and vice-versa. In fact, it is often more efficient to test “not intersects” than to test “disjoint” because the intersects tests can be spatially indexed, while the disjoint test cannot.

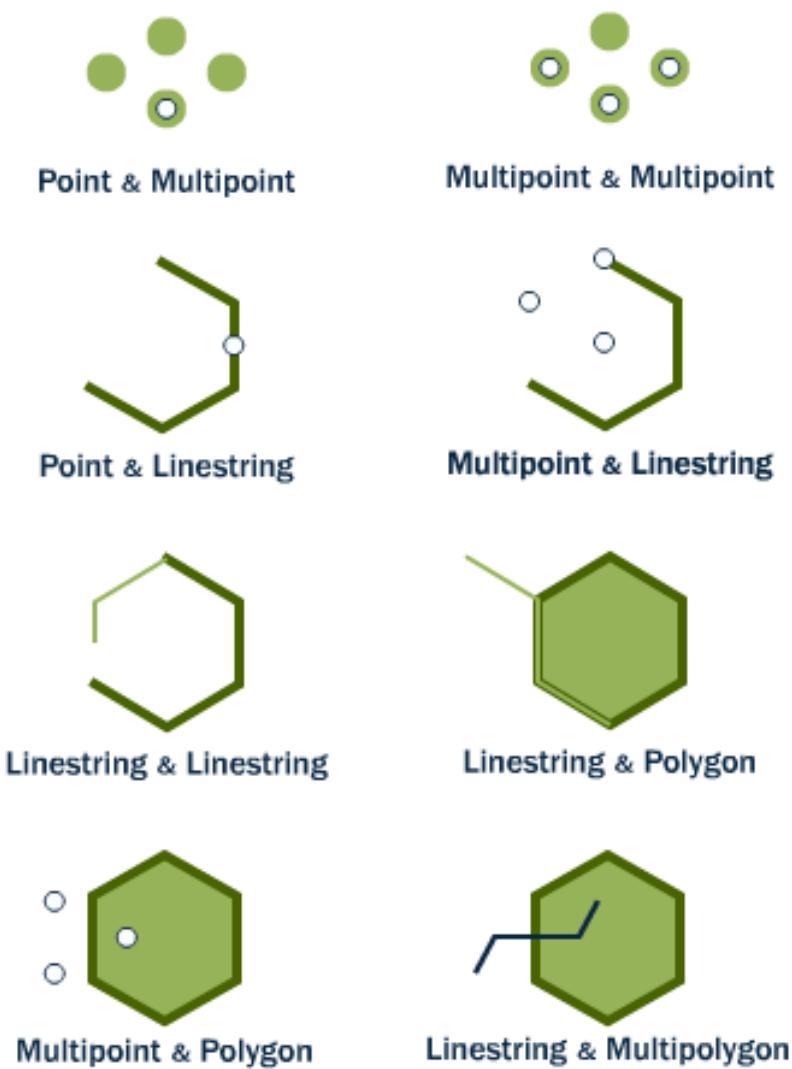
For multipoint/polygon, multipoint/linestring, linestring/linestring, linestring/polygon, and linestring/multipolygon comparisons, **ST\_Crosses(geometry A, geometry B)** returns t (TRUE) if the intersection results in a geometry whose dimension is one less than the maximum dimension of the two source geometries and the intersection set is interior to both source geometries.

**ST\_Overlaps(geometry A, geometry B)** compares two geometries of the same dimension and returns TRUE if their intersection set results in a geometry different from both but of the same dimension.

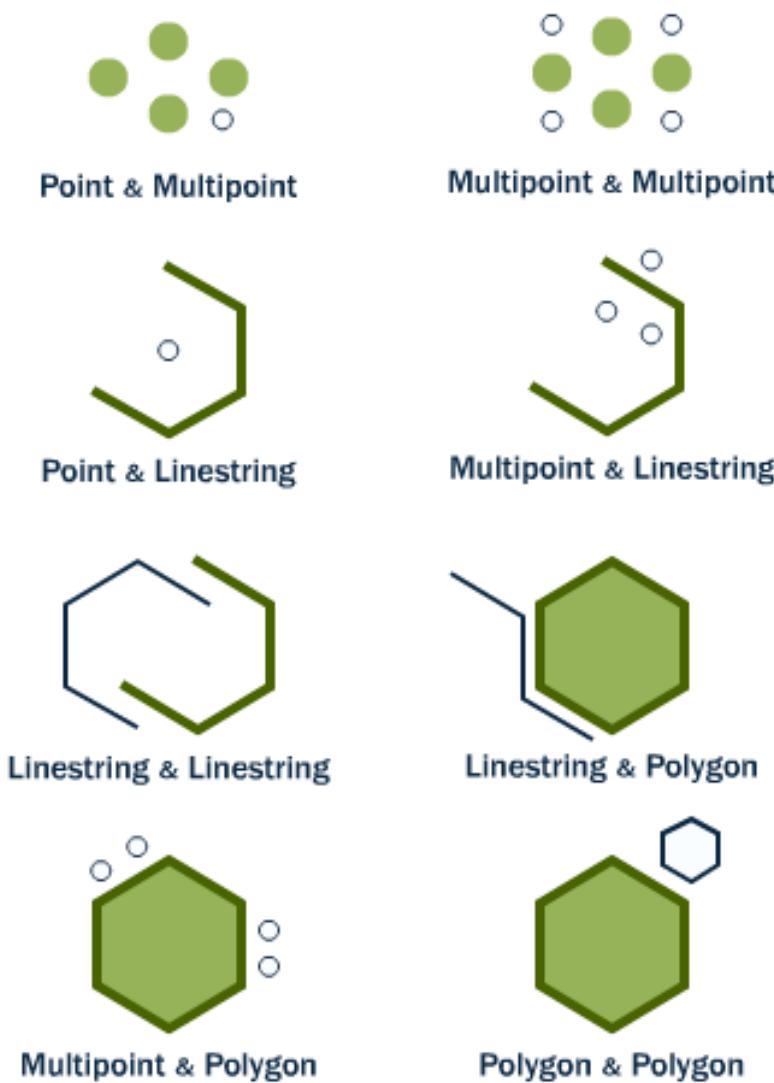
Let's take our Broad Street subway station and determine its neighborhood using the **ST\_Intersects** function:

```
SELECT name, ST_AsText(geom)
FROM nyc_subway_stations
WHERE name = 'Broad St';
```

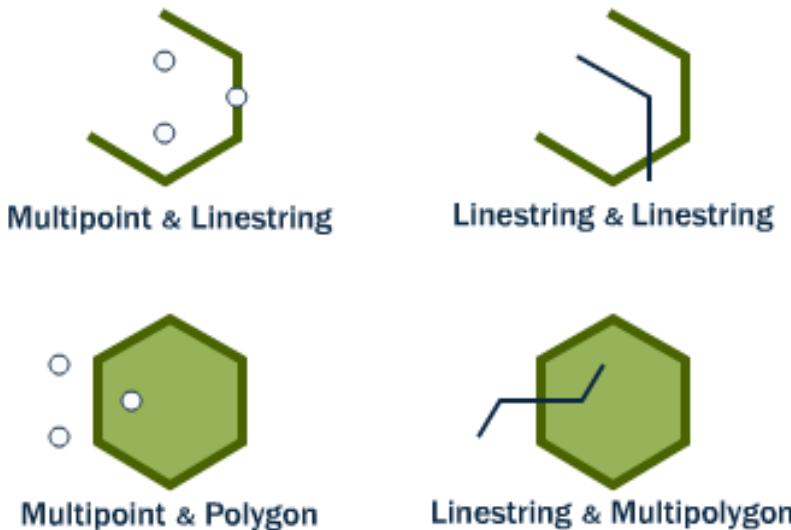
## Intersects



## Disjoint



## Cross



```
POINT(583571 4506714)
```

```
SELECT name, boroname
FROM nyc_neighborhoods
WHERE ST_Intersects(geom, ST_GeomFromText('POINT(583571 4506714)', 26918));
```

name	boroname
Financial District	Manhattan

### 4.5.3 ST\_Touches

**ST\_Touches** tests whether two geometries touch at their boundaries, but do not intersect in their interiors

**ST\_Touches(geometry A, geometry B)** returns TRUE if either of the geometries' boundaries intersect or if only one of the geometry's interiors intersects the other's boundary.

### 4.5.4 ST\_Within and ST\_Contains

**ST\_Within** and **ST\_Contains** test whether one geometry is fully within the other.

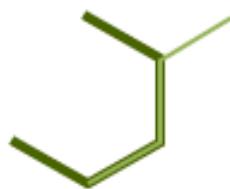
**ST\_Within(geometry A, geometry B)** returns TRUE if the first geometry is completely within the second geometry. ST\_Within tests for the exact opposite result of ST\_Contains.

**ST\_Contains(geometry A, geometry B)** returns TRUE if the second geometry is completely

## Overlap



Multipoint & Multipoint

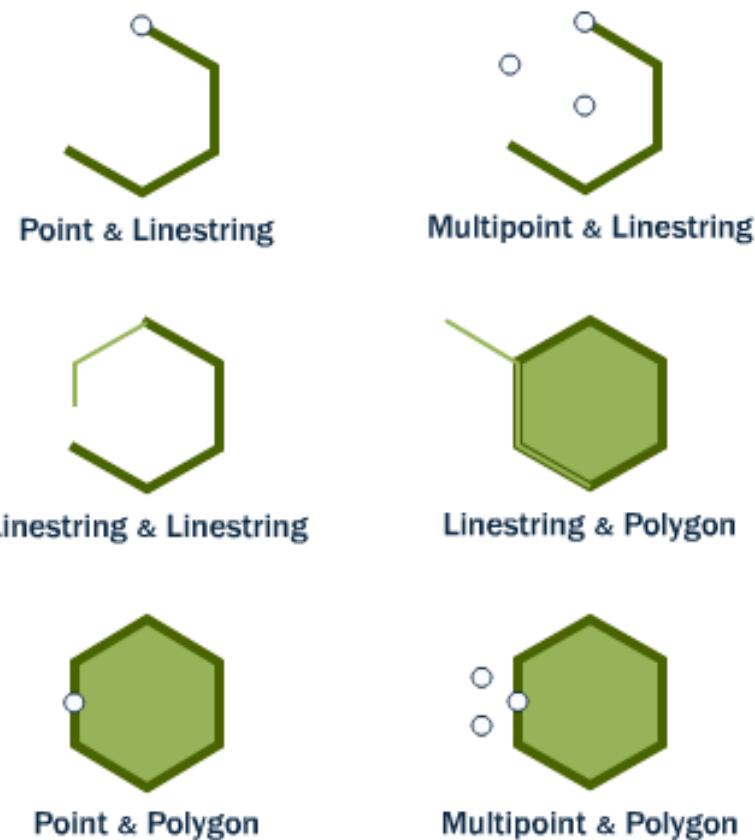


Linestring & Linestring



Polygon & Polygon

## Touch



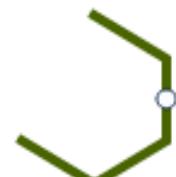
## Within/Contains



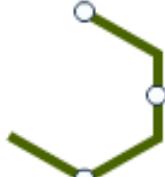
Point & Multipoint



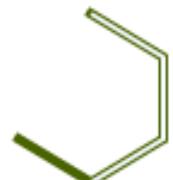
Multipoint & Multipoint



Point & Linestring



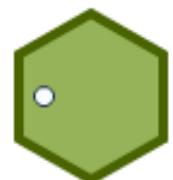
Multipoint & Linestring



Linestring & Linestring



Linestring & Polygon



Point & Polygon



Multipoint & Polygon

contained by the first geometry.

#### 4.5.5 ST\_Distance and ST\_DWithin

An extremely common GIS question is “find all the stuff within distance X of this other stuff”.

The **ST\_Distance (geometry A, geometry B)** calculates the *shortest* distance between two geometries and returns it as a float. This is useful for actually reporting back the distance between objects.

```
SELECT ST_Distance(
    ST_GeometryFromText('POINT(0 5)'),
    ST_GeometryFromText('LINESTRING(-2 2, 2 2)'));
```

3

For testing whether two objects are within a distance of one another, the **ST\_DWithin** function provides an index-accelerated true/false test. This is useful for questions like “how many trees are within a 500 meter buffer of the road?”. You don’t have to calculate an actual buffer, you just have to test the distance relationship.

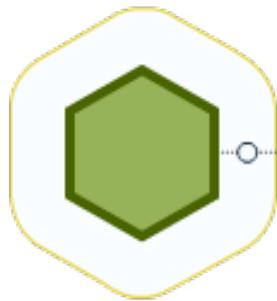
## ST\_Dwithin



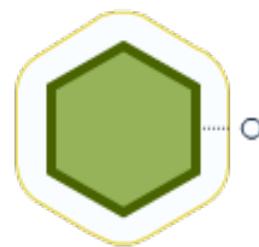
Point & Point (True)



Point & Point (False)



Polygon & Point (True)



Polygon & Point (False)

Using our Broad Street subway station again, we can find the streets nearby (within 10 meters of) the subway stop:

```
SELECT name
FROM nyc_streets
WHERE ST_DWithin(
    geom,
    ST_GeomFromText('POINT(583571 4506714)', 26918),
    10
);
```

```
name
-----
Wall St
Broad St
Nassau St
```

And we can verify the answer on a map. The Broad St station is actually at the intersection of Wall, Broad and Nassau Streets.



### 4.5.6 Function List

**ST\_Contains(geometry A, geometry B)**: Returns true if and only if no points of B lie in the exterior of A, and at least one point of the interior of B lies in the interior of A.

**ST\_Crosses(geometry A, geometry B)**: Returns TRUE if the supplied geometries have some, but not all, interior points in common.

**ST\_Disjoint(geometry A , geometry B):** Returns TRUE if the Geometries do not “spatially intersect” - if they do not share any space together.

**ST\_Distance(geometry A, geometry B):** Returns the 2-dimensional cartesian minimum distance (based on spatial ref) between two geometries in projected units.

**ST\_DWithin(geometry A, geometry B, radius):** Returns true if the geometries are within the specified distance (radius) of one another.

**ST\_Equals(geometry A, geometry B):** Returns true if the given geometries represent the same geometry. Directionality is ignored.

**ST\_Intersects(geometry A, geometry B):** Returns TRUE if the Geometries/Geography “spatially intersect” - (share any portion of space) and FALSE if they don’t (they are Disjoint).

**ST\_Overlaps(geometry A, geometry B):** Returns TRUE if the Geometries share space, are of the same dimension, but are not completely contained by each other.

**ST\_Touches(geometry A, geometry B):** Returns TRUE if the geometries have at least one point in common, but their interiors do not intersect.

**ST\_Within(geometry A , geometry B):** Returns true if the geometry A is completely inside geometry B

## 4.6 Spatial Relationships Exercises

Here’s a reminder of the functions we saw in the last section. They should be useful for the exercises!

- **sum(expression)** aggregate to return a sum for a set of records
- **count(expression)** aggregate to return the size of a set of records
- **ST\_Contains(geometry A, geometry B)** returns true if geometry A contains geometry B
- **ST\_Crosses(geometry A, geometry B)** returns true if geometry A crosses geometry B
- **ST\_Disjoint(geometry A , geometry B)** returns true if the geometries do not “spatially intersect”
- **ST\_Distance(geometry A, geometry B)** returns the minimum distance between geometry A and geometry B
- **ST\_DWithin(geometry A, geometry B, radius)** returns true if geometry A is radius distance or less from geometry B
- **ST\_Equals(geometry A, geometry B)** returns true if geometry A is the same as geometry B
- **ST\_Intersects(geometry A, geometry B)** returns true if geometry A intersects geometry B
- **ST\_Overlaps(geometry A, geometry B)** returns true if geometry A and geometry B share space, but are not completely contained by each other.
- **ST\_Touches(geometry A, geometry B)** returns true if the boundary of geometry A touches geometry B
- **ST\_Within(geometry A , geometry B)** returns true if geometry A is within geometry B

Also remember the tables we have available:

- nyc\_census\_blocks
  - blkid, popn\_total, boroname, geom
- nyc\_streets
  - name, type, geom
- nyc\_subway\_stations
  - name, geom
- nyc\_neighborhoods
  - name, boroname, geom

### 4.6.1 Exercises

- “What is the geometry value for the street named ‘Atlantic Commons’?”

```
SELECT ST_AsText(geom)
  FROM nyc_streets
 WHERE name = 'Atlantic Commons';
```

```
MULTILINESTRING((586781.701577724 4504202.15314339, 586863.51964484,
                  ↪4504215.9881701))
```

- “What neighborhood and borough is Atlantic Commons in?”

```
SELECT name, boroname
  FROM nyc_neighborhoods
 WHERE ST_Intersects(
    geom,
    ST_GeomFromText('LINESTRING(586782 4504202,586864 4504216)', 26918)
 );
```

name	boroname
Fort Green	Brooklyn

**Note:** “Hey, why did you change from a ‘MULTILINESTRING’ to a ‘LINESTRING’?” Spatially they describe the same shape, so going from a single-item multi-geometry to a singleton saves a few keystrokes.

More importantly, we also rounded the coordinates to make them easier to read, which does actually change results: we couldn’t use the ST\_Touches() predicate to find out which roads join Atlantic Commons, because the coordinates are not exactly the same anymore.

- “What streets does Atlantic Commons join with?”

```
SELECT name
  FROM nyc_streets
 WHERE ST_DWithin(
    geom,
    ST_GeomFromText('LINESTRING(586782 4504202,586864 4504216)', 26918),
```

```
    0.1
);
```

```
name
-----
Cumberland St
Atlantic Commons
```



- “Approximately how many people live on (within 50 meters of) Atlantic Commons?”

```
SELECT Sum(popn_total)
  FROM nyc_census_blocks
 WHERE ST_DWithin(
   geom,
   ST_GeomFromText('LINESTRING(586782 4504202,586864 4504216)',_
→26918),
   50
);
```

```
1438
```

## 4.7 Spatial Joins

Spatial joins are the bread-and-butter of spatial databases. They allow you to combine information from different tables by using spatial relationships as the join key. Much of what we think of as “standard GIS analysis” can be expressed as spatial joins.

To review how joins work, here’s a step-by-step look at joining the neighborhoods and subway stations tables and then subsequently adding a simple where clause. Joining two small tables effectively multiplies the number of rows into a much larger dataset with all possible combinations of the two tables. Then we use a where clause (or two) to narrow our focus.

```
--129 neighborhoods
SELECT name, boroname FROM nyc_neighborhoods

--491 subway stations
SELECT name, borough FROM nyc_subway_stations

--join produces 63339 rows (129 * 491)
SELECT n.name, n.boroname,
       s.name, s.borough
  FROM nyc_neighborhoods n, nyc_subway_stations s

--a simple where clause narrows this to 11058 rows
SELECT n.name, n.boroname,
       s.name, s.borough
  FROM nyc_neighborhoods n, nyc_subway_stations s
 WHERE n.boroname = s.borough
```

In the previous section, we explored spatial relationships using a two-step process: first we extracted a subway station point for ‘Broad St’; then, we used that point to ask further questions such as “what neighborhood is the ‘Broad St’ station in?”

Using a spatial join, we can answer the question in one step, retrieving information about the subway station and the neighborhood that contains it:

```
SELECT
    subways.name AS subway_name,
    neighborhoods.name AS neighborhood_name,
    neighborhoods.boroname AS borough
  FROM nyc_neighborhoods AS neighborhoods
  JOIN nyc_subway_stations AS subways
    ON ST_Contains(neighborhoods.geom, subways.geom)
 WHERE subways.name = 'Broad St';
```

subway_name	neighborhood_name	borough
Broad St	Financial District	Manhattan

We could have joined every subway station to its containing neighborhood, but in this case we wanted information about just one. Any function that provides a true/false relationship between two tables can be used to drive a spatial join, but the most commonly used ones are: **ST\_Intersects**, **ST\_Contains**, and **ST\_DWithin**.

### 4.7.1 Join and Summarize

The combination of a `JOIN` with a `GROUP BY` provides the kind of analysis that is usually done in a GIS system.

For example: **“What is the population and racial make-up of the neighborhoods of Manhattan?”** Here we have a question that combines information from about population from the census with the boundaries of neighborhoods, with a restriction to just one borough of Manhattan.

```
SELECT
    neighborhoods.name AS neighborhood_name,
    Sum(census.popn_total) AS population,
    100.0 * Sum(census.popn_white) / Sum(census.popn_total) AS white_pct,
    100.0 * Sum(census.popn_black) / Sum(census.popn_total) AS black_pct
```

```
FROM nyc_neighborhoods AS neighborhoods
JOIN nyc_census_blocks AS census
ON ST_Intersects(neighborhoods.geom, census.geom)
WHERE neighborhoods.boroname = 'Manhattan'
GROUP BY neighborhoods.name
ORDER BY white_pct DESC;
```

neighborhood_name	population	white_pct	black_pct
Carnegie Hill	18763	90.1	1.4
North Sutton Area	22460	87.6	1.6
West Village	26718	87.6	2.2
Upper East Side	203741	85.0	2.7
Soho	15436	84.6	2.2
Greenwich Village	57224	82.0	2.4
Central Park	46600	79.5	8.0
Tribeca	20908	79.1	3.5
Gramercy	104876	75.5	4.7
Murray Hill	29655	75.0	2.5
Chelsea	61340	74.8	6.4
Upper West Side	214761	74.6	9.2
Midtown	76840	72.6	5.2
Battery Park	17153	71.8	3.4
Financial District	34807	69.9	3.8
Clinton	32201	65.3	7.9
East Village	82266	63.3	8.8
Garment District	10539	55.2	7.1
Morningside Heights	42844	52.7	19.4
Little Italy	12568	49.0	1.8
Yorkville	58450	35.6	29.7
Inwood	50047	35.2	16.8
Washington Heights	169013	34.9	16.8
Lower East Side	96156	33.5	9.1
East Harlem	60576	26.4	40.4
Hamilton Heights	67432	23.9	35.8
Chinatown	16209	15.2	3.8
Harlem	134955	15.1	67.1

What's going on here? Notionally (the actual evaluation order is optimized under the covers by the database) this is what happens:

1. The JOIN clause creates a virtual table that includes columns from both the neighborhoods and census tables.
2. The WHERE clause filters our virtual table to just rows in Manhattan.
3. The remaining rows are grouped by the neighborhood name and fed through the aggregation function to `Sum()` the population values.
4. After a little arithmetic and formatting (e.g., GROUP BY, ORDER BY) on the final numbers, our query spits out the percentages.

---

**Note:** The JOIN clause combines two FROM items. By default, we are using an INNER JOIN, but there are four other types of joins. For further information see the [join\\_type](#) definition in the PostgreSQL documentation.

---

We can also use distance tests as a join key, to create summarized “all items within a radius” queries. Let’s explore the racial geography of New York using distance queries.

First, let’s get the baseline racial make-up of the city.

```
SELECT
  100.0 * Sum(popn_white) / Sum(popn_total) AS white_pct,
  100.0 * Sum(popn_black) / Sum(popn_total) AS black_pct,
  Sum(popn_total) AS popn_total
FROM nyc_census_blocks;
```

white_pct	black_pct	popn_total
44.0039500762811	25.5465789002416	8175032

So, of the 8M people in New York, about 44% are recorded as “white” and 26% are recorded as “black”.

Duke Ellington once sang that “You / must take the A-train / To / go to Sugar Hill way up in Harlem.” As we saw earlier, Harlem has far and away the highest African-American population in Manhattan (80.5%). Is the same true of Duke’s A-train?

First, note that the contents of the `nyc_subway_stations` table `routes` field is what we are interested in to find the A-train. The values in there are a little **complex**.

```
SELECT DISTINCT routes FROM nyc_subway_stations;
```

```
4,5
[null]
N,Q,R,W
J
B,M,Q,R
D,F,N,Q
J,M
```

---

**Note:** The `DISTINCT` keyword eliminates duplicate rows from the result. Without the `DISTINCT` keyword, the query above identifies 491 results instead of 73.

---

So to find the A-train, we will want any row in `routes` that has an ‘A’ in it. We can do this a number of ways, but today we will use the fact that `strpos(routes, 'A')` will return a non-zero number only if ‘A’ is in the `routes` field.

```
SELECT DISTINCT routes
FROM nyc_subway_stations AS subways
WHERE strpos(subways.routes, 'A') > 0;
```

```
A,C
A,B,C,D
A,C,E,L
A,C,F
A,B,C
A,S
A,C,E
A,C,G
A
```

Let's summarize the racial make-up of within 200 meters of the A-train line.

```
SELECT
  100.0 * Sum(popn_white) / Sum(popn_total) AS white_pct,
  100.0 * Sum(popn_black) / Sum(popn_total) AS black_pct,
  Sum(popn_total) AS popn_total
FROM nyc_census_blocks AS census
JOIN nyc_subway_stations AS subways
ON ST_DWithin(census.geom, subways.geom, 200)
WHERE strpos(subways.routes, 'A') > 0;
```

white_pct	black_pct	popn_total
45.5901255900202	22.0936235670937	189824

So the racial make-up along the A-train isn't radically different from the make-up of New York City as a whole.

### 4.7.2 Advanced Join

In the last section we saw that the A-train didn't serve a population that differed much from the racial make-up of the rest of the city. Are there any trains that have a non-average racial make-up?

To answer that question, we'll add another join to our query, so that we can simultaneously calculate the make-up of many subway lines at once. To do that, we'll need to create a new table that enumerates all the lines we want to summarize.

```
CREATE TABLE subway_lines ( route char(1) );
INSERT INTO subway_lines (route) VALUES
  ('A'), ('B'), ('C'), ('D'), ('E'), ('F'), ('G'),
  ('J'), ('L'), ('M'), ('N'), ('Q'), ('R'), ('S'),
  ('Z'), ('1'), ('2'), ('3'), ('4'), ('5'), ('6'),
  ('7');
```

Now we can join the table of subway lines onto our original query.

```
SELECT
  lines.route,
  100.0 * Sum(popn_white) / Sum(popn_total) AS white_pct,
  100.0 * Sum(popn_black) / Sum(popn_total) AS black_pct,
  Sum(popn_total) AS popn_total
FROM nyc_census_blocks AS census
JOIN nyc_subway_stations AS subways
ON ST_DWithin(census.geom, subways.geom, 200)
JOIN subway_lines AS lines
ON strpos(subways.routes, lines.route) > 0
GROUP BY lines.route
ORDER BY black_pct DESC;
```

route	white_pct	black_pct	popn_total
S	39.8	46.5	33301
3	42.7	42.1	223047
5	33.8	41.4	218919
2	39.3	38.4	291661
C	46.9	30.6	224411

4	37.6	27.4	174998
B	40.0	26.9	256583
A	45.6	22.1	189824
J	37.6	21.6	132861
Q	56.9	20.6	127112
Z	38.4	20.2	87131
D	39.5	19.4	234931
L	57.6	16.8	110118
G	49.6	16.1	135012
6	52.3	15.7	260240
1	59.1	11.3	327742
F	60.9	7.5	229439
M	56.5	6.4	174196
E	66.8	4.7	90958
R	58.5	4.0	196999
7	35.7	3.5	102401
N	59.7	3.5	147792

As before, the joins create a virtual table of all the possible combinations available within the constraints of the `JOIN ON` restrictions, and those rows are then fed into a `GROUP` summary. The spatial magic is in the `ST_DWithin` function, that ensures only census blocks close to the appropriate subway stations are included in the calculation.

### 4.7.3 Function List

`ST_Contains(geometry A, geometry B)`: Returns true if and only if no points of B lie in the exterior of A, and at least one point of the interior of B lies in the interior of A.

`ST_DWithin(geometry A, geometry B, radius)`: Returns true if the geometries are within the specified distance of one another.

`ST_Intersects(geometry A, geometry B)`: Returns TRUE if the Geometries/Geography “spatially intersect” - (share any portion of space) and FALSE if they don’t (they are Disjoint).

`round(v numeric, s integer)`: PostgreSQL math function that rounds to s decimal places

`strpos(string, substring)`: PostgreSQL string function that returns an integer location of a specified substring.

`sum(expression)`: PostgreSQL aggregate function that returns the sum of records in a set of records.

## 4.8 Spatial Joins Exercises

Here’s a reminder of some of the functions we have seen. Hint: they should be useful for the exercises!

- `sum(expression)`: aggregate to return a sum for a set of records
- `count (expression)`: aggregate to return the size of a set of records
- `ST_Area (geometry)` returns the area of the polygons
- `ST_AsText (geometry)` returns WKT text
- `ST_Contains (geometry A, geometry B)` returns the true if geometry A contains geometry B
- `ST_Distance (geometry A, geometry B)` returns the minimum distance between geometry A and geometry B

- **ST\_DWithin(geometry A, geometry B, radius)** returns the true if geometry A is within radius distance or less from geometry B
- **ST\_GeomFromText(text)** returns geometry
- **ST\_Intersects(geometry A, geometry B)** returns the true if geometry A intersects geometry B
- **ST\_Length(linestring)** returns the length of the linestring
- **ST\_Touches(geometry A, geometry B)** returns the true if the boundary of geometry A touches geometry B
- **ST\_Within(geometry A, geometry B)** returns the true if geometry A is within geometry B

Also remember the tables we have available:

- nyc\_census\_blocks
  - name, popn\_total, boroname, geom
- nyc\_streets
  - name, type, geom
- nyc\_subway\_stations
  - name, routes, geom
- nyc\_neighborhoods
  - name, boroname, geom

#### 4.8.1 Exercises

- **“What subway station is in ‘Little Italy’? What subway route is it on?”**

```
SELECT s.name, s.routes
FROM nyc_subway_stations AS s
JOIN nyc_neighborhoods AS n
ON ST_Contains(n.geom, s.geom)
WHERE n.name = 'Little Italy';
```

---

**Note:** Recall: the function AS is used to give a table another name by using an alias, which can make queries easier to read and write. In this case, s is an alias for nyc\_subway\_stations, n is an alias for nyc\_neighborhoods, s.name refers to the name column in the nyc\_subway\_stations table, etc.

---

name		routes
Spring St		6

- **“What are all the neighborhoods served by the 6-train?”** (Hint: The routes column in the nyc\_subway\_stations table has values like ‘B,D,6,V’ and ‘C,6’)

```
SELECT DISTINCT n.name, n.boroname
FROM nyc_subway_stations AS s
```

```
JOIN nyc_neighborhoods AS n
ON ST_Contains(n.geom, s.geom)
WHERE strpos(s.routes, '6') > 0;
```

name	boroname
Chinatown	Manhattan
East Harlem	Manhattan
Financial District	Manhattan
Gramercy	Manhattan
Greenwich Village	Manhattan
Hunts Point	The Bronx
Little Italy	Manhattan
Midtown	Manhattan
Mott Haven	The Bronx
Murray Hill	Manhattan
Parkchester	The Bronx
Soundview	The Bronx
South Bronx	The Bronx
Upper East Side	Manhattan
Yorkville	Manhattan

**Note:** We used the DISTINCT keyword to remove duplicate values from our result set where there were more than one subway station in a neighborhood.

---

- “After 9/11, the ‘Battery Park’ neighborhood was off limits for several days. How many people had to be evacuated?”

```
SELECT Sum(popn_total)
FROM nyc_neighborhoods AS n
JOIN nyc_census_blocks AS c
ON ST_Intersects(n.geom, c.geom)
WHERE n.name = 'Battery Park';
```

17153

- “What are the population density (people / km<sup>2</sup>) of the ‘Upper West Side’ and ‘Upper East Side’?” (Hint: There are 1000000 m<sup>2</sup> in one km<sup>2</sup>.)

```
SELECT
    n.name,
    Sum(c.popn_total) / (ST_Area(n.geom) / 1000000.0) AS popn_per_sqkm
FROM nyc_census_blocks AS c
JOIN nyc_neighborhoods AS n
ON ST_Intersects(c.geom, n.geom)
WHERE n.name = 'Upper West Side'
OR n.name = 'Upper East Side'
GROUP BY n.name, n.geom;
```

name	popn_per_sqkm
Upper East Side	48524.4877489857
Upper West Side	40152.4896080024

## 4.9 Spatial Indexing

Recall that spatial index is one of the three key features of a spatial database. Indexes are what make using a spatial database for large data sets possible. Without indexing, any search for a feature would require a “sequential scan” of every record in the database. Indexing speeds up searching by organizing the data into a search tree which can be quickly traversed to find a particular record.

Spatial indices are one of the greatest assets of PostGIS. In the previous example building spatial joins requires comparing whole tables with each other. This can get very costly: joining two tables of 10,000 records each without indexes would require 100,000,000 comparisons; with indexes the cost could be as low as 20,000 comparisons.

When we loaded the `nyc_census_blocks` table, the QGIS DB Manager was set to automatically create a spatial index called `sidx_nyc_census_blocks_geom`

To demonstrate how important indexes are for performance, let’s search `nyc_census_blocks` **without** our spatial index.

Our first step is to remove the index.

```
DROP INDEX sidx_nyc_census_blocks_geom;
```

---

**Note:** The `DROP INDEX` statement drops an existing index from the database system. For more information, see the PostgreSQL [documentation](#).

---

Now, watch the “Timing” meter at the lower right-hand corner of the pgAdmin query window and run the following. Our query searches through every single census block in order to identify the Broad Street entry.

```
SELECT blocks.blkid
  FROM nyc_census_blocks blocks
  JOIN nyc_subway_stations subways
    ON ST_Contains(blocks.geom, subways.geom)
   WHERE subways.name = 'Broad St';
```

```
blkid
-----
360610007001009
```

The `nyc_census_blocks` table is very small (only a few thousand records) so even without an index, the query only takes **103 ms** on my test computer.

Now add the spatial index back in and run the query again.

```
CREATE INDEX sidx_nyc_census_blocks_geom
  ON nyc_census_blocks
  USING GIST (geom);
```

---

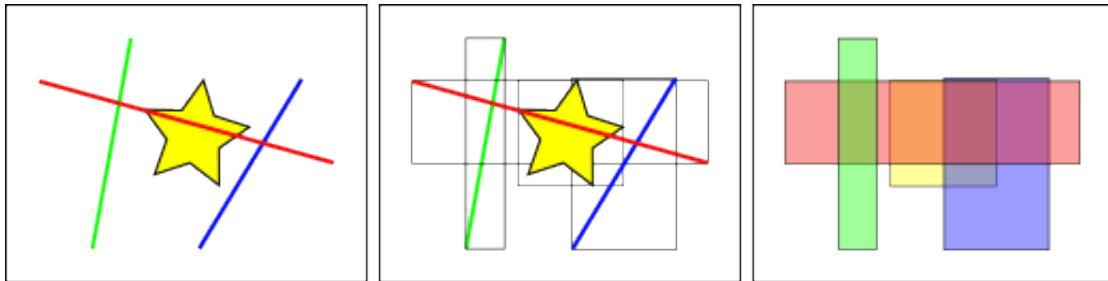
**Note:** The `USING GIST` clause tells PostgreSQL to use the generic index structure (GIST) when building the index. If you receive an error that looks like `ERROR: index row requires 11340 bytes, maximum size is 8191` when creating your index, you have likely neglected to add the `USING GIST` clause.

---

On my test computer the time drops to **66 ms**. The larger your table, the larger the relative speed improvement of an indexed query will be.

### 4.9.1 How Spatial Indexes Work

Standard database indexes create a hierarchical tree based on the values of the column being indexed. Spatial indexes are a little different – they are unable to index the geometric features themselves and instead index the bounding boxes of the features.



In the figure above, the number of lines that intersect the yellow star is **one**, the red line. But the bounding boxes of features that intersect the yellow box is **two**, the red and blue ones.

The way the database efficiently answers the question “what lines intersect the yellow star” is to first answer the question “what boxes intersect the yellow box” using the index (which is very fast) and then do an exact calculation of “what lines intersect the yellow star” **only for those features returned by the first test**.

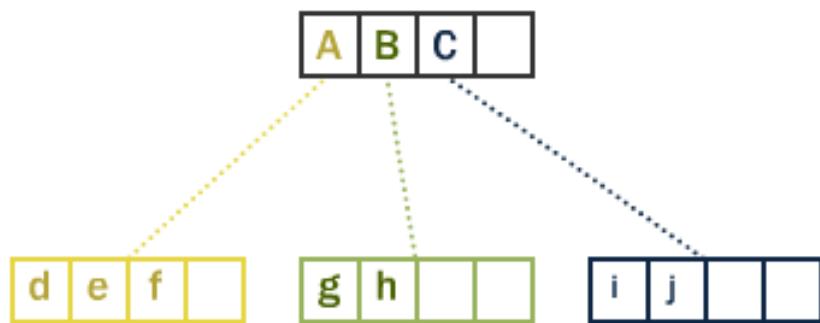
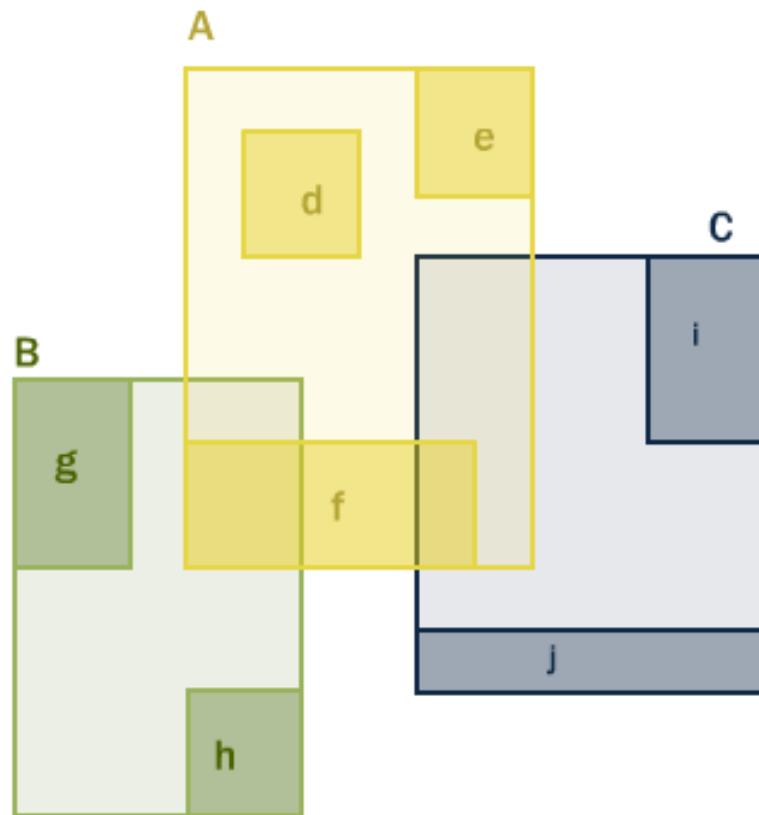
For a large table, this “two pass” system of evaluating the approximate index first, then carrying out an exact test can radically reduce the amount of calculations necessary to answer a query.

Both PostGIS and Oracle Spatial share the same “R-Tree”<sup>2</sup> spatial index structure. R-Trees break up data into rectangles, and sub-rectangles, and sub-sub rectangles, etc. It is a self-tuning index structure that automatically handles variable data density and object size.

---

<sup>2</sup> <http://postgis.org/support/rtree.pdf>

## R-tree Hierarchy



### 4.9.2 Index-Only Queries

Most of the commonly used functions in PostGIS (`ST_Contains`, `ST_Intersects`, `ST_DWithin`, etc) include an index filter automatically. But some functions (e.g., `ST_Relate`) do not include an index filter.

To do a bounding-box search using the index (and no filtering), make use of the `&&` operator. For geometries, the `&&` operator means “bounding boxes overlap or touch” in the same way that for numbers the `=` operator means “values are the same”.

Let's compare an index-only query for the population of the 'West Village' to a more exact query. Using `&&` our index-only query looks like the following:

```
SELECT Sum(popn_total)
FROM nyc_neighborhoods neighborhoods
JOIN nyc_census_blocks blocks
ON neighborhoods.geom && blocks.geom
WHERE neighborhoods.name = 'West Village';
```

```
49821
```

Now let's do the same query using the more exact `ST_Intersects` function.

```
SELECT Sum(popn_total)
FROM nyc_neighborhoods neighborhoods
JOIN nyc_census_blocks blocks
ON ST_Intersects(neighborhoods.geom, blocks.geom)
WHERE neighborhoods.name = 'West Village';
```

```
26718
```

A much lower answer! The first query summed up every block that intersected the neighborhood's bounding box; the second query only summed up those blocks that intersected the neighborhood itself.

### 4.9.3 Analyzing

The PostgreSQL query planner intelligently chooses when to use or not to use indexes to evaluate a query. Counter-intuitively, it is not always faster to do an index search: if the search is going to return every record in the table, traversing the index tree to get each record will actually be slower than just linearly reading the whole table from the start.

In order to figure out what situation it is dealing with (reading a small part of the table versus reading a large portion of the table), PostgreSQL keeps statistics about the distribution of data in each indexed table column. By default, PostgreSQL gathers statistics on a regular basis. However, if you dramatically change the make-up of your table within a short period of time, the statistics will not be up-to-date.

To ensure your statistics match your table contents, it is wise to run the `ANALYZE` command after bulk data loads and deletes in your tables. This forces the statistics system to gather data for all your indexed columns.

The `ANALYZE` command asks PostgreSQL to traverse the table and update its internal statistics used for query plan estimation (query plan analysis will be discussed later).

```
ANALYZE nyc_census_blocks;
```

### 4.9.4 Vacuuming

It's worth stressing that just creating an index is not enough to allow PostgreSQL to use it effectively. VACUUMing must be performed whenever a new index is created or after a large number of UPDATES, INSERTs or DELETEs are issued against a table. The `VACUUM` command asks PostgreSQL to reclaim any unused space in the table pages left by updates or deletes to records.

Vacuuming is so critical for the efficient running of the database that PostgreSQL provides an "autovacuum" option.

Enabled by default, autovacuum both vacuums (recovers space) and analyzes (updates statistics) on your tables at sensible intervals determined by the level of activity. While this is essential for highly transactional databases, it is not advisable to wait for an autovacuum run after adding indices or bulk-loading data. If a large batch update is performed, you should manually run VACUUM.

Vacuuming and analyzing the database can be performed separately as needed. Issuing VACUUM command will not update the database statistics; likewise issuing an ANALYZE command will not recover unused table rows. Both commands can be run against the entire database, a single table, or a single column.

```
VACUUM ANALYZE nyc_census_blocks;
```

#### 4.9.5 Function List

`geometry_a && geometry_b`: Returns TRUE if A's bounding box overlaps B's.

`geometry_a = geometry_b`: Returns TRUE if A's bounding box is the same as B's.

`ST_Intersects(geometry_a, geometry_b)`: Returns TRUE if the Geometries/Geography “spatially intersect” - (share any portion of space) and FALSE if they don't (they are Disjoint).

### 4.10 Projecting Data

The earth is not flat, and there is no simple way of putting it down on a flat paper map (or computer screen), so people have come up with all sorts of ingenious solutions, each with pros and cons. Some projections preserve area, so all objects have a relative size to each other; other projections preserve angles (conformal) like the Mercator projection; some projections try to find a good intermediate mix with only little distortion on several parameters. Common to all projections is that they transform the (spherical) world onto a flat Cartesian coordinate system, and which projection to choose depends on how you will be using the data.

We've already encountered projections when we [loaded our nyc data](#). (Recall that pesky SRID 26918). Sometimes, however, you need to transform and re-project between spatial reference systems. PostGIS includes built-in support for changing the projection of data, using the `ST_Transform(geometry, srid)` function. For managing the spatial reference identifiers on geometries, PostGIS provides the `ST_SRID(geometry)` and `ST_SetSRID(geometry, srid)` functions.

We can confirm the SRID of our data with the `ST_SRID` command:

```
SELECT ST_SRID(geom) FROM nyc_streets LIMIT 1;
```

```
26918
```

And what is definition of “26918”? As we saw in “[loading data section](#)”, the definition is contained in the `spatial_ref_sys` table. In fact, **two** definitions are there. The “well-known text” ([WKT](#)) definition is in the `srtext` column, and there is a second definition in “proj.4” format in the `proj4text` column.

```
SELECT * FROM spatial_ref_sys WHERE srid = 26918;
```

In fact, for the internal PostGIS re-projection calculations, it is the contents of the `proj4text` column that are used. For our 26918 projection, here is the proj.4 text:

```
SELECT proj4text FROM spatial_ref_sys WHERE srid = 26918;
```

```
+proj=utm +zone=18 +ellps=GRS80 +datum=NAD83 +units=m +no_defs
```

In practice, both the `srtext` and the `proj4text` columns are important: the `srtext` column is used by external programs like GeoServer, uDig, and FME and others; the `proj4text` column is used internally.

### 4.10.1 Comparing Data

Taken together, a coordinate and an SRID define a location on the globe. Without an SRID, a coordinate is just an abstract notion. A “Cartesian” coordinate plane is defined as a “flat” coordinate system placed on the surface of Earth. Because PostGIS functions work on such a plane, comparison operations require that both geometries be represented in the same SRID.

If you feed in geometries with differing SRIDs you will just get an error:

```
SELECT ST_Equals(
    ST_GeomFromText('POINT(0 0)', 4326),
    ST_GeomFromText('POINT(0 0)', 26918)
);
```

```
ERROR: Operation on mixed SRID geometries
CONTEXT: SQL function "st_equals" statement 1
```

---

**Note:** Be careful of getting too happy with using `ST_Transform` for on-the-fly conversion. Spatial indexes are built using SRID of the stored geometries. If comparison are done in a different SRID, spatial indexes are (often) not used. It is best practice to choose **one SRID** for all the tables in your database. Only use the transformation function when you are reading or writing data to external applications.

---

### 4.10.2 Transforming Data

If we return to our `proj4` definition for SRID 26918, we can see that our working projection is UTM (Universal Transverse Mercator) of zone 18, with meters as the unit of measurement.

```
+proj=utm +zone=18 +ellps=GRS80 +datum=NAD83 +units=m +no_defs
```

Let’s convert some data from our working projection to geographic coordinates – also known as “longitude/latitude”.

To convert data from one SRID to another, you must first verify that your geometry has a valid SRID. Since we have already confirmed a valid SRID, we next need the SRID of the projection to transform into. In other words, what is the SRID of geographic coordinates?

The most common SRID for geographic coordinates is 4326, which corresponds to “longitude/latitude on the WGS84 spheroid”. You can see the definition at the [spatialreference.org](http://spatialreference.org/ref/epsg/4326/) site:

<http://spatialreference.org/ref/epsg/4326/>

You can also pull the definitions from the `spatial_ref_sys` table:

```
SELECT srtext FROM spatial_ref_sys WHERE srid = 4326;
```

```
GEOGCS["WGS 84",
  DATUM["WGS_1984",
    SPHEROID["WGS 84", 6378137, 298.257223563, AUTHORITY["EPSG", "7030"]],
    AUTHORITY["EPSG", "6326"]],
  PRIMEM["Greenwich", 0, AUTHORITY["EPSG", "8901"]],
  UNIT["degree", 0.01745329251994328, AUTHORITY["EPSG", "9122"]],
  AUTHORITY["EPSG", "4326"]]
```

Let's convert the coordinates of the 'Broad St' subway station into geographics:

```
SELECT ST_AsText(ST_Transform(geom, 4326))
FROM nyc_subway_stations
WHERE name = 'Broad St';
```

```
POINT(-74.0106714688735 40.7071048155841)
```

If you load data or create a new geometry without specifying an SRID, the SRID value will be 0. Recall in geometries, that when we created our geometries table we didn't specify an SRID. If we query our database, we should expect all the nyc\_ tables to have an SRID of 26918, while the geometries table defaulted to an SRID of 0.

To view a table's SRID assignment, query the database's geometry\_columns table.

```
SELECT f_table_name AS name, srid
FROM geometry_columns;
```

name	srid
nyc_census_blocks	26918
nyc_neighborhoods	26918
nyc_streets	26918
nyc_subway_stations	26918
geometries	0

However, if you know what the SRID of the coordinates is supposed to be, you can set it post-facto, using **ST\_SetSRID** on the geometry. Then you will be able to transform the geometry into other systems.

```
SELECT ST_AsText(
  ST_Transform(
    ST_SetSRID(geom, 26918),
    4326)
)
FROM geometries;
```

### 4.10.3 Function List

**ST\_AsText**: Returns the Well-Known Text (WKT) representation of the geometry/geography without SRID metadata.

**ST\_SetSRID(geometry, srid)**: Sets the SRID on a geometry to a particular integer value.

**ST\_SRID(geometry)**: Returns the spatial reference identifier for the ST\_Geometry as defined in spatial\_ref\_sys table.

**ST\_Transform(geometry, srid)**: Returns a new geometry with its coordinates transformed to the SRID referenced by the integer parameter.

## 4.11 Projection Exercises

Here's a reminder of some of the functions we have seen. Hint: they should be useful for the exercises!

- `sum(expression)` aggregate to return a sum for a set of records
- `ST_Length(linestring)` returns the length of the linestring
- `ST_SRID(geometry, srid)` returns the SRID of the geometry
- `ST_Transform(geometry, srid)` converts geometries into different spatial reference systems
- `ST_GeomFromText(text)` returns geometry
- `ST_AsText(geometry)` returns WKT text
- `ST_AsGML(geometry)` returns GML text

Remember the online resources that are available to you:

- <http://spatialreference.org>
- <http://prj2epsg.org>

Also remember the tables we have available:

- `nyc_census_blocks`
  - name, popn\_total, boroname, geom
- `nyc_streets`
  - name, type, geom
- `nyc_subway_stations`
  - name, geom
- `nyc_neighborhoods`
  - name, boroname, geom

### 4.11.1 Exercises

- “What is the length of all streets in New York, as measured in UTM 18?”

```
SELECT Sum(ST_Length(geom))  
FROM nyc_streets;
```

```
10418904.7172
```

- “What is the WKT definition of SRID 2831?”

```
SELECT srtext FROM spatial_ref_sys  
WHERE SRID = 2831;
```

Or, via [prj2epsg](http://prj2epsg.org)

```
PROJCS["NAD83(HARN) / New York Long Island",  
    GEOGCS["NAD83(HARN)",  
        DATUM["NAD83 (High Accuracy Regional Network)",
```

```

SPHEROID["GRS 1980", 6378137.0, 298.257222101,
    AUTHORITY["EPSG", "7019"]],
TOWGS84[-0.991, 1.9072, 0.5129, 0.0257899075194932, -0.,
    →009650098960270402, -0.011659943232342112, 0.0],
    AUTHORITY["EPSG", "6152"]],
PRIMEM["Greenwich", 0.0,
    AUTHORITY["EPSG", "8901"]],
UNIT["degree", 0.017453292519943295],
AXIS["Geodetic longitude", EAST],
AXIS["Geodetic latitude", NORTH],
AUTHORITY["EPSG", "4152"]],
PROJECTION["Lambert Conic Conformal (2SP)",
    AUTHORITY["EPSG", "9802"]],
PARAMETER["central_meridian", -74.0],
PARAMETER["latitude_of_origin", 40.16666666666664],
PARAMETER["standard_parallel_1", 41.03333333333333],
PARAMETER["false_easting", 300000.0],
PARAMETER["false_northing", 0.0],
PARAMETER["scale_factor", 1.0],
PARAMETER["standard_parallel_2", 40.66666666666664],
UNIT["m", 1.0],
AXIS["Easting", EAST],
AXIS["Northing", NORTH],
AUTHORITY["EPSG", "2831"]]

```

- “What is the length of all streets in New York, as measured in SRID 2831?”

```
SELECT Sum(ST_Length(ST_Transform(geom, 2831)))
FROM nyc_streets;
```

```
10421993.7063767
```

**Note:** The difference between the UTM 18 and the State Plane Long Island measurements is  $(10421993 - 10418904)/10418904$ , or 0.02%. Calculated on the spheroid using geography the total street length is 10421999, which is closer to the State Plane value. This is not surprising, since the State Plane Long Island projection is precisely calibrated for a very small area (New York City) while UTM 18 has to provide reasonable results for a large regional area.

- “What is the KML representation of the point at ‘Broad St’ subway station?”

```
SELECT ST_AsKML(geom)
FROM nyc_subway_stations
WHERE name = 'Broad St';
```

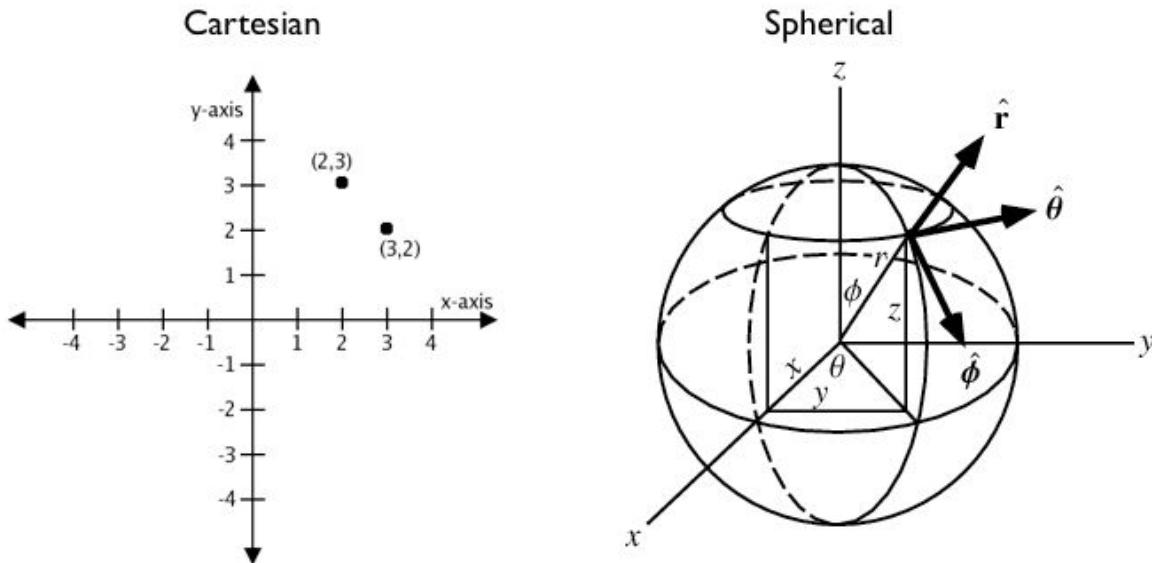
```
<Point>
<coordinates>
-74.010671468873412,40.707104815587613
</coordinates>
</Point>
```

Hey! The coordinates are in geographics even though we didn’t call **ST\_Transform**, why? Because the KML standard dictates that all coordinates *must* be in geographics (EPSG:4326, in fact) so the **ST\_AsKML** function does the transformation automatically.

## 4.12 Geography

It is very common to have data in which the coordinates are “geographics” or “latitude/longitude”.

Unlike coordinates in Mercator, UTM, or Stateplane, geographic coordinates are **not Cartesian coordinates**. Geographic coordinates do not represent a linear distance from an origin as plotted on a plane. Rather, these **spherical coordinates** describe angular coordinates on a globe. In spherical coordinates a point is specified by the angle of rotation from a reference meridian (longitude), and the angle from the equator (latitude).



You can treat geographic coordinates as approximate Cartesian coordinates and continue to do spatial calculations. However, measurements of distance, length and area will be nonsensical. Since spherical coordinates measure **angular** distance, the units are in “degrees.” Further, the approximate results from indexes and true/false tests like intersects and contains can become terribly wrong. The distance between points get larger as problem areas like the poles or the international dateline are approached.

For example, here are the coordinates of Los Angeles and Paris.

- Los Angeles: POINT (-118.4079 33.9434)
- Paris: POINT (2.3490 48.8533)

The following calculates the distance between Los Angeles and Paris using the standard PostGIS Cartesian **ST\_Distance** (**geometry, geometry**). Note that the SRID of 4326 declares a geographic spatial reference system.

```
SELECT ST_Distance(
    ST_GeometryFromText('POINT(-118.4079 33.9434)', 4326), -- Los Angeles
    ST_GeometryFromText('POINT(2.3490 48.8533)', 4326)      -- Paris (CDG)
);
```

```
121.898285970107
```

Aha! 121! But, what does that mean?

The units for spatial reference 4326 are degrees. So our answer is 121 degrees. But (again), what does that mean?

On a sphere, the size of one “degree square” is quite variable, becoming smaller as you move away from the equator. Think of the meridians (vertical lines) on the globe getting closer to each other as you go towards the poles. So, a distance of 121 degrees doesn’t *mean* anything. It is a nonsense number.

In order to calculate a meaningful distance, we must treat geographic coordinates not as approximate Cartesian coordinates but rather as true spherical coordinates. We must measure the distances between points as true paths over a sphere – a portion of a great circle.

Starting with version 1.5, PostGIS provides this functionality through the `geography` type.

**Note:** Different spatial databases have different approaches for “handling geographics”

- Oracle attempts to paper over the differences by transparently doing geographic calculations when the SRID is geographic.
- SQL Server uses two spatial types, “`STGeometry`” for Cartesian data and “`STGeography`” for geographics.
- Informix Spatial is a pure Cartesian extension to Informix, while Informix Geodetic is a pure geographic extension.
- Similar to SQL Server, PostGIS uses two types, “`geometry`” and “`geography`”.

Using the `geography` instead of `geometry` type, let’s try again to measure the distance between Los Angeles and Paris. Instead of `ST_GeometryFromText(text)`, we will use `ST_GeographyFromText(text)`.

```
SELECT ST_Distance(
    ST_GeographyFromText('POINT(-118.4079 33.9434)'), -- Los Angeles (LAX)
    ST_GeographyFromText('POINT(2.5559 49.0083)')        -- Paris (CDG)
);
```

9124665.27317673

A big number! All return values from `geography` calculations are in meters, so our answer is 9124km.

Older versions of PostGIS supported very basic calculations over the sphere using the `ST_Distance_Spheroid(point, point, measurement)` function. However, `ST_Distance_Spheroid` is substantially limited. The function only works on points and provides no support for indexing across the poles or international dateline.

The need to support non-point geometries becomes very clear when posing a question like “How close will a flight from Los Angeles to Paris come to Iceland?”



Working with geographic coordinates on a Cartesian plane (the purple line) yields a *very* wrong answer indeed! Using great circle routes (the red lines) gives the right answer. If we convert our LAX-CDG flight into a line string and calculate the distance to a point in Iceland using geography we'll get the right answer (recall) in meters.

```
SELECT ST_Distance(
    ST_GeographyFromText ('LINESTRING (-118.4079 33.9434, 2.5559 49.0083)'), --
    ↪ LAX-CDG
    ST_GeographyFromText ('POINT (-22.6056 63.9850)') --
    ↪ Iceland (KEF)
);
```

```
502454.90667289
```

So the closest approach to Iceland (as measured from its international airport) on the LAX-CDG route is a relatively small 502km.

The Cartesian approach to handling geographic coordinates breaks down entirely for features that cross the international dateline. The shortest great-circle route from Los Angeles to Tokyo crosses the Pacific Ocean. The shortest Cartesian route crosses the Atlantic and Indian Oceans.



```
SELECT ST_Distance(
    ST_GeometryFromText('Point(-118.4079 33.9434)'), -- LAX
    ST_GeometryFromText('Point(139.733 35.567)')      -- NRT (Tokyo/Narita)
) AS geometry_distance,
ST_Distance(
    ST_GeographyFromText('Point(-118.4079 33.9434)'), -- LAX
    ST_GeographyFromText('Point(139.733 35.567)')      -- NRT (Tokyo/Narita)
) AS geography_distance;
```

geometry_distance	geography_distance
258.146005837336	8833954.77277118

### 4.12.1 Using Geography

In order to load geometry data into a geography table, the geometry first needs to be projected into EPSG:4326 (longitude/latitude), then it needs to be changed into geography. The **ST\_Transform(geometry, srid)** function converts coordinates to geographics and the **Geography(geometry)** function “casts” them from geometry to geography.

```
CREATE TABLE nyc_subway_stations_geog AS
SELECT
    Geography(ST_Transform(geom, 4326)) AS geog,
    name,
    routes
FROM nyc_subway_stations;
```

Building a spatial index on a geography table is exactly the same as for geometry:

```
CREATE INDEX sidx_nyc_subway_stations_geog
ON nyc_subway_stations_geog USING GIST (geog);
```

The difference is under the covers: the geography index will correctly handle queries that cover the poles or the **international date-line**, while the geometry one will not.

There are only a small number of native functions for the geography type:

- **ST\_AsText(geography)** returns text
- **ST\_GeographyFromText(text)** returns geography
- **ST\_AsBinary(geography)** returns bytea
- **ST\_GeogFromWKB(bytea)** returns geography
- **ST\_AsSVG(geography)** returns text
- **ST\_AsGML(geography)** returns text
- **ST\_AsKML(geography)** returns text
- **ST\_AsGeoJson(geography)** returns text
- **ST\_Distance(geography, geography)** returns double
- **ST\_DWithin(geography, geography, float8)** returns boolean
- **ST\_Area(geography)** returns double
- **ST\_Length(geography)** returns double

- **ST\_Covers (geography, geography)** returns boolean
- **ST\_CoveredBy (geography, geography)** returns boolean
- **ST\_Intersects (geography, geography)** returns boolean
- **ST\_Buffer (geography, float8)** returns geography<sup>3</sup>
- **ST\_Intersection (geography, geography)** returns geography<sup>3</sup>

### 4.12.2 Creating a Geography Table

The SQL for creating a new table with a geography column is much like that for creating a geometry table. For example:

```
CREATE TABLE airports (
    code VARCHAR(3),
    geog GEOGRAPHY(Point)
);

INSERT INTO airports VALUES ('LAX', 'POINT(-118.4079 33.9434)');
INSERT INTO airports VALUES ('CDG', 'POINT(2.5559 49.0083)');
INSERT INTO airports VALUES ('KEF', 'POINT(-22.6056 63.9850)');
```

In the table definition, the `GEOGRAPHY (Point)` specifies our airport data type as points. The new geography fields don't get registered in the `geometry_columns` view. Instead, they are registered in a view called `geography_columns`.

```
SELECT * FROM geography_columns;
```

f_table_name	f_geography_column	srid	type
nyc_subway_stations_geog	geog	0	Geometry
airports	geog	4326	Point

---

**Note:** Some columns were omitted from the above output.

---

### 4.12.3 Casting to Geometry

While the basic functions for geography types can handle many use cases, there are times when you might need access to other functions only supported by the geometry type. Fortunately, you can convert objects back and forth from geography to geometry.

The PostgreSQL syntax convention for casting is to append `::typename` to the end of the value you wish to cast. So, `2::text` will convert a numeric two to a text string '2'. And '`POINT(0 0) ::geometry`' will convert the text representation of point into a geometry point.

The `ST_X (point)` function only supports the geometry type. How can we read the X coordinate from our geographies?

---

<sup>3</sup> The buffer and intersection functions are actually wrappers on top of a cast to geometry, and are not carried out natively in spherical coordinates. As a result, they may fail to return correct results for objects with very large extents that cannot be cleanly converted to a planar representation.

For example, the `ST_Buffer (geography, distance)` function transforms the geography object into a “best” projection, buffers it, and then transforms it back to geographics. If there is no “best” projection (the object is too large), the operation can fail or return a malformed buffer.

```
SELECT code, ST_X(geog::geometry) AS longitude FROM airports;
```

code	longitude
LAX	-118.4079
CDG	2.5559
KEF	-22.6056

By appending `::geometry` to our geography value, we convert the object to a geometry with an SRID of 4326. From there we can use as many geometry functions as strike our fancy. But, remember – now that our object is a geometry, the coordinates will be interpreted as Cartesian coordinates, not spherical ones.

#### 4.12.4 Why (Not) Use Geography

Geographics are universally accepted coordinates – everyone understands what latitude/longitude mean, but very few people understand what UTM coordinates mean. Why not use geography all the time?

- First, as noted earlier, there are far fewer functions available (right now) that directly support the geography type. You may spend a lot of time working around geography type limitations.
- Second, the calculations on a sphere are computationally far more expensive than Cartesian calculations. For example, the Cartesian formula for distance (Pythagoras) involves one call to `sqrt()`. The spherical formula for distance (Haversine) involves two `sqrt()` calls, an `arctan()` call, four `sin()` calls and two `cos()` calls. Trigonometric functions are very costly, and spherical calculations involve a lot of them.

The conclusion?

If your data is geographically compact (contained within a state, county or city), use the geometry type with a Cartesian projection that makes sense with your data. See the <http://spatialreference.org> site and type in the name of your region for a selection of possible reference systems.

If you need to measure distance with a dataset that is geographically dispersed (covering much of the world), use the geography type. The application complexity you save by working in geography will offset any performance issues. And casting to `geometry` can offset most functionality limitations.

#### 4.12.5 Function List

`ST_Distance(geometry, geometry)`: For geometry type Returns the 2-dimensional Cartesian minimum distance (based on spatial ref) between two geometries in projected units. For geography type defaults to return spheroidal minimum distance between two geographies in meters.

`ST_GeographyFromText(text)`: Returns a specified geography value from Well-Known Text representation or extended (WKT).

`ST_Transform(geometry, srid)`: Returns a new geometry with its coordinates transformed to the SRID referenced by the integer parameter.

`ST_X(point)`: Returns the X coordinate of the point, or NULL if not available. Input must be a point.

### 4.13 Geometry Constructing Functions

All the functions we have seen so far work with geometries “as they are” and returns

- analyses of the objects (`ST_Length(geometry)`, `ST_Area(geometry)`),

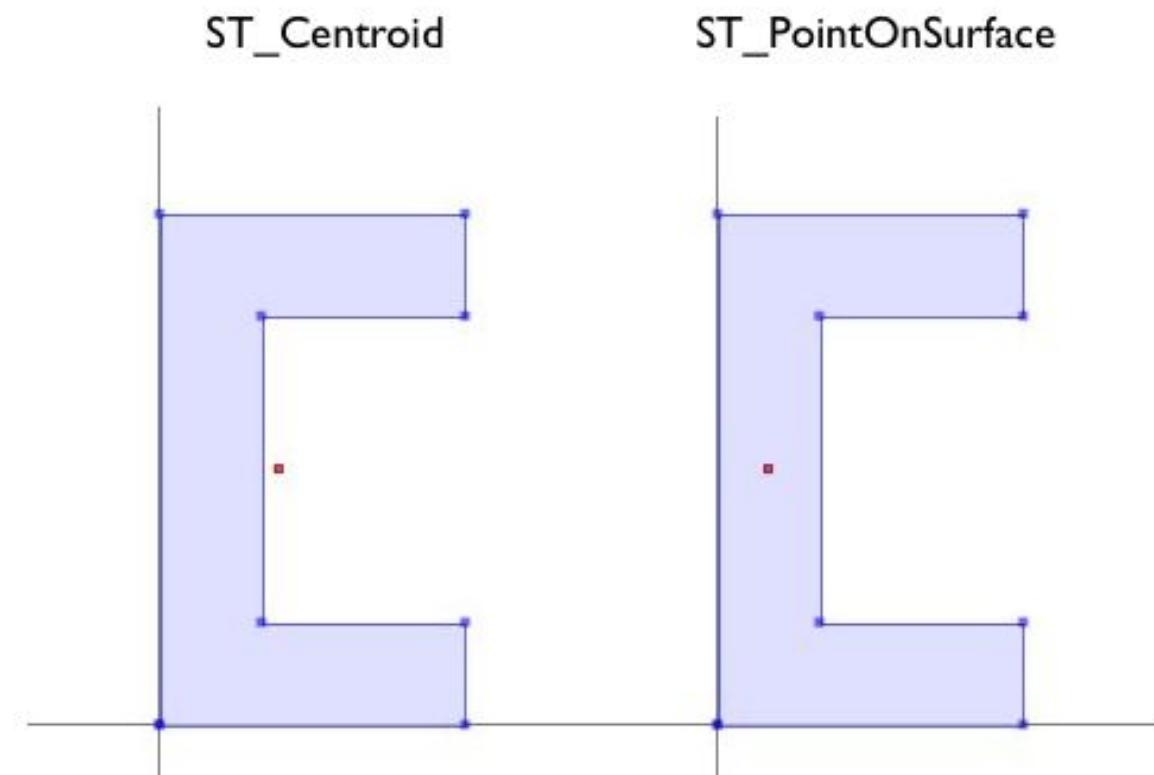
- serializations of the objects (**ST\_AsText(geometry)**, **ST\_AsGML(geometry)**),
- parts of the object (**ST\_RingN(geometry, n)**) or
- true/false tests (**ST\_Contains(geometry, geometry)**, **ST\_Intersects(geometry, geometry)**).

“Geometry constructing functions” take geometries as inputs and output new shapes.

### 4.13.1 ST\_Centroid / ST\_PointOnSurface

A common need when composing a spatial query is to replace a polygon feature with a point representation of the feature. This is useful for spatial joins (as discussed in polypolyjoins) because using **ST\_Intersects(geometry, geometry)** on two polygon layers often results in double-counting: a polygon on a boundary will intersect an object on both sides; replacing it with a point forces it to be on one side or the other, not both.

- **ST\_Centroid(geometry)** returns a point that is approximately on the center of mass of the input argument. This simple calculation is very fast, but sometimes not desirable, because the returned point is not necessarily in the feature itself. If the input feature has a convexity (imagine the letter ‘C’) the returned centroid might not be in the interior of the feature.
- **ST\_PointOnSurface(geometry)** returns a point that is guaranteed to be inside the input argument. It is substantially more computationally expensive than the centroid operation.



### 4.13.2 ST\_Buffer

The buffering operation is common in GIS workflows, and is also available in PostGIS. **ST\_Buffer(geometry, distance)** takes in a buffer distance and geometry type and outputs a polygon with a boundary the buffer distance away from the input geometry.

## ST\_Buffer



Buffering a point



Buffering a multipoint

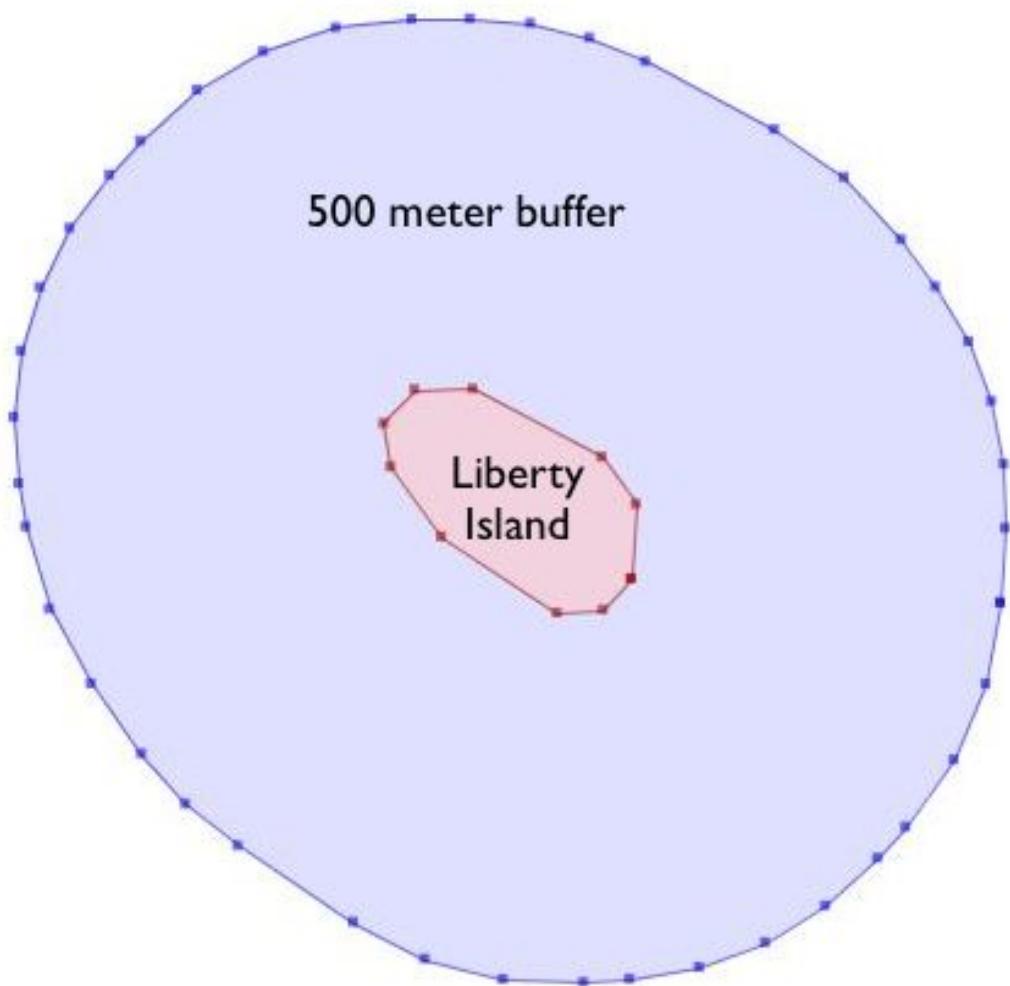


Buffering a linestring

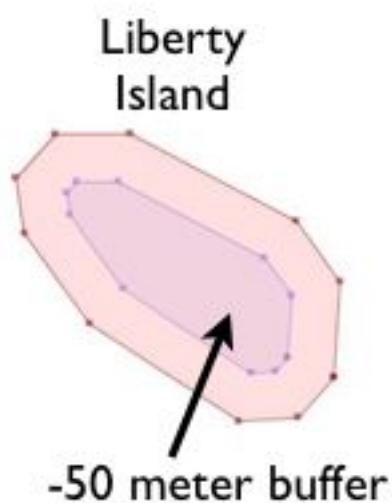
Buffering a polygon  
with one interior ring

For example, if the US Park Service wanted to enforce a marine traffic zone around Liberty Island, they might build a 500 meter buffer polygon around the island. Liberty Island is a single census block in our `nyc_census_blocks` table, so we can easily extract and buffer it.

```
-- Make a new table with a Liberty Island 500m buffer zone
CREATE TABLE liberty_island_zone AS
SELECT ST_Buffer(geom, 500) ::geometry(Polygon, 26918) AS geom
FROM nyc_census_blocks
WHERE blkid = '360610001001001';
```



The **ST\_Buffer** function also accepts negative distances and builds inscribed polygons within polygonal inputs. For lines and points you will just get an empty return.



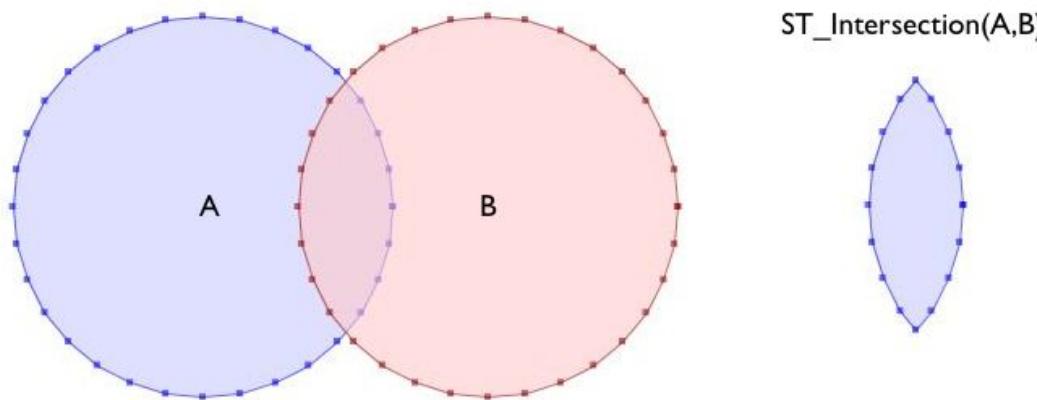
### 4.13.3 ST\_Intersection

Another classic GIS operation – the “overlay” – creates a new coverage by calculating the intersection of two superimposed polygons. The resultant has the property that any polygon in either of the parents can be built by merging polygons in the resultant.

The **ST\_Intersection(geometry A, geometry B)** function returns the spatial area (or line, or point) that both arguments have in common. If the arguments are disjoint, the function returns an empty geometry.

```
-- What is the area these two circles have in common?
-- Using ST_Buffer to make the circles!
```

```
SELECT ST_AsText(ST_Intersection(
    ST_Buffer('POINT(0 0)', 2),
    ST_Buffer('POINT(3 0)', 2)
));
```



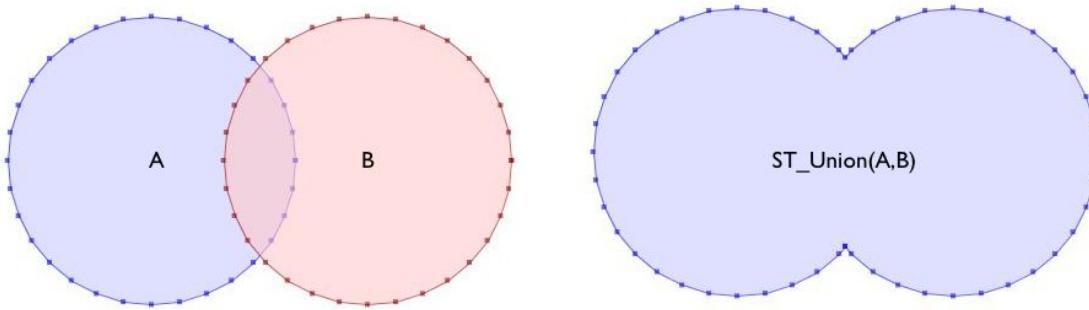
### 4.13.4 ST\_Union

In the previous example we intersected geometries, creating a new geometry that had lines from both the inputs. The **ST\_Union** does the reverse; it takes inputs and removes common lines. There are two forms of the **ST\_Union** function:

- **ST\_Union(geometry, geometry)**: A two-argument version that takes in two geometries and returns the merged union. For example, our two-circle example from the previous section looks like this when you replace the intersection with a union.

```
-- What is the total area these two circles cover?
-- Using ST_Buffer to make the circles!
```

```
SELECT ST_AsText(ST_Union(
    ST_Buffer('POINT(0 0)', 2),
    ST_Buffer('POINT(3 0)', 2)
));
```



- **ST\_Union([geometry])**: An aggregate version that takes in a set of geometries and returns the merged geometry for the entire group. The aggregate ST\_Union can be used with the GROUP BY SQL statement to create carefully merged subsets of basic geometries. It is very powerful,

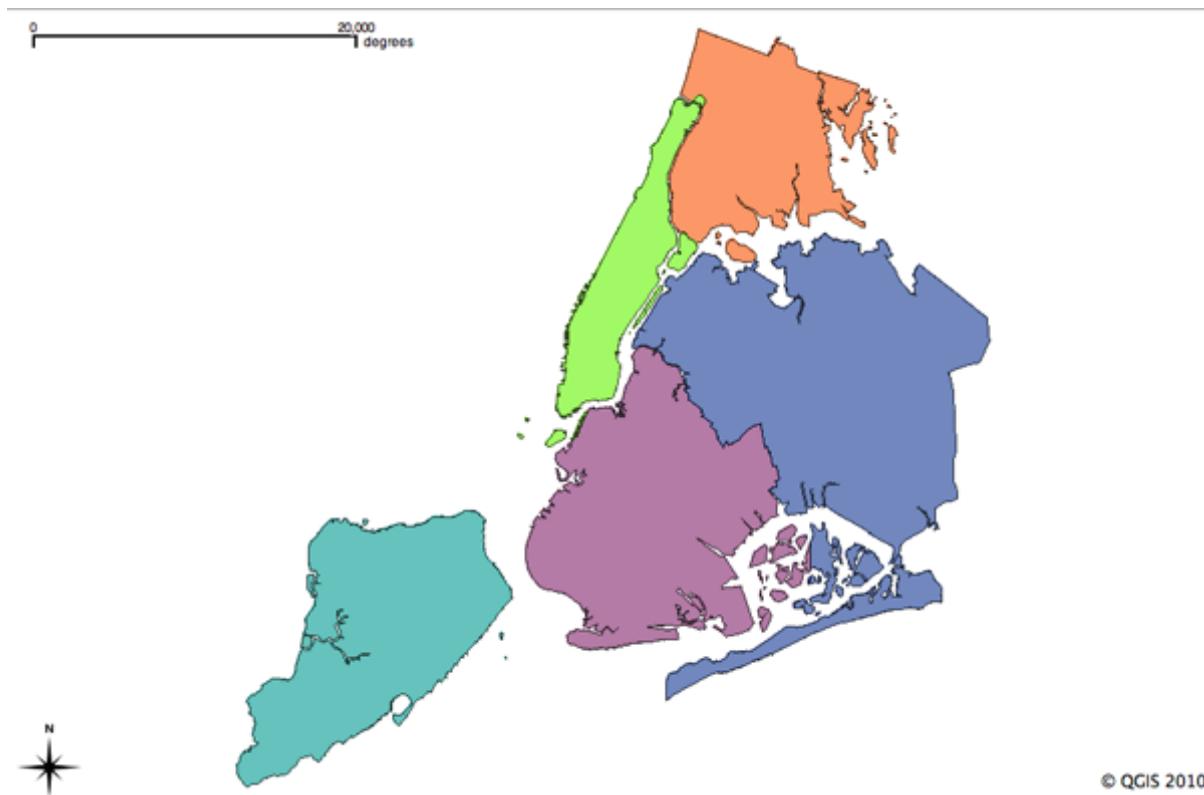
As an example of **ST\_Union** aggregation, consider our `nyc_census_blocks` table. Census geography is carefully constructed so that larger geographies can be built up from smaller ones. So, we can create a census tracts map by merging the blocks that form each tract (as we do later in creatingtractstable). Or, we can create a county map by merging blocks that fall within each county.

To carry out the merge, note that the unique key `blkid` actually embeds information about the higher level geographies. Here are the parts of the key for Liberty Island we used earlier:

```
360610001001001 = 36 061 000100 1 001
36      = State of New York
061     = New York County (Manhattan)
000100 = Census Tract
1       = Census Block Group
001     = Census Block
```

So, we can create a county map by merging all geometries that share the same first 5 digits of their `blkid`. Be patient; this is computationally expensive and can take a minute or two.

```
-- Create a nyc_census_counties table by merging census blocks
CREATE TABLE nyc_census_counties AS
SELECT
    ST_Union(geom)::Geometry(MultiPolygon,26918) AS geom,
    SubStr(blkid,1,5) AS countyid
FROM nyc_census_blocks
GROUP BY countyid;
```



An area test can confirm that our union operation did not lose any geometry. First, we calculate the area of each individual census block, and sum those areas grouping by census county id.

```
SELECT SubStr(blkid,1,5) AS countyid, Sum(ST_Area(geom)) AS area
FROM nyc_census_blocks
GROUP BY countyid
ORDER BY countyid;
```

countyid	area
36005	110196022.906506
36047	181927497.678368
36061	59091860.6261323
36081	283194473.613692
36085	150758328.111199

Then we calculate the area of each of our new county polygons from the county table:

```
SELECT countyid, ST_Area(geom) AS area
FROM nyc_census_counties;
```

countyid	area
36005	110196022.906507
36047	181927497.678367
36061	59091860.6261324
36081	283194473.593646
36085	150758328.111199

The same answer! We have successfully built an NYC county table from our census blocks data.

### 4.13.5 Function List

**ST\_AsText(text):** Returns the Well-Known Text (WKT) representation of the geometry/geography without SRID metadata.

**ST\_Buffer(geometry, distance):** For geometry: Returns a geometry that represents all points whose distance from this Geometry is less than or equal to distance. Calculations are in the Spatial Reference System of this Geometry. For geography: Uses a planar transform wrapper.

**ST\_Intersection(geometry A, geometry B):** Returns a geometry that represents the shared portion of geomA and geomB. The geography implementation does a transform to geometry to do the intersection and then transform back to WGS84.

**ST\_Union():** Returns a geometry that represents the point set union of the Geometries.

**substring(string [from int] [for int]):** PostgreSQL string function to extract substring matching SQL regular expression.

**sum(expression):** PostgreSQL aggregate function that returns the sum of records in a set of records.

## ADVANCED

### 5.1 More Spatial Joins

In the last section we saw the `ST_Centroid(geometry)` and `ST_Union([geometry])` functions, and some simple examples. In this section we will do some more elaborate things with them.

#### 5.1.1 Creating a Census Tracts Table

In the workshop `\data\` directory, is a file that includes attribute data, but no geometry, `nyc_census_sociodata.sql`. The table includes interesting socioeconomic data about New York: commute times, incomes, and education attainment. There is just one problem. The data are summarized by “census tract” and we have no census tract spatial data!

In this section we will

- Load the `nyc_census_sociodata.sql` table
- Create a spatial table for census tracts
- Join the attribute data to the spatial data
- Carry out some analysis using our new data

#### Loading `nyc_census_sociodata.sql`

1. Open the SQL query window in PgAdmin
2. Select **File->Open** from the menu and browse to the `nyc_census_sociodata.sql` file
3. (*Alternative*) In case your postgres user does not have access to access the files in the directory you can use a text editor to open the `nyc_census_sociodata.sql` file and copy its content into the pgAdmin query window
4. Press the “Run Query” button
5. If you press the “Refresh” button in pgAdmin, the list of tables should now include at `nyc_census_sociodata` table

#### Creating a Census Tracts Table

As we saw in the previous section, we can build up higher level geometries from the census block by summarizing on substrings of the `blkid` key. In order to get census tracts, we need to summarize grouping on the first 11 characters of the `blkid`.

```
360610001001001 = 36 061 000100 1 001  
  
36      = State of New York  
061     = New York County (Manhattan)  
000100  = Census Tract  
1       = Census Block Group  
001     = Census Block
```

Create the new table using the **ST\_Union** aggregate:

```
-- Make the tracts table  
CREATE TABLE nyc_census_tract_geoms AS  
SELECT  
    ST_Union(geom) AS geom,  
    SubStr(blkid,1,11) AS tractid  
FROM nyc_census_blocks  
GROUP BY tractid;  
  
-- Index the tractid  
CREATE INDEX idx_nyc_census_tract_geoms_tractid  
    ON nyc_census_tract_geoms (tractid);
```

### Join the Attributes to the Spatial Data

Join the table of tract geometries to the table of tract attributes with a standard attribute join

```
-- Make the tracts table  
CREATE TABLE nyc_census_tracts AS  
SELECT  
    g.geom,  
    a.*  
FROM nyc_census_tract_geoms g  
JOIN nyc_census_sociodata a  
ON g.tractid = a.tractid;  
  
-- Index the geometries  
CREATE INDEX sidx_nyc_census_tract  
    ON nyc_census_tracts USING GIST (geom);
```

### Answer an Interesting Question

Answer an interesting question! “List top 10 New York neighborhoods ordered by the proportion of people who have graduate degrees.”

```
SELECT  
    100.0 * Sum(t.edu_graduate_dipl) / Sum(t.edu_total) AS graduate_pct,  
    n.name, n.boroname  
FROM nyc_neighborhoods n  
JOIN nyc_census_tracts t  
ON ST_Intersects(n.geom, t.geom)  
WHERE t.edu_total > 0  
GROUP BY n.name, n.boroname  
ORDER BY graduate_pct DESC  
LIMIT 10;
```

We sum up the statistics we are interested, then divide them together at the end. In order to avoid divide-by-zero errors, we don't bother bringing in tracts that have a population count of zero.

graduate_pct	name	boroname
47.6	Carnegie Hill	Manhattan
42.2	Upper West Side	Manhattan
41.1	Battery Park	Manhattan
39.6	Flatbush	Brooklyn
39.3	Tribeca	Manhattan
39.2	North Sutton Area	Manhattan
38.7	Greenwich Village	Manhattan
38.6	Upper East Side	Manhattan
37.9	Murray Hill	Manhattan
37.4	Central Park	Manhattan

**Note:** New York geographers will be wondering at the presence of “Flatbush” in this list of over-educated neighborhoods. The answer is discussed in the next section.

### 5.1.2 Polygon/Polygon Joins

In our interesting query (in `interestingquestion`) we used the `ST_Intersects(geometry_a, geometry_b)` function to determine which census tract polygons to include in each neighborhood summary. Which leads to the question: what if a tract falls on the border between two neighborhoods? It will intersect both, and so will be included in the summary statistics for **both**.



To avoid this kind of double counting there are two methods:

- The simple method is to ensure that each tract only falls in **one** summary area (using

`ST_Centroid(geometry)`

- The complex method is to divide crossing tracts at the borders (using `ST_Intersection(geometry, geometry)`)

Here is an example of using the simple method to avoid double counting in our graduate education query:

```
SELECT
  100.0 * Sum(t.edu_graduate_dipl) / Sum(t.edu_total) AS graduate_pct,
  n.name, n.boroname
FROM nyc_neighborhoods n
JOIN nyc_census_tracts t
ON ST_Contains(n.geom, ST_Centroid(t.geom))
WHERE t.edu_total > 0
GROUP BY n.name, n.boroname
ORDER BY graduate_pct DESC
LIMIT 10;
```

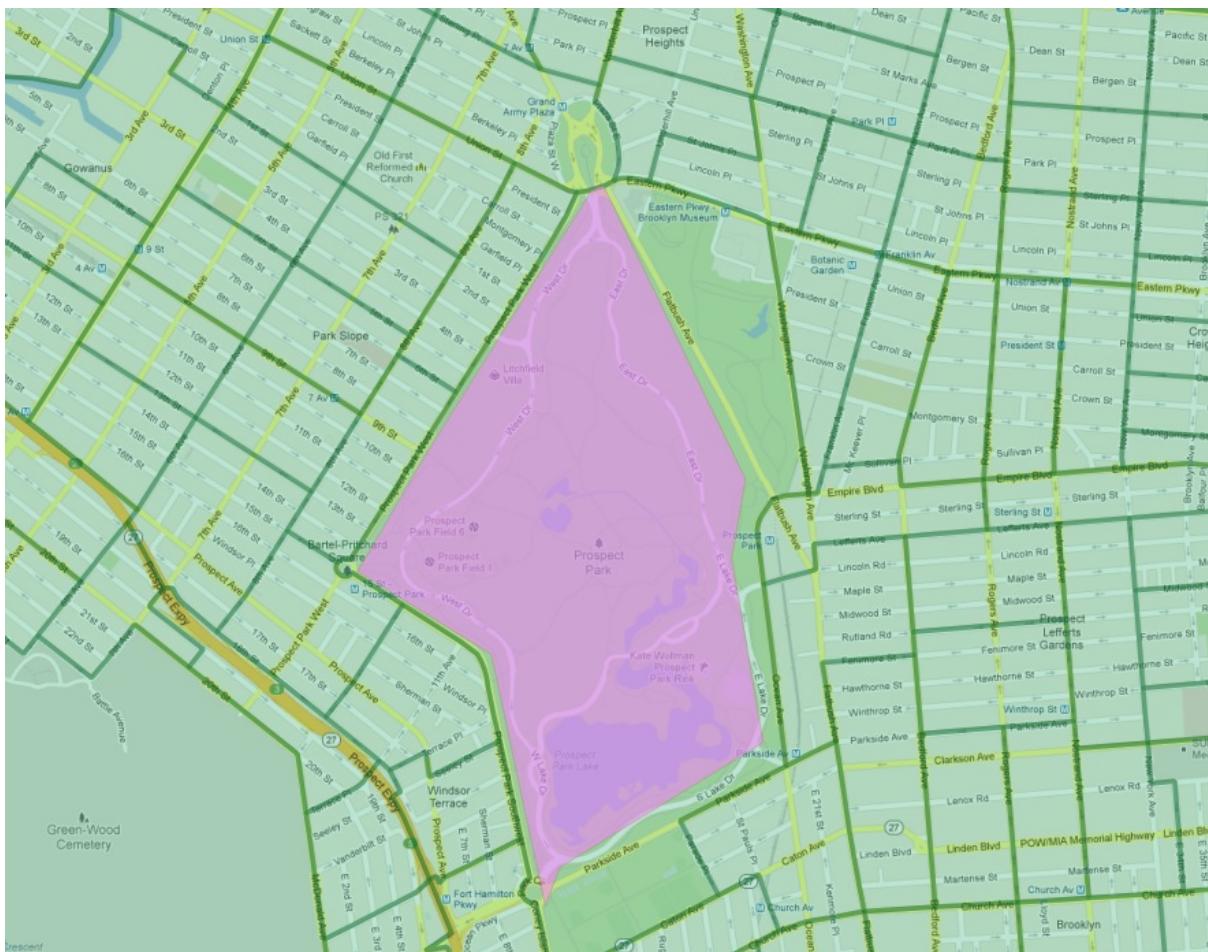
Note that the query takes longer to run now, because the `ST_Centroid` function has to be run on every census tract.

graduate_pct	name	boroname
48.0	Carnegie Hill	Manhattan
44.2	Morningside Heights	Manhattan
42.1	Greenwich Village	Manhattan
42.0	Upper West Side	Manhattan
41.4	Tribeca	Manhattan
40.7	Battery Park	Manhattan
39.5	Upper East Side	Manhattan
39.3	North Sutton Area	Manhattan
37.4	Cobble Hill	Brooklyn
37.4	Murray Hill	Manhattan

Avoiding double counting changes the results!

### What about Flatbush?

In particular, the Flatbush neighborhood has dropped off the list. The reason why can be seen by looking more closely at the map of the Flatbush neighborhood in our table.



As defined by our data source, Flatbush is not really a neighborhood in the conventional sense, since it just covers the area of Prospect Park. The census tract for that area records, naturally, zero residents. However, the neighborhood boundary does scrape one of the expensive census tracts bordering the north side of the park (in the gentrified Park Slope neighborhood). When using polygon/polygon tests, this single tract was added to the otherwise empty Flatbush, resulting in the very high score for that query.

### 5.1.3 Large Radius Distance Joins

A query that is fun to ask is “How do the commute times of people near (within 500 meters) subway stations differ from those of people far away from subway stations?”

However, the question runs into some problems of double counting: many people will be within 500 meters of multiple subway stations. Compare the population of New York:

```
SELECT Sum(popn_total)
FROM nyc_census_blocks;
```

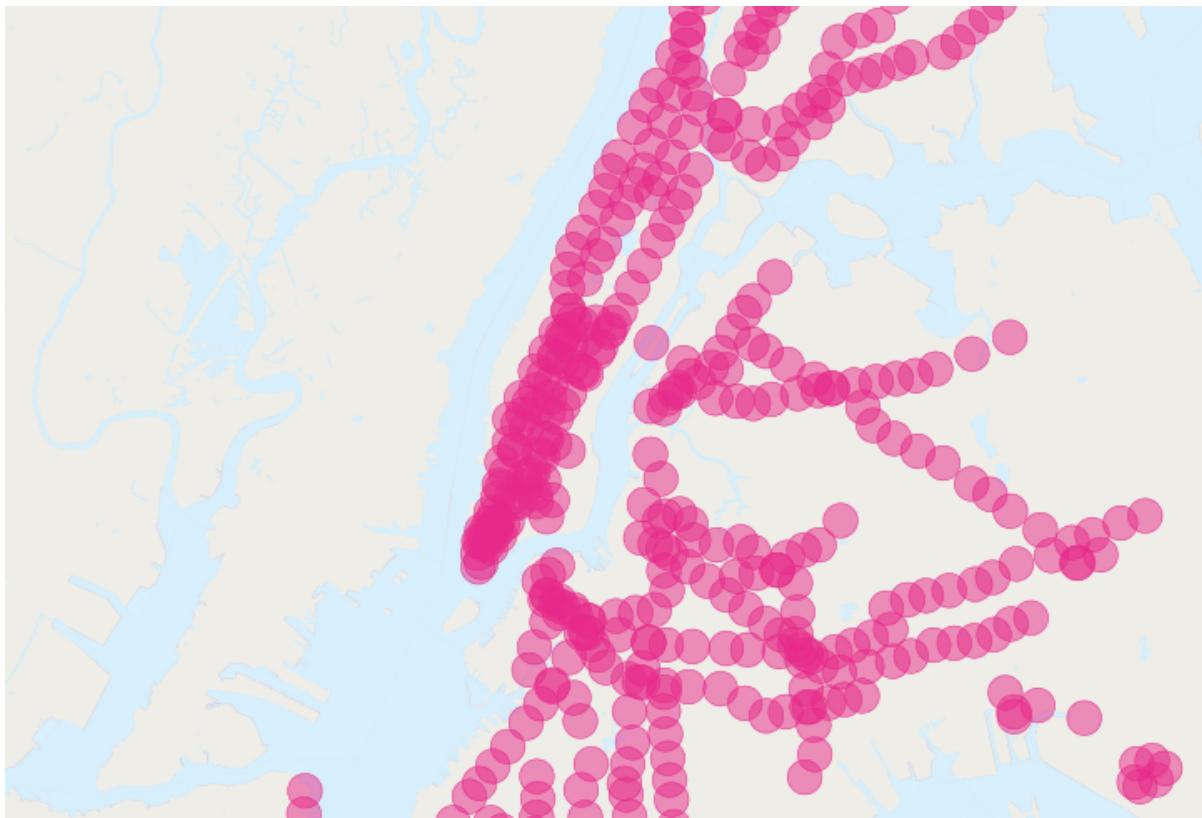
```
8175032
```

With the population of the people in New York within 500 meters of a subway station:

```
SELECT Sum(popn_total)
FROM nyc_census_blocks census
JOIN nyc_subway_stations subway
ON ST_DWithin(census.geom, subway.geom, 500);
```

10855873

There's more people close to the subway than there are people! Clearly, our simple SQL is making a big double-counting error. You can see the problem looking at the picture of the buffered subways.



The solution is to ensure that we have only distinct census blocks before passing them into the summarization portion of the query. We can do that by breaking our query up into a subquery that finds the distinct blocks, wrapped in a summarization query that returns our answer:

```
WITH distinct_blocks AS (
  SELECT DISTINCT ON (blkid) popn_total
  FROM nyc_census_blocks census
  JOIN nyc_subway_stations subway
  ON ST_DWithin(census.geom, subway.geom, 500)
)
SELECT Sum(popn_total)
FROM distinct_blocks;
```

5005743

That's better! So a bit over half the population of New York is within 500m (about a 5-7 minute walk) of the subway.

## 5.2 Validity

In 90% of the cases the answer to the question, “why is my query giving me a ‘TopologyException’ error” is “one or more of the inputs are invalid”. Which begs the question: what does it mean to be invalid, and why should we care?

### 5.2.1 What is Validity

Validity is most important for polygons, which define bounded areas and require a good deal of structure. Lines are very simple and cannot be invalid, nor can points.

Some of the rules of polygon validity feel obvious, and others feel arbitrary (and in fact, are arbitrary).

- Polygon rings must close.
- Rings that define holes should be inside rings that define exterior boundaries.
- Rings may not self-intersect (they may neither touch nor cross one another).
- Rings may not touch other rings, except at a point.

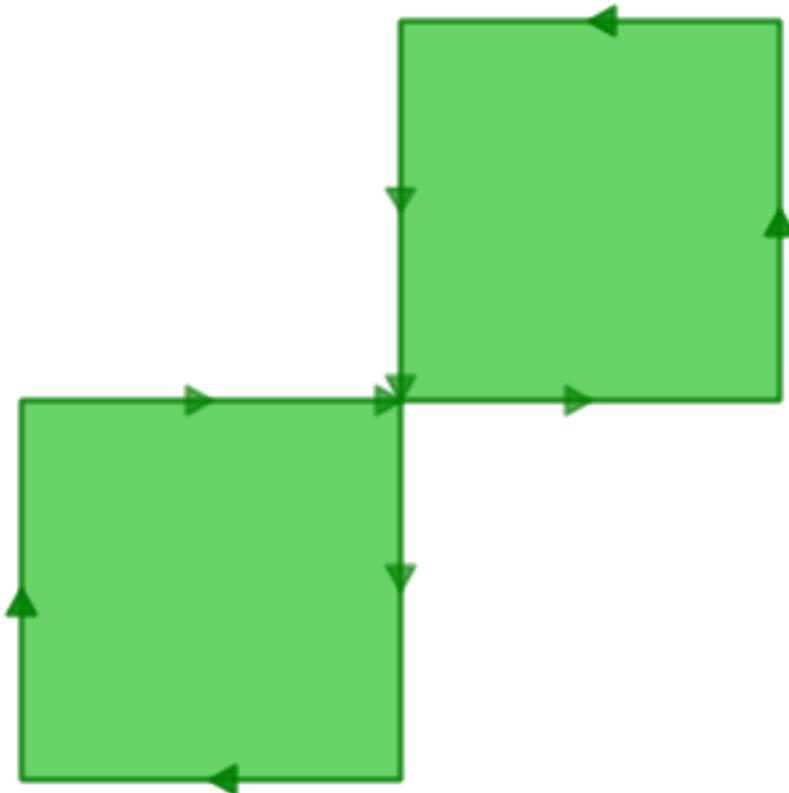
The last two rules are in the arbitrary category. There are other ways to define polygons that are equally self-consistent but the rules above are the ones used by the [OGC SFSQL](#) standard that PostGIS conforms to.

The reason the rules are important is because algorithms for geometry calculations depend on consistent structure in the inputs. It is possible to build algorithms that have no structural assumptions, but those routines tend to be very slow, because the first step in any structure-free routine is to *analyze the inputs and build structure into them*.

Here's an example of why structure matters. This polygon is invalid:

```
POLYGON((0 0, 0 1, 2 1, 2 2, 1 2, 1 0, 0 0));
```

You can see the invalidity a little more clearly in this diagram:



The outer ring is actually a figure-eight, with a self-intersection in the middle. Note that the graphic routines successfully render the polygon fill, so that visually it appears to be an “area”: two one-unit

squares, so a total area of two units of area.

Let's see what the database thinks the area of our polygon is:

```
SELECT ST_Area(ST_GeometryFromText(
    'POLYGON((0 0, 0 1, 1 1, 2 1, 2 2, 1 2, 1 1, 1 0, 0 0))'
));
```

```
st_area
-----
0
```

What's going on here? The algorithm that calculates area assumes that rings do not self-intersect. A well-behaved ring will always have the area that is bounded (the interior) on one side of the bounding line (it doesn't matter which side, just that it is on *one* side). However, in our (poorly behaved) figure-eight, the bounded area is to the right of the line for one lobe and to the left for the other. This causes the areas calculated for each lobe to cancel out (one comes out as 1, the other as -1) hence the "zero area" result.

### 5.2.2 Detecting Validity

In the previous example we had one polygon that we knew was invalid. How do we detect invalidity in a table with millions of geometries? With the **ST\_IsValid(geometry)** function. Used against our figure-eight, we get a quick answer:

```
SELECT ST_IsValid(ST_GeometryFromText(
    'POLYGON((0 0, 0 1, 1 1, 2 1, 2 2, 1 2, 1 1, 1 0, 0 0))'
));
```

```
f
```

Now we know that the feature is invalid, but we don't know why. We can use the **ST\_IsValidReason(geometry)** function to find out the source of the invalidity:

```
SELECT ST_IsValidReason(ST_GeometryFromText(
    'POLYGON((0 0, 0 1, 1 1, 2 1, 2 2, 1 2, 1 1, 1 0, 0 0))'
));
```

```
Self-intersection[1 1]
```

Note that in addition to the reason (self-intersection) the location of the invalidity (coordinate (1 1)) is also returned.

We can use the **ST\_IsValid(geometry)** function to test our tables too:

```
-- Find all the invalid polygons and what their problem is
SELECT name, boroname, ST_IsValidReason(geom)
FROM nyc_neighborhoods
WHERE NOT ST_IsValid(geom);
```

name	boroname	st_isvalidreason
Howard Beach	Queens	Self-intersection[597264.]
083368305	4499924.54228856]	

Corona	Queens	Self-intersection[595483.]
→058764138 4513817.95350787]		
Red Hook	Brooklyn	Self-intersection[584306.]
→820375986 4502360.51774956]		
Steinway	Queens	Self-intersection[593545.]
→572199759 4514735.20870587]		

### 5.2.3 Repairing Invalidity

First the bad news: there is no 100% guaranteed way to fix invalid geometries. The worst case scenario is identifying them with the **ST\_IsValid(geometry)** function, moving them to a side table, exporting that table, and repairing them externally.

Here's an example of SQL to move invalid geometries out of the main table into a side table suitable for dumping to an external cleaning process.

```
-- Side table of invalids
CREATE TABLE nyc_neighborhoods_invalid AS
SELECT * FROM nyc_neighborhoods
WHERE NOT ST_IsValid(geom);

-- Remove them from the main table
DELETE FROM nyc_neighborhoods
WHERE NOT ST_IsValid(geom);
```

A good tool for visually repairing invalid geometry is **OpenJump** (<http://openjump.org>) which includes a validation routine under **Tools->QA->Validate Selected Layers**.

Now the good news: a large proportion of invalidities **can be fixed inside the database** using either:

- the **ST\_MakeValid** function or,
- the **ST\_Buffer** function.

#### ST\_MakeValid

**ST\_MakeValid** attempts to repair invalidities without only minimal alterations to the input geometries. No vertices are dropped or moved, the structure of the object is simply re-arranged. This is a good thing for clean, but invalid data, and a bad thing for messy and invalid data.

```
-- Fix the invalid figure-8 polygon
SELECT ST_AsText(ST_MakeValid(
    'POLYGON((0 0, 0 1, 1 1, 2 1, 2 2, 1 2, 1 1, 1 0, 0 0))'
));
```

```
MULTIPOLYGON(
    ((0 0,0 1,1 1,1 0,0 0)),
    ((1 1,1 2,2 2,2 1,1 1))
)
```

**ST\_MakeValid** successfully converts the figure-8 into a multi-polygon that represents the same area.

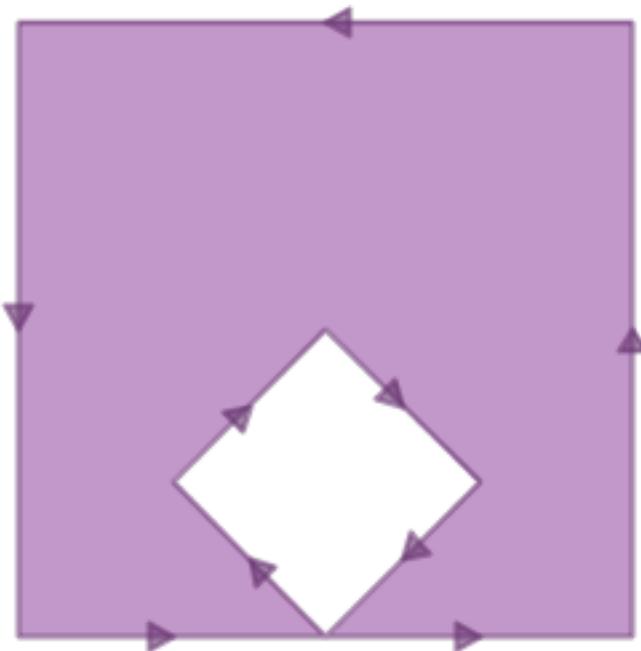
#### ST\_Buffer

Cleaning using the buffer trick takes advantage of the way buffers are built: a buffered geometry is a brand new geometry, constructed by offsetting lines from the original geometry. If you offset the

original lines by **nothing** (zero) then the new geometry will be structurally identical to the original one, but because it is built using the *OGC* topology rules, it will be valid.

For example, here's a classic invalidity – the “banana polygon” – a single ring that encloses an area but bends around to touch itself, leaving a “hole” which is not actually a hole.

```
POLYGON((0 0, 2 0, 1 1, 2 2, 3 1, 2 0, 4 0, 4 4, 0 4, 0 0))
```



Running the zero-offset buffer on the polygon returns a valid *OGC* polygon, consisting of an outer and inner ring that touch at one point.

```
SELECT ST_AsText(
    ST_Buffer(
        ST_GeometryFromText('POLYGON((0 0, 2 0, 1 1, 2 2, 3 1, 2 0, 4 0,
        ↵ 4 4, 0 4, 0 0))'),
        0.0
    )
);
```

```
POLYGON((0 0,0 4,4 4,4 0,2 0,0 0),(2 0,3 1,2 2,1 1,2 0))
```

---

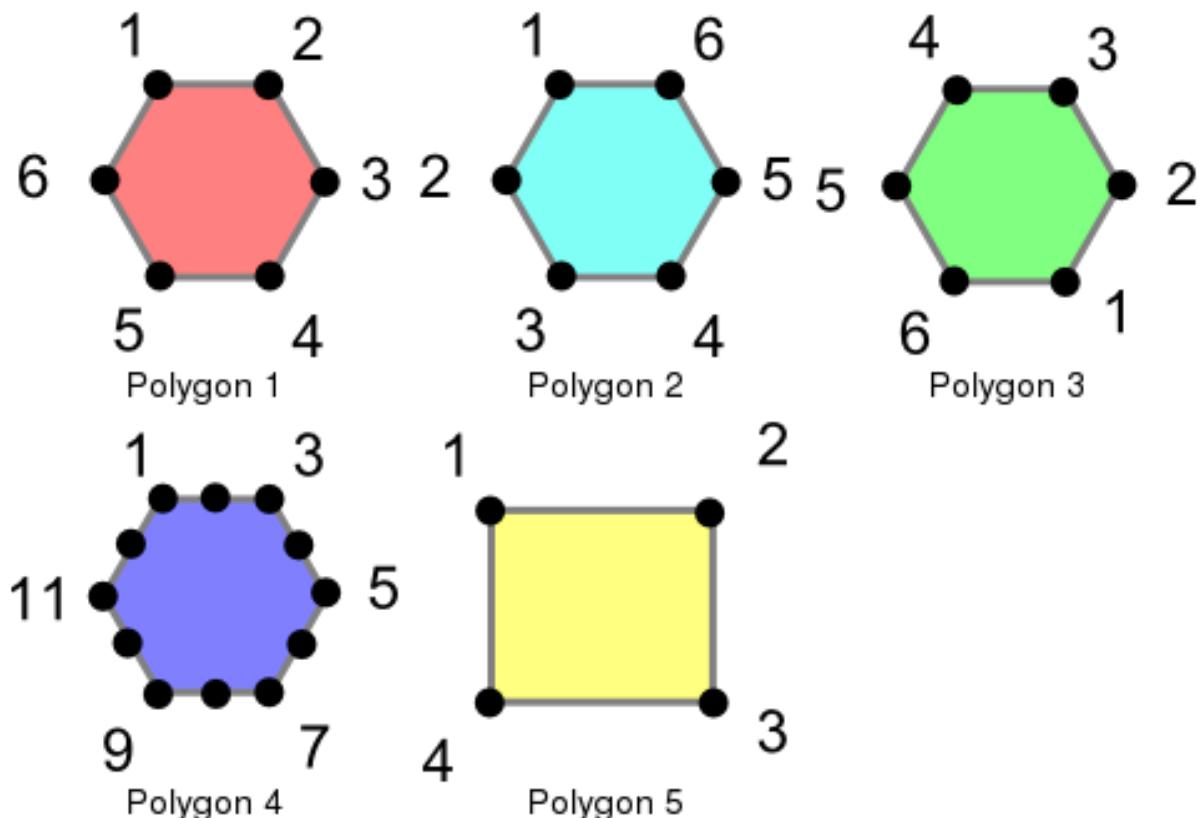
**Note:** The “banana polygon” (or “inverted shell”) is a case where the *OGC* topology model for valid geometry and the model used internally by ESRI differ. The ESRI model considers rings that touch to be invalid, and prefers the banana form for this kind of shape. The OGC model is the reverse. Neither is “correct”, they are just different ways to model the same situation.

---

## 5.3 Equality

### 5.3.1 Equality

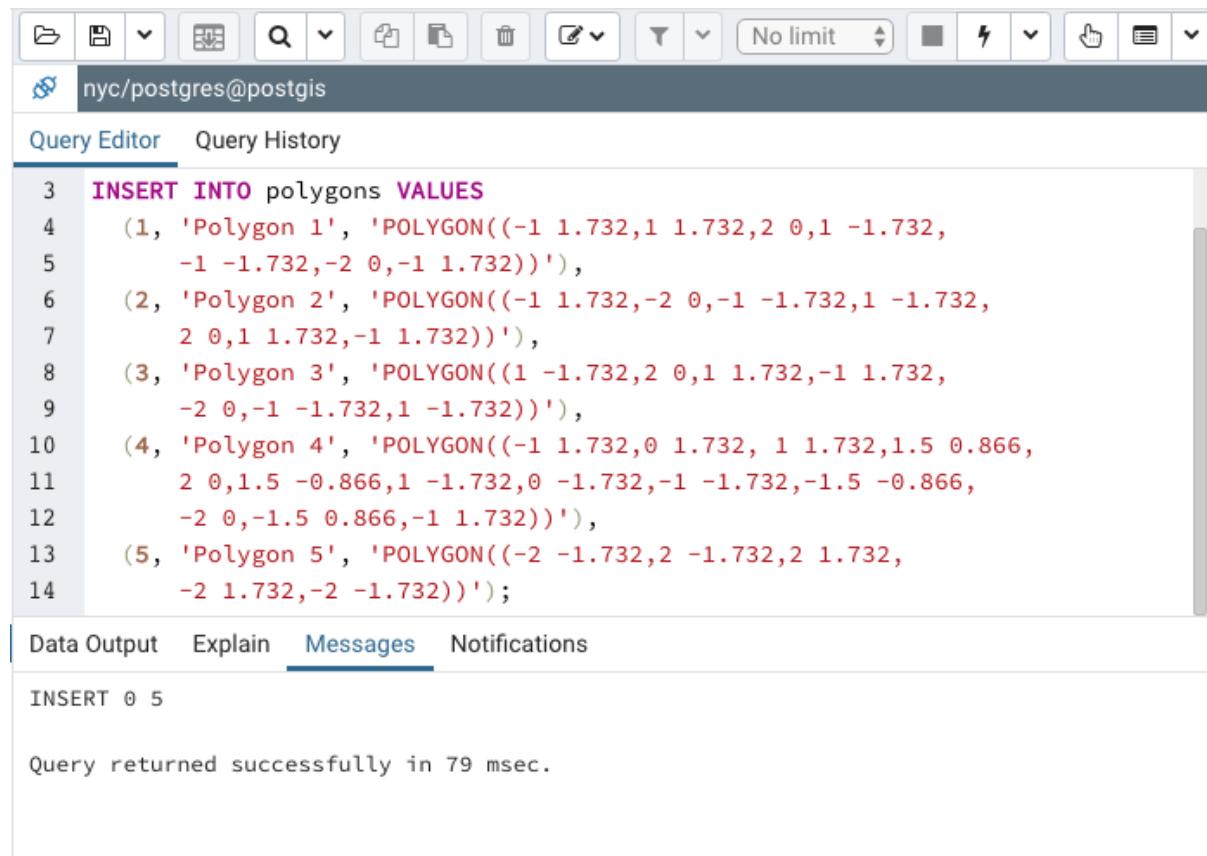
Determining equality when dealing with geometries can be tricky. PostGIS supports three different functions that can be used to determine different levels of equality, though for clarity we will use the definitions below. To illustrate these functions, we will use the following polygons.



These polygons are loaded using the following commands.

```
CREATE TABLE polygons (id integer, name varchar, poly geometry);

INSERT INTO polygons VALUES
(1, 'Polygon 1', 'POLYGON((-1 1.732,1 1.732,2 0,1 -1.732,
-1 -1.732,-2 0,-1 1.732))'),
(2, 'Polygon 2', 'POLYGON((-1 1.732,-2 0,-1 -1.732,1 -1.732,
2 0,1 1.732,-1 1.732))'),
(3, 'Polygon 3', 'POLYGON((1 -1.732,2 0,1 1.732,-1 1.732,
-2 0,-1 -1.732,1 -1.732))'),
(4, 'Polygon 4', 'POLYGON((-1 1.732,0 1.732, 1 1.732,1.5 0.866,
2 0,1.5 -0.866,1 -1.732,0 -1.732,-1 -1.732,-1.5 -0.866,
-2 0,-1.5 0.866,-1 1.732))'),
(5, 'Polygon 5', 'POLYGON((-2 -1.732,2 -1.732,2 1.732,
-2 1.732,-2 -1.732))');
```



The screenshot shows a PostgreSQL query editor interface. The top bar has various icons for file operations, search, and database management. The connection bar shows 'nyc/postgres@postgis'. Below the toolbar, tabs for 'Query Editor' and 'Query History' are visible. The main area contains a code snippet:

```

3  INSERT INTO polygons VALUES
4    (1, 'Polygon 1', 'POLYGON((-1 1.732,1 1.732,2 0,1 -1.732,
5      -1 -1.732,-2 0,-1 1.732))'),
6    (2, 'Polygon 2', 'POLYGON((-1 1.732,-2 0,-1 -1.732,1 -1.732,
7      2 0,1 1.732,-1 1.732))'),
8    (3, 'Polygon 3', 'POLYGON((1 -1.732,2 0,1 1.732,-1 1.732,
9      -2 0,-1 -1.732,1 -1.732))'),
10   (4, 'Polygon 4', 'POLYGON((-1 1.732,0 1.732, 1 1.732,1.5 0.866,
11     2 0,1.5 -0.866,1 -1.732,0 -1.732,-1 -1.732,-1.5 -0.866,
12     -2 0,-1.5 0.866,-1 1.732))'),
13   (5, 'Polygon 5', 'POLYGON((-2 -1.732,2 -1.732,2 1.732,
14     -2 1.732,-2 -1.732))');

```

Below the code, there are tabs for 'Data Output', 'Explain', 'Messages', and 'Notifications'. The 'Messages' tab is selected, showing the output of the query:

```

INSERT 0 5

```

Query returned successfully in 79 msec.

Exact equality is determined by comparing two geometries, vertex by vertex, in order, to ensure they are identical in position. The following examples show how this method can be limited in its effectiveness.

```

SELECT a.name, b.name, CASE WHEN ST_OrderingEquals(a.poly, b.poly)
    THEN 'Exactly Equal' ELSE 'Not Exactly Equal' end
    FROM polygons AS a, polygons AS b;

```

name	name	case
Polygon 1	Polygon 1	Exactly Equal
Polygon 1	Polygon 2	Not Exactly Equal
Polygon 1	Polygon 3	Not Exactly Equal
Polygon 1	Polygon 4	Not Exactly Equal
Polygon 1	Polygon 5	Not Exactly Equal
Polygon 2	Polygon 1	Not Exactly Equal
Polygon 2	Polygon 2	Exactly Equal
Polygon 2	Polygon 3	Not Exactly Equal
Polygon 2	Polygon 4	Not Exactly Equal
Polygon 2	Polygon 5	Not Exactly Equal
Polygon 3	Polygon 1	Not Exactly Equal
Polygon 3	Polygon 2	Not Exactly Equal
Polygon 3	Polygon 3	Exactly Equal
Polygon 3	Polygon 4	Not Exactly Equal
Polygon 3	Polygon 5	Not Exactly Equal
Polygon 4	Polygon 1	Not Exactly Equal
Polygon 4	Polygon 2	Not Exactly Equal
Polygon 4	Polygon 3	Not Exactly Equal
Polygon 4	Polygon 4	Exactly Equal
Polygon 4	Polygon 5	Not Exactly Equal
Polygon 5	Polygon 1	Not Exactly Equal

Polygon 5	Polygon 2	Not Exactly Equal
Polygon 5	Polygon 3	Not Exactly Equal
Polygon 5	Polygon 4	Not Exactly Equal
Polygon 5	Polygon 5	Exactly Equal

In this example, the polygons are only equal to themselves, not to other seemingly equivalent polygons (as in the case of Polygons 1 through 3). In the case of Polygons 1, 2, and 3, the vertices are in identical positions but are defined in differing orders. Polygon 4 has colinear (and thus redundant) vertices on the hexagon edges causing inequality with Polygon 1.

As we saw above, exact equality does not take into account the spatial nature of the geometries. There is a function, aptly named **ST\_Equals**, available to test the spatial equality or equivalence of geometries.

```
SELECT a.name, b.name, CASE WHEN ST_Equals(a.poly, b.poly)
    THEN 'Spatially Equal' ELSE 'Not Equal' end
FROM polygons AS a, polygons AS b;
```

name	name	case
Polygon 1	Polygon 1	Spatially Equal
Polygon 1	Polygon 2	Spatially Equal
Polygon 1	Polygon 3	Spatially Equal
Polygon 1	Polygon 4	Spatially Equal
Polygon 1	Polygon 5	Not Equal
Polygon 2	Polygon 1	Spatially Equal
Polygon 2	Polygon 2	Spatially Equal
Polygon 2	Polygon 3	Spatially Equal
Polygon 2	Polygon 4	Spatially Equal
Polygon 2	Polygon 5	Not Equal
Polygon 3	Polygon 1	Spatially Equal
Polygon 3	Polygon 2	Spatially Equal
Polygon 3	Polygon 3	Spatially Equal
Polygon 3	Polygon 4	Spatially Equal
Polygon 3	Polygon 5	Not Equal
Polygon 4	Polygon 1	Spatially Equal
Polygon 4	Polygon 2	Spatially Equal
Polygon 4	Polygon 3	Spatially Equal
Polygon 4	Polygon 4	Spatially Equal
Polygon 4	Polygon 5	Not Equal
Polygon 5	Polygon 1	Not Equal
Polygon 5	Polygon 2	Not Equal
Polygon 5	Polygon 3	Not Equal
Polygon 5	Polygon 4	Not Equal
Polygon 5	Polygon 5	Spatially Equal

These results are more in line with our intuitive understanding of equality. Polygons 1 through 4 are considered equal, since they enclose the same area. Note that neither the direction of the polygon is drawn, the starting point for defining the polygon, nor the number of points used are important here. What is important is that the polygons contain the same space.

Exact equality requires, in the worst case, comparison of each and every vertex in the geometry to determine equality. This can be slow, and may not be appropriate for comparing huge numbers of geometries. To allow for speedier comparison, the equal bounds operator, `=`, is provided. This operates only on the bounding box (rectangle), ensuring that the geometries occupy the same two dimensional extent, but not necessarily the same space.

```
SELECT a.name, b.name, CASE WHEN a.poly = b.poly
    THEN 'Equal Bounds' ELSE 'Non-equal Bounds' end
FROM polygons AS a, polygons AS b;
```

name	name	case
Polygon 1	Polygon 1	Equal Bounds
Polygon 1	Polygon 2	Non-equal Bounds
Polygon 1	Polygon 3	Non-equal Bounds
Polygon 1	Polygon 4	Non-equal Bounds
Polygon 1	Polygon 5	Non-equal Bounds
Polygon 2	Polygon 1	Non-equal Bounds
Polygon 2	Polygon 2	Equal Bounds
Polygon 2	Polygon 3	Non-equal Bounds
Polygon 2	Polygon 4	Non-equal Bounds
Polygon 2	Polygon 5	Non-equal Bounds
Polygon 3	Polygon 1	Non-equal Bounds
Polygon 3	Polygon 2	Non-equal Bounds
Polygon 3	Polygon 3	Equal Bounds
Polygon 3	Polygon 4	Non-equal Bounds
Polygon 3	Polygon 5	Non-equal Bounds
Polygon 4	Polygon 1	Non-equal Bounds
Polygon 4	Polygon 2	Non-equal Bounds
Polygon 4	Polygon 3	Non-equal Bounds
Polygon 4	Polygon 4	Equal Bounds
Polygon 4	Polygon 5	Non-equal Bounds
Polygon 5	Polygon 1	Non-equal Bounds
Polygon 5	Polygon 2	Non-equal Bounds
Polygon 5	Polygon 3	Non-equal Bounds
Polygon 5	Polygon 4	Non-equal Bounds
Polygon 5	Polygon 5	Equal Bounds

As you can see, all of our spatially equal geometries also have equal bounds. Unfortunately, Polygon 5 is also returned as equal under this test, because it shares the same bounding box as the other geometries. Why is this useful, then? Although this will be covered in detail later, the short answer is that this enables the use of spatial indexing that can quickly reduce huge comparison sets into more manageable blocks when joining or filtering data.

## 5.4 Linear Referencing

Linear referencing is a means of representing features that can be described by referencing a base set of linear features. Common examples of features that are modelled using linear referencing are:

- Highway assets, which are referenced using miles along a highway network
- Road maintenance operations, which are referenced as occurring along a road network between a pair of mile measurements.
- Aquatic inventories, where fish presence is recorded as existing between a pair of mileage-upstream measurements.
- Hydrologic characterizations (“reaches”) of streams, recorded with a from- and to-mileage.

The benefit of linear referencing models is that the dependent spatial observations do not need to be separately recorded from the base observations, and updates to the base observation layer can be carried out knowing that the dependent observations will automatically track the new geometry.

---

**Note:** The ESRI convention for linear referencing is to have a base table of linear spatial features, and a non-spatial table of “events” which includes a foreign key reference to the spatial feature and a measure along the referenced feature. We will use the term “event table” to refer to the non-spatial tables we build.

---

### 5.4.1 Creating Linear References

If you have an existing point table that you want to reference to a linear network, use the **ST\_LineLocatePoint** function, which takes a line and point, and returns the proportion along the line that the point can be found.

```
-- Simple example of locating a point half-way along a line
SELECT ST_LineLocatePoint('LINESTRING(0 0, 2 2)', 'POINT(1 1)');
-- Answer 0.5

-- What if the point is not on the line? It projects to closest point
SELECT ST_LineLocatePoint('LINESTRING(0 0, 2 2)', 'POINT(0 2)');
-- Answer 0.5
```

We can convert the **nyc\_subway\_stations** into an “event table” relative to the streets by using **ST\_LineLocatePoint**.

```
-- All the SQL below is in aid of creating the new event table
CREATE TABLE nyc_subway_station_events AS
-- We first need to get a candidate set of maybe-closest
-- streets, ordered by id and distance...
WITH ordered_nearest AS (
SELECT
    ST_GeometryN(streets.geom, 1) AS streets_geom,
    streets.id AS streets_id,
    subways.geom AS subways_geom,
    subways.id AS subways_id,
    ST_Distance(streets.geom, subways.geom) AS distance
FROM nyc_streets streets
JOIN nyc_subway_stations subways
ON ST_DWithin(streets.geom, subways.geom, 200)
ORDER BY subways_id, distance ASC
)
-- We use the 'distinct on' PostgreSQL feature to get the first
-- street (the nearest) for each unique street gid. We can then
-- pass that one street into ST_LineLocatePoint along with
-- its candidate subway station to calculate the measure.
SELECT
    DISTINCT ON (subways_id)
    subways_id,
    streets_id,
    ST_LineLocatePoint(streets_geom, subways_geom) AS measure,
    distance
FROM ordered_nearest;

-- Primary keys are useful for visualization softwares
ALTER TABLE nyc_subway_station_events ADD PRIMARY KEY (subways_id);
```

Once we have an event table, it’s fun to turn it back into a spatial view, so we can visualize the events relative to the original points they were derived from.

To go from a measure to a point, we use the **ST\_LineInterpolatePoint** function. Here's our previous simple examples reversed:

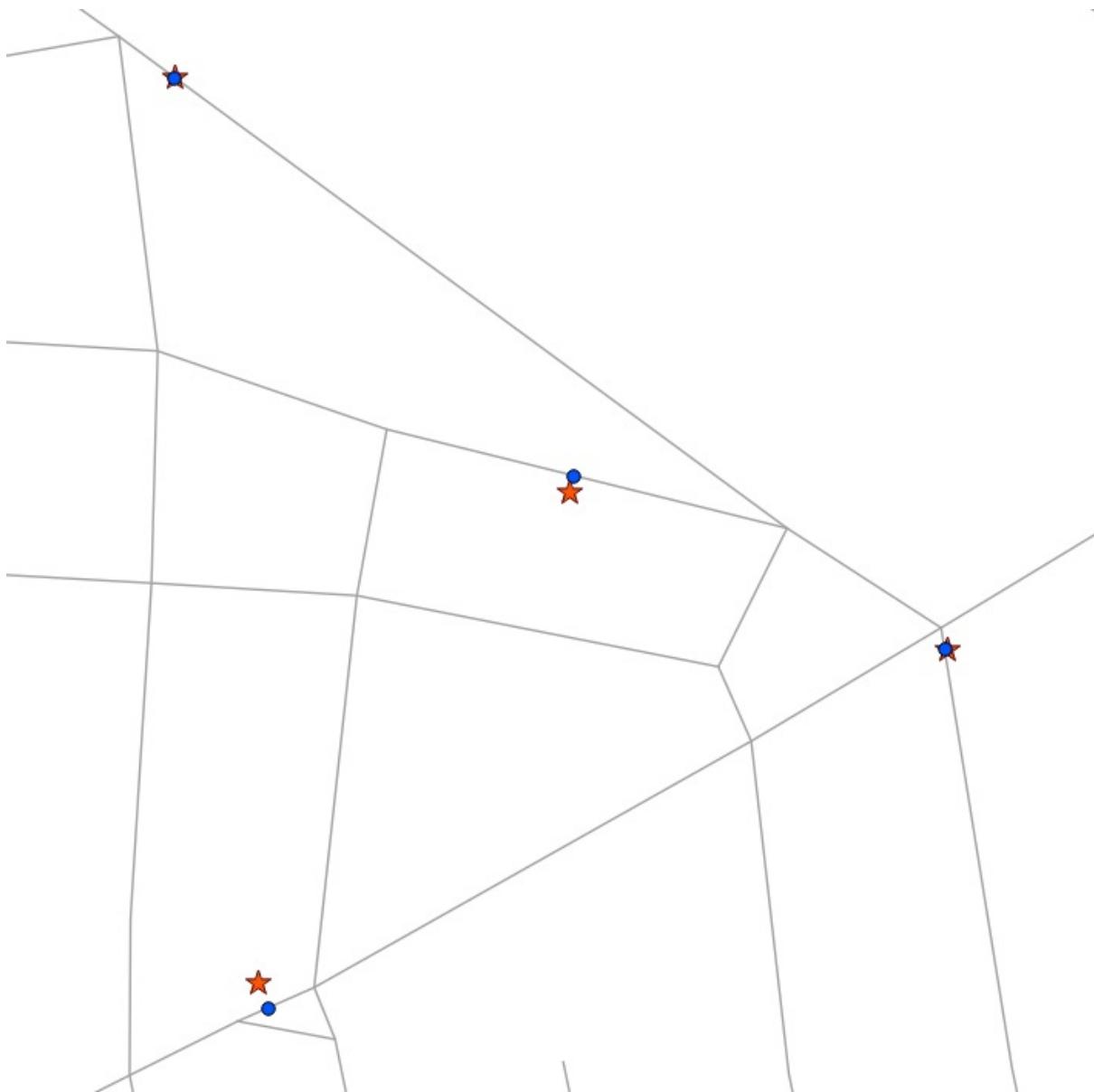
```
-- Simple example of locating a point half-way along a line
SELECT ST_AsText(ST_LineInterpolatePoint('LINESTRING(0 0, 2 2)', 0.5));

-- Answer POINT(1 1)
```

And we can join the **nyc\_subway\_station\_events** tables back to the **nyc\_streets** table and use the **measure** attribute to generate the spatial event points, without referencing the original **nyc\_subway\_stations** table.

```
-- New view that turns events back into spatial objects
CREATE OR REPLACE VIEW nyc_subway_stations_lrs AS
SELECT
    events.subways_id,
    ST_LineInterpolatePoint(ST_GeometryN(streets.geom, 1), events.measure) AS ↗geom,
    events.streets_id
FROM nyc_subway_station_events events
JOIN nyc_streets streets
ON (streets.id = events.streets_id);
```

Viewing the original (red star) and event (blue circle) points with the streets, you can see how the events are snapped directly to the closest street lines.



**Note:** One surprising use of the linear referencing functions has nothing to do with linear referencing models. As shown above, it's possible to use the functions to snap points to linear features. For use cases like GPS tracks or other inputs that are expected to reference a linear network, snapping is a handy feature to have available.

#### 5.4.2 Function List

- `ST_LineInterpolatePoint(geometry A, double measure)`: Returns a point interpolated along a line.
- `ST_LineLocatePoint(geometry A, geometry B)`: Returns a float between 0 and 1 representing the location of the closest point on LineString to the given Point.
- `ST_Line_Substring(geometry A, double from, double to)`: Return a linestring being a substring of the input one starting and ending at the given fractions of total 2d length.
- `ST_LocateAlong(geometry A, double measure, double offset)`: Return a derived geometry collection value with elements that match the specified measure.

- `ST_Locate_Between(geometry A, double from, double to, double offset)`: Return a derived geometry collection value with elements that match the specified range of measures inclusively.
- `ST_AddMeasure(geometry A, double from, double to)`: Return a derived geometry with measure elements linearly interpolated between the start and end points. If the geometry has no measure dimension, one is added.

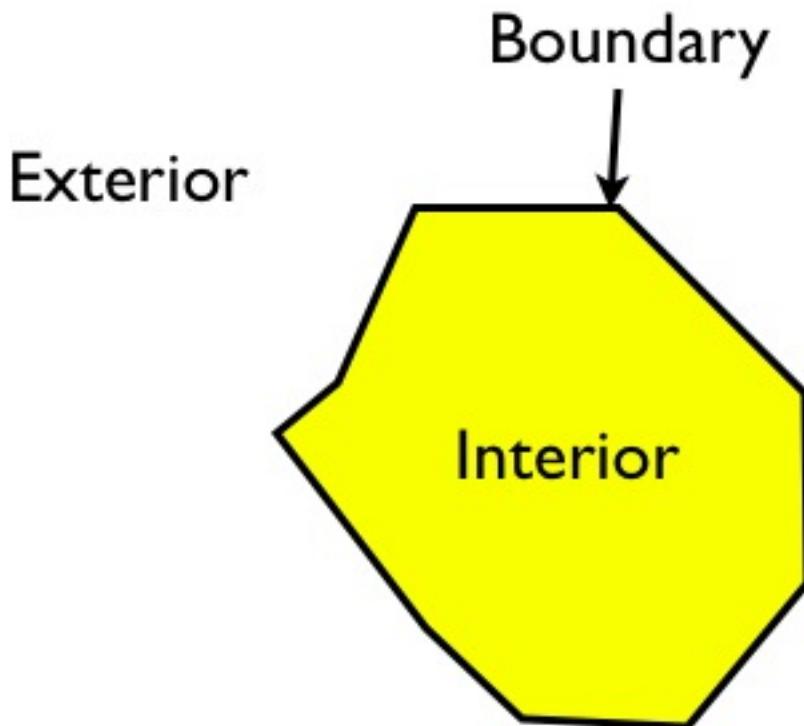
## 5.5 Dimensionally Extended 9-Intersection Model

The “Dimensionally Extended 9-Intersection Model” (DE9IM) is a framework for modelling how two spatial objects interact.

First, every spatial object has:

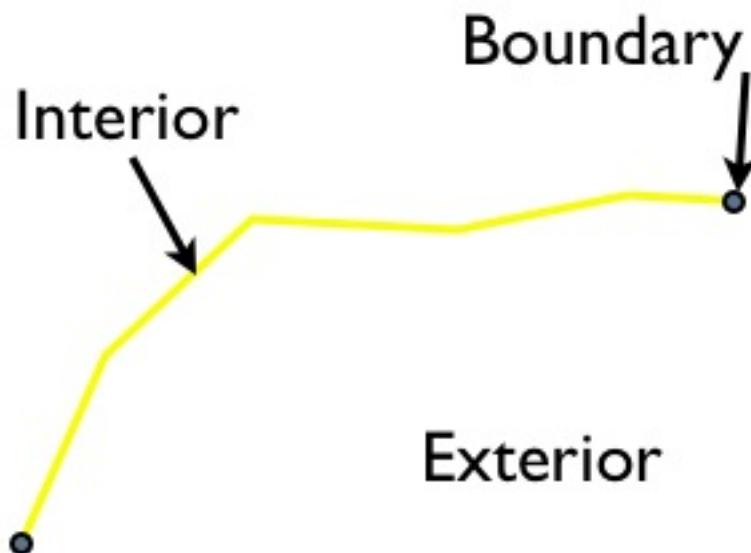
- An interior
- A boundary
- An exterior

For polygons, the interior, boundary and exterior are obvious:



The interior is the part bounded by the rings; the boundary is the rings themselves; the exterior is everything else in the plane.

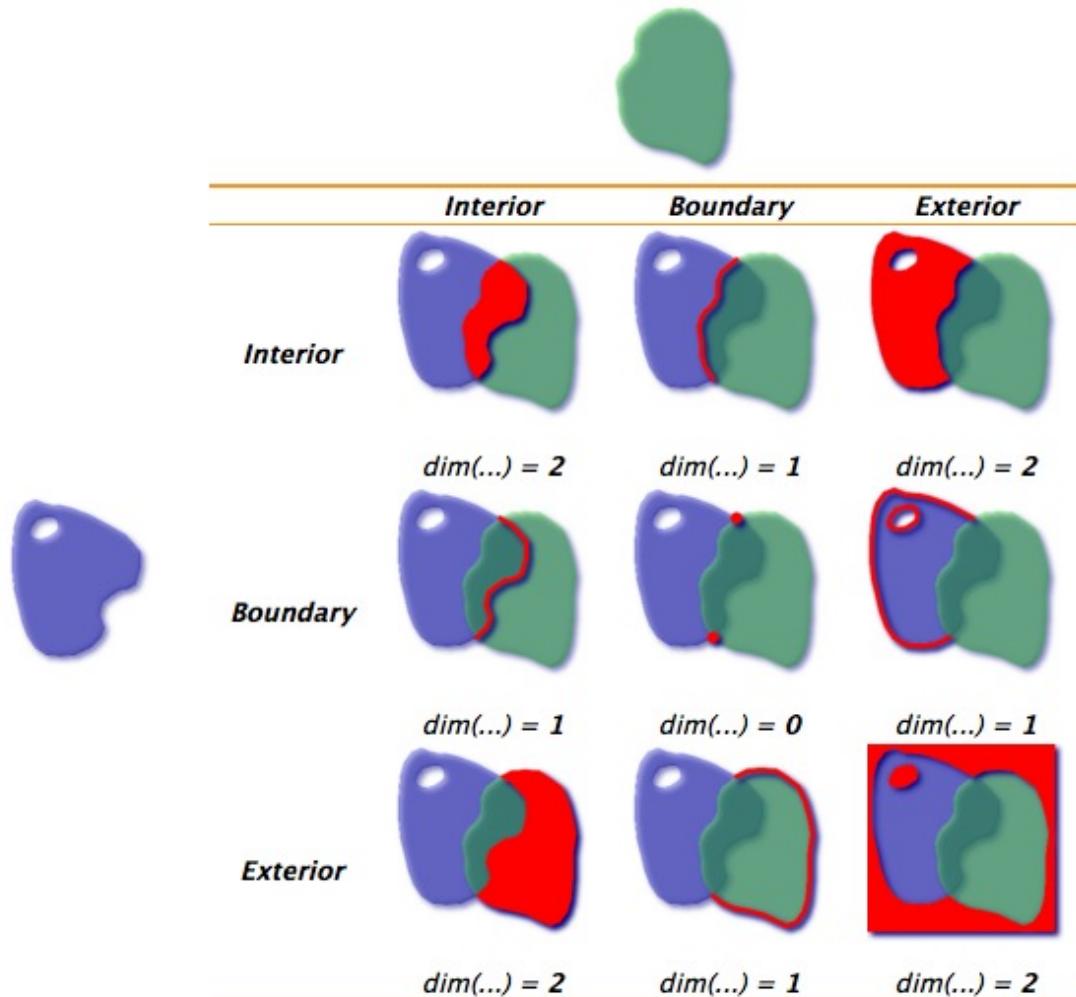
For linear features, the interior, boundary and exterior are less well-known:



The interior is the part of the line bounded by the ends; the boundary is the ends of the linear feature, and the exterior is everything else in the plane.

For points, things are even stranger: the interior is the point; the boundary is the empty set and the exterior is everything else in the plane.

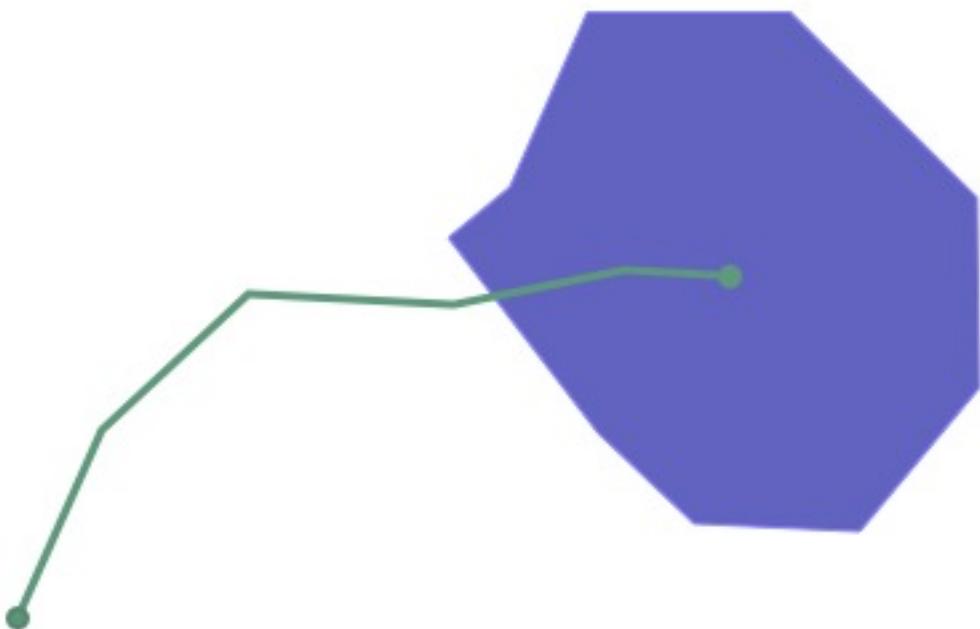
Using these definitions of interior, exterior and boundary, the relationships between any pair of spatial features can be characterized using the dimensionality of the nine possible intersections between the interiors/boundaries/exterior of a pair of objects.



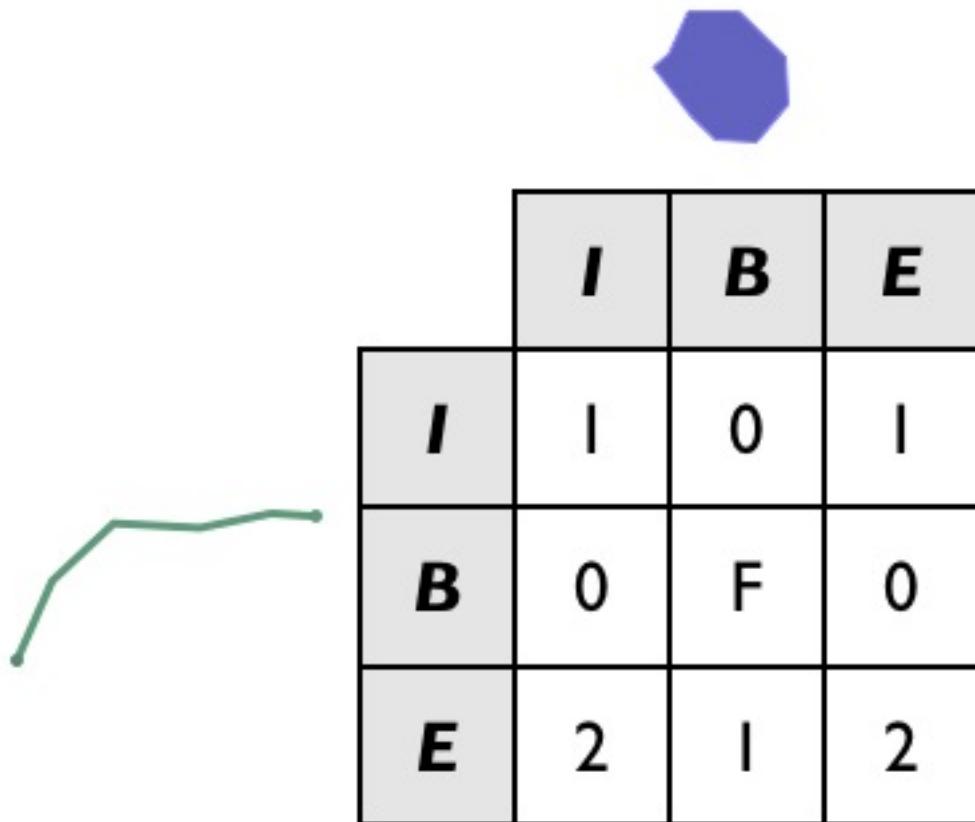
For the polygons in the example above, the intersection of the interiors is a 2-dimensional area, so that portion of the matrix is filled out with a “2”. The boundaries only intersect at points, which are zero-dimensional, so that portion of the matrix is filled out with a 0.

When there is no intersection between components, the square the matrix is filled out with an “F”.

Here's another example, of a linestring partially entering a polygon:



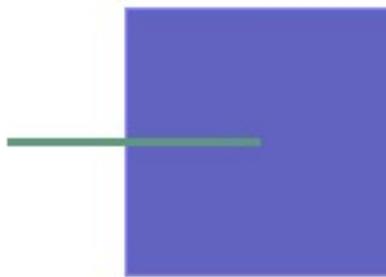
The DE9IM matrix for the interaction is this:



Note that the boundaries of the two objects don't actually intersect at all (the end point of the line interacts with the interior of the polygon, not the boundary, and vice versa), so the B/B cell is filled in with an "F".

While it's fun to visually fill out DE9IM matrices, it would be nice if a computer could do it, and that's what the **ST\_Relate** function is for.

The previous example can be simplified using a simple box and line, with the same spatial relationship as our polygon and linestring:



And we can generate the DE9IM information in SQL:

```
SELECT ST_Relate(
    'LINESTRING(0 0, 2 0)',
    'POLYGON((1 -1, 1 1, 3 1, 3 -1, 1 -1))'
);
```

The answer (1010F0212) is the same as we calculated visually, but returned as a 9-character string, with the first row, second row and third row of the table appended together.

```
101
0F0
212
```

However, the power of DE9IM matrices is not in generating them, but in using them as a matching key to find geometries with very specific relationships to one another.

```
CREATE TABLE lakes ( id serial primary key, geom geometry(POLYGON) );
CREATE TABLE docks ( id serial primary key, good boolean, geom
↪geometry(LINESTRING) );

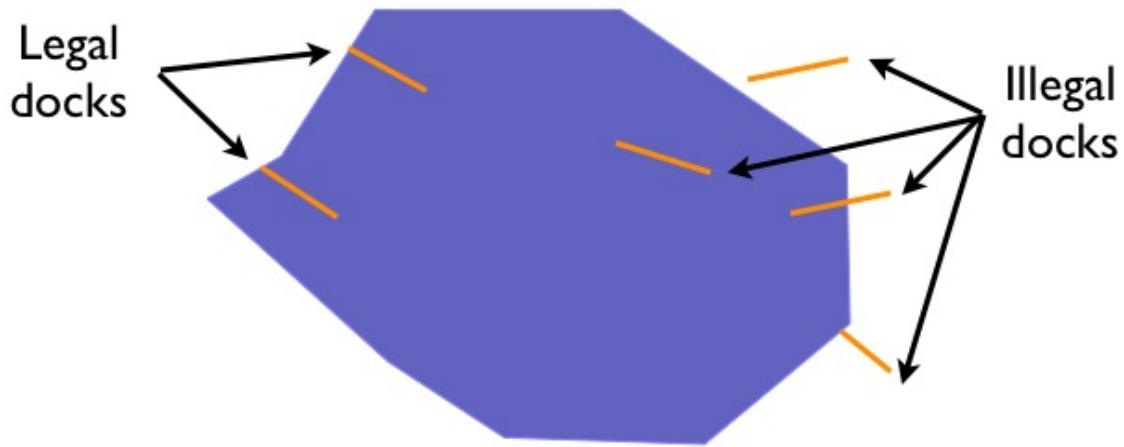
INSERT INTO lakes ( geom )
VALUES ( 'POLYGON ((100 200, 140 230, 180 310, 280 310, 390 270, 400 210,
↪ 320 140, 215 141, 150 170, 100 200))' );

INSERT INTO docks ( geom, good )
VALUES
('LINESTRING (170 290, 205 272)',true),
('LINESTRING (120 215, 176 197)',true),
('LINESTRING (290 260, 340 250)',false),
('LINESTRING (350 300, 400 320)',false),
('LINESTRING (370 230, 420 240)',false),
('LINESTRING (370 180, 390 160)',false);
```

**Note:** In versions of PostGIS >= 2.0 the behavior to add geometry columns changed and it is now

possible to specify the type of geometry and srid while creating a table as done in the previous query to create the **Lakes** and **Docks** tables. Before, the `AddGeometryColumn` function had to be used to specify the geometry column after the creation of the table, both ways are now possible.

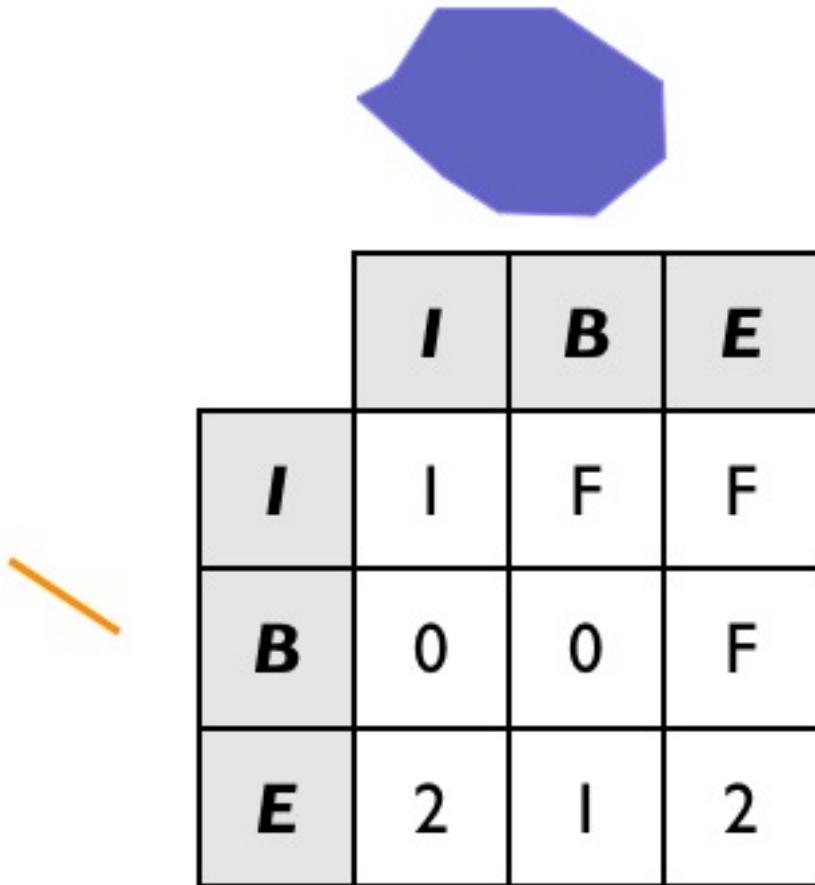
Suppose we have a data model that includes **Lakes** and **Docks**, and suppose further that Docks must be inside lakes, and must touch the boundary of their containing lake at one end. Can we find all the docks in our database that obey that rule?



Our legal docks have the following characteristics:

- Their interiors have a linear (1D) intersection with the lake interior
- Their boundaries have a point (0D) intersection with the lake interior
- Their boundaries **also** have a point (0D) intersection with the lake boundary
- Their interiors have no intersection (F) with the lake exterior

So their DE9IM matrix looks like this:



So to find all the legal docks, we would want to find all the docks that intersect lakes (a super-set of **potential** candidates we use for our join key), and then find all the docks in that set which have the legal relate pattern.

```
SELECT docks.*  
FROM docks JOIN lakes ON ST_Intersects(docks.geom, lakes.geom)  
WHERE ST_Relate(docks.geom, lakes.geom, '1FF00F212');  
  
-- Answer: our two good docks
```

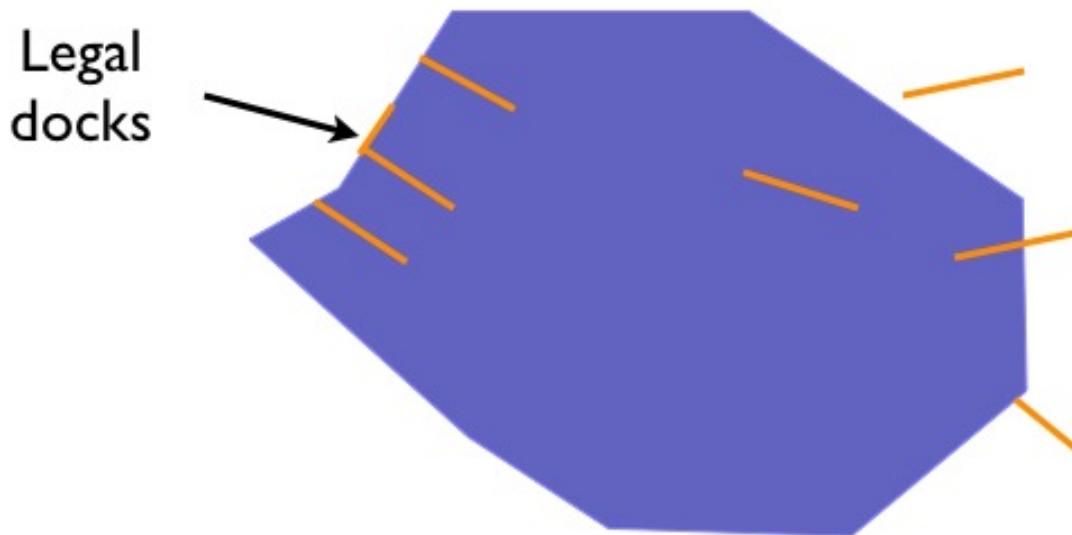
Note the use of the three-parameter version of **ST\_Relate**, which returns true if the pattern matches or false if it does not. For a fully-defined pattern like this one, the three-parameter version is not needed – we could have just used a string equality operator.

However, for looser pattern searches, the three-parameter allows substitution characters in the pattern string:

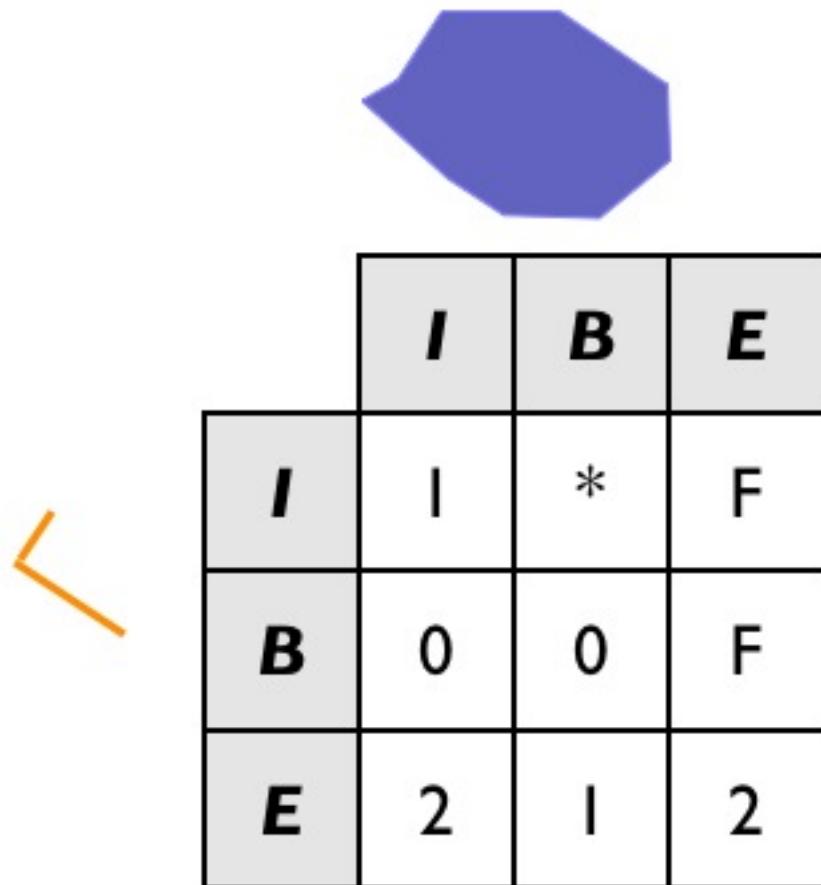
- “\*” means “any value in this cell is acceptable”
- “T” means “any non-false value (0, 1 or 2) is acceptable”

So for example, one possible dock we did not include in our example graphic is a dock with a two-dimensional intersection with the lake boundary:

```
INSERT INTO docks ( geom, good )  
VALUES ('LINESTRING (140 230, 150 250, 210 230)',true);
```



If we are to include this case in our set of “legal” docks, we need to change the relate pattern in our query. In particular, the intersection of the dock interior lake boundary can now be either 1 (our new case) or F (our original case). So we use the “\*” catchall in the pattern.

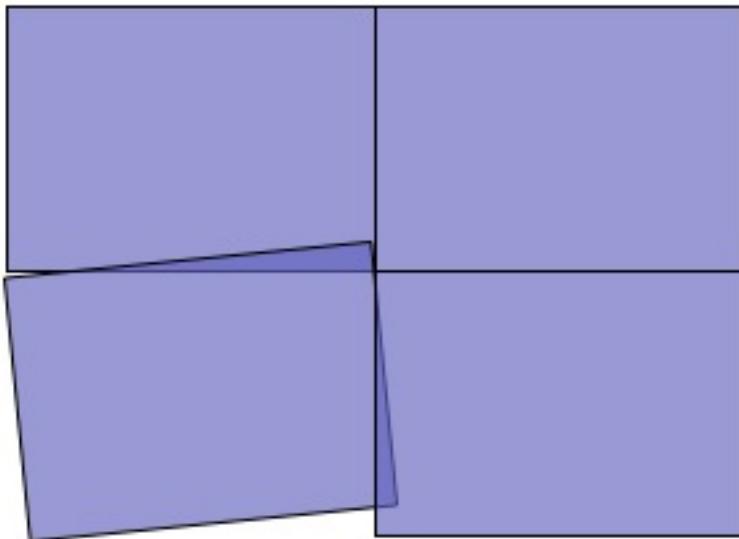


And the SQL looks like this:

```
SELECT docks.*  
FROM docks JOIN lakes ON ST_Intersects(docks.geom, lakes.geom)  
WHERE ST_Relate(docks.geom, lakes.geom, '1*F00F212');  
  
-- Answer: our (now) three good docks
```

Confirm that the stricter SQL from the previous example does *not* return the new dock.

The TIGER data is carefully quality controlled when it is prepared, so we expect our data to meet strict standards. For example: no census block should overlap any other census block. Can we test for that?

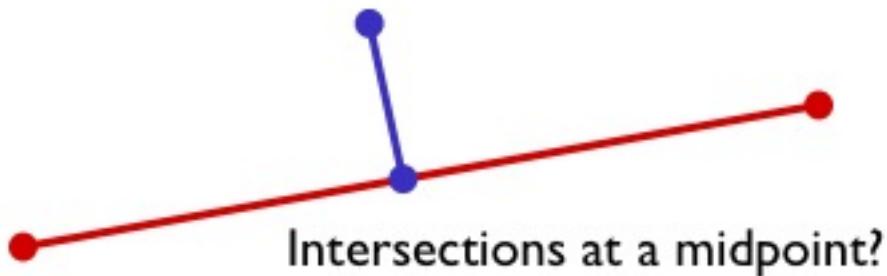


Tracts with an overlap?

Sure!

```
SELECT a.id, b.id  
FROM nyc_census_blocks a, nyc_census_blocks b  
WHERE ST_Intersects(a.geom, b.geom)  
  AND ST_Relate(a.geom, b.geom, '2*****')  
  AND a.id != b.id  
LIMIT 10;  
  
-- Answer: 10, there's some funny business
```

Similarly, we would expect that the roads data is all end-noded. That is, we expect that intersections only occur at the ends of lines, not at the mid-points.



We can test for that by looking for streets that intersect (so we have a join) but where the intersection between the boundaries is not zero-dimensional (that is, the end points don't touch):

```
SELECT a.id, b.id
FROM nyc_streets a, nyc_streets b
WHERE ST_Intersects(a.geom, b.geom)
  AND NOT ST_Relate(a.geom, b.geom, '*****0*****')
  AND a.id != b.id
LIMIT 10;
```

-- Answer: 10, this happens, so the data is not end-noded.

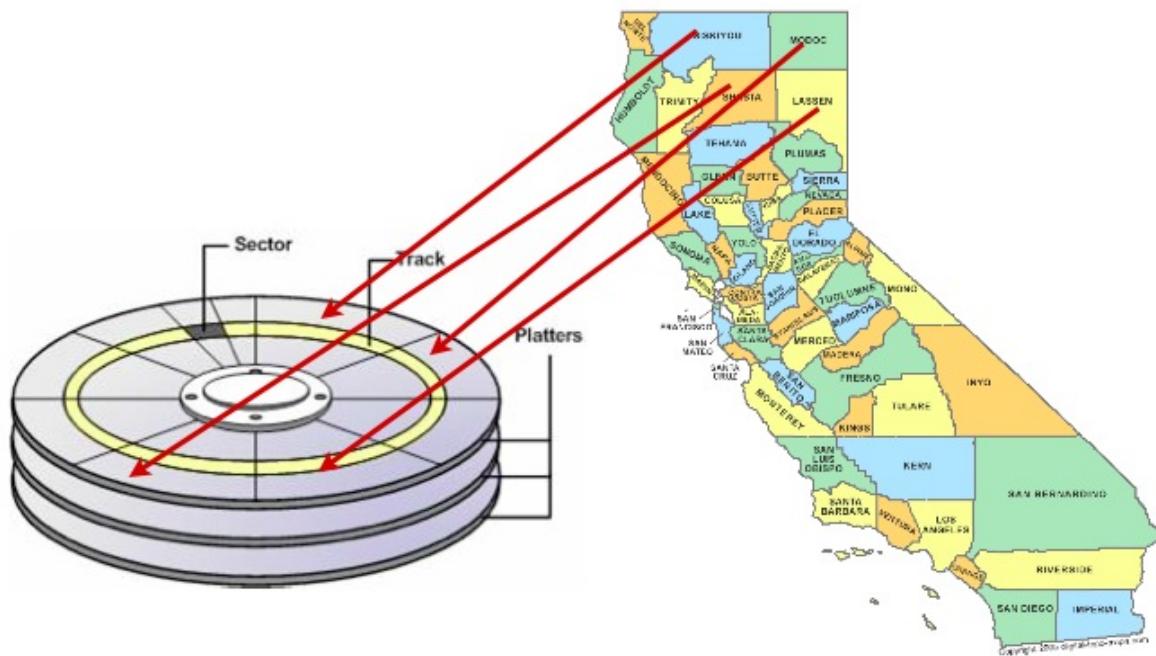
### 5.5.1 Function List

**ST\_Relate(geometry A, geometry B):** Returns a text string representing the DE9IM relationship between the geometries.

## 5.6 Clustering on Indices

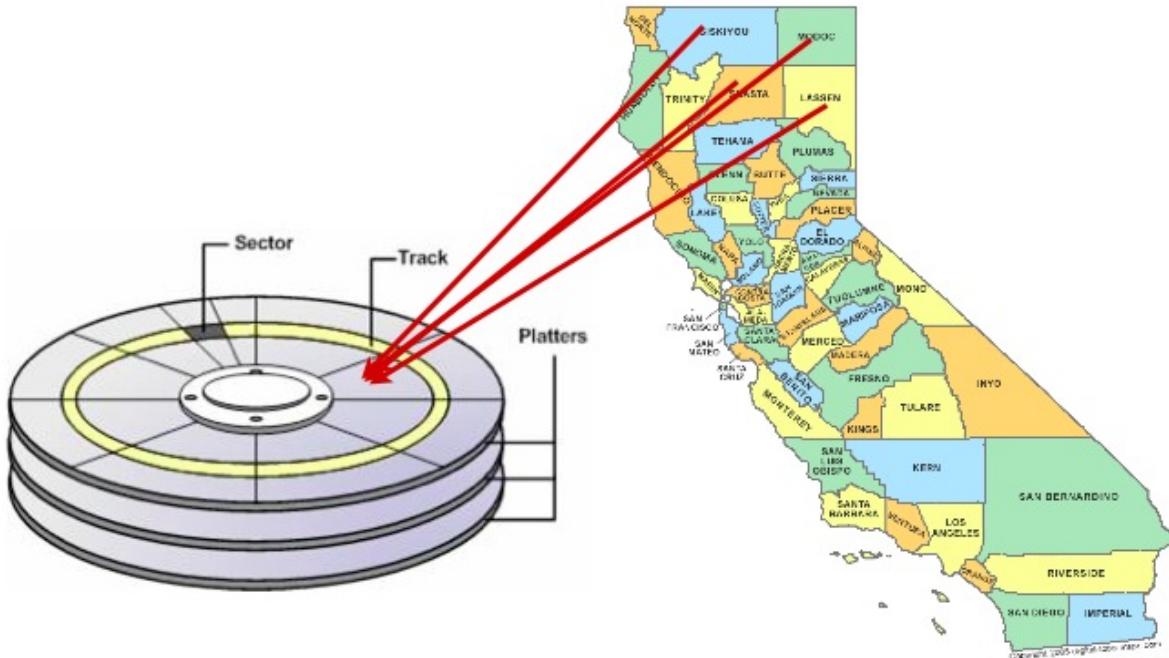
Databases can only retrieve information as fast as they can get it off of disk. Small databases will float up entirely into RAM cache, and get away from physical disk limitations, but for large databases, access to the physical disk will be a limiting stop in disk access speed.

Data is written to disk **opportunistically**, so there **isn't** necessarily any correlation between the order data is stored on the disk and the way it will be **accessed** or organized by applications.



One way to speed up access to data is to ensure that records which is likely to be retrieved together in the same result set are located in similar physical locations on the hard disk platters. This is called “clustering”.

The right clustering scheme to use can be tricky, but a general rule applies: indexes define a natural ordering scheme for data which is similar to the access pattern that will be used in retrieving the data.



Because of this, ordering the data on the disk in the same order as the index can provide a speed advantage in some cases.

### 5.6.1 Clustering on the R-Tree

Spatial data tends to be accessed in spatially correlated windows: think of the map window in a web or desktop application. All the data in the windows has similar location value (or it wouldn't be in the window!)

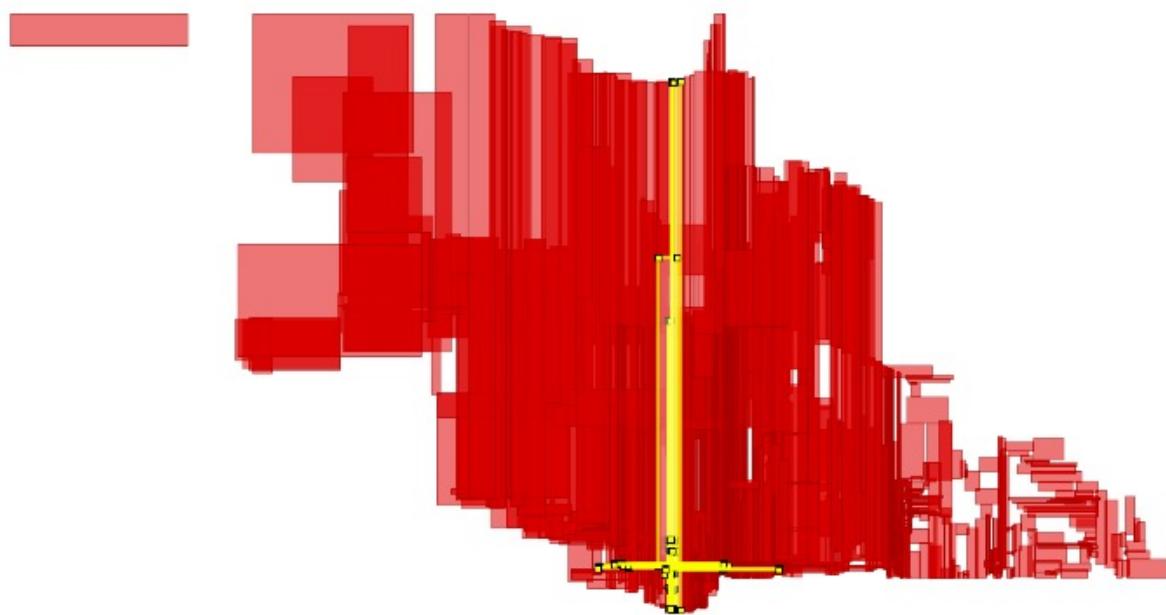
So, clustering based on a spatial index makes sense for spatial data that is going to be accessed with spatial queries: similar things tend to have similar locations.

Let's cluster our `nyc_census_blocks` based on their spatial index:

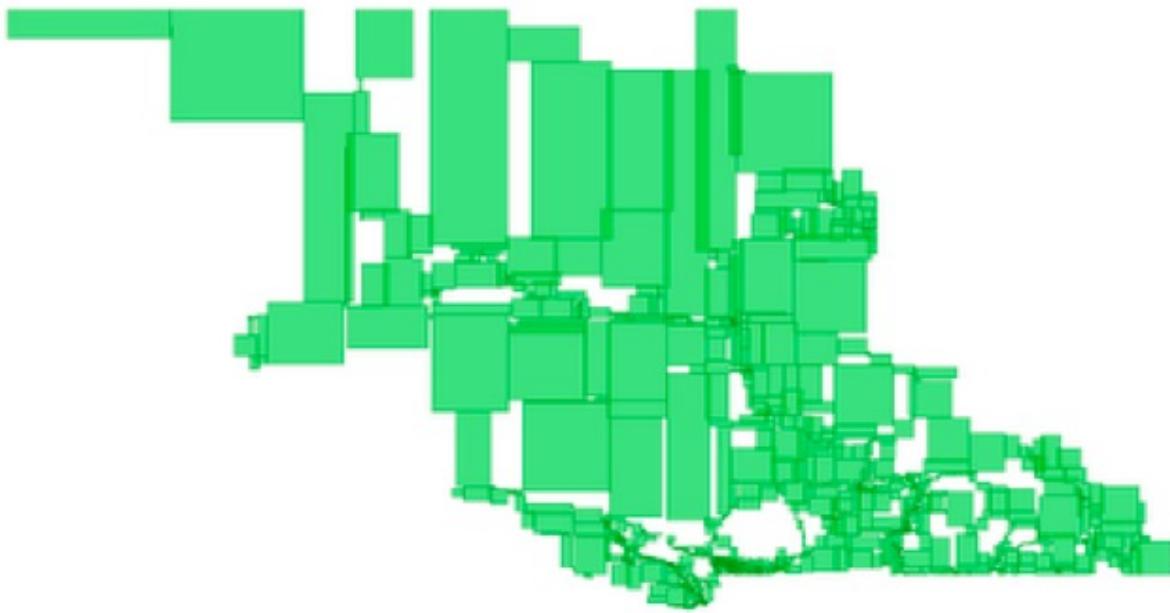
```
-- Cluster the blocks based on their spatial index
CLUSTER nyc_census_blocks USING sidx_nyc_census_blocks_geom;
```

The command re-writes the `nyc_census_blocks` in the order defined by the spatial index `sidx_nyc_census_blocks_geom`. Can you perceive a speed difference? Probably not, because the table is quite small and easily fits into memory, so disk access overhead doesn't affect performance.

One of the surprises of the R-Tree is that an R-Tree built incrementally on spatial data might not have high spatial coherence of the leaves. For example, see this visualization of the spatial index leaves of an index on roads in the province of British Columbia.



We would prefer to cluster using a more spatially compact tree, like this balanced R-Tree.



We don't have a balanced R-Tree algorithm available in PostGIS, but we do have a useful proxy that puts spatial data into a spatially autocorrelated order, the **ST\_GeoHash()** function.

### 5.6.2 Clustering on Geohash

To cluster on the ST\_GeoHash() function you first need to have a geohash index on your data. Fortunately, they are easy to build.

The geohash algorithm only works on data in geographic (longitude/latitude) coordinates, so we need to transform the geometries (to EPSG:4326, which is longitude/latitude) at the same time as we hash them.

```
CREATE INDEX geohash_nyc_census_blocks ON nyc_census_blocks (ST_GeoHash(ST_Transform(geom, 4326)));
```

Once you have a geohash index, clustering on it uses the same syntax as the R-Tree clustering.

```
CLUSTER nyc_census_blocks USING geohash_nyc_census_blocks;
```

Now your data is nicely arranged in spatially correlated order!

### 5.6.3 Function List

**ST\_GeoHash(geometry A)**: Returns a text string representing the GeoHash of the bounds of the object.

## 5.7 3-D

---

**Note:** This section refers to many features that are only available with PostGIS 2.0 and higher.

---

### 5.7.1 3-D Geometries

So far, we have been working with 2-D geometries, with only X and Y coordinates. But PostGIS supports additional dimensions on all geometry types, a "Z" dimension to add height information and

a “M” dimension for additional dimensional information (commonly time, or road-mile, or upstream-distance information) for each coordinate.

For 3-D and 4-D geometries, the extra dimensions are added as extra coordinates for each vertex in the geometry, and the geometry type is enhanced to indicate how to interpret the extra dimensions. Adding the extra dimensions results in three extra possible geometry types for each geometry primitive:

- Point (a 2-D type) is joined by PointZ, PointM and PointZM types.
- Linestring (a 2-D type) is joined by LineStringZ, LineStringM and LineStringZM types.
- Polygon (a 2-D type) is joined by PolygonZ, PolygonM and PolygonZM types.
- And so on.

For well-known text ([WKT](#)) representation, the format for higher dimensional geometries is given by the ISO SQL/MM specification. The extra dimensionality information is simply added to the text string after the type name, and the extra coordinates added after the X/Y information. For example:

- POINT ZM (1 2 3 4)
- LINESTRING M (1 1 0, 1 2 0, 1 3 1, 2 2 0)
- POLYGON Z ((0 0 0, 0 1 0, 1 1 0, 1 0 0, 0 0 0))

The ST\_AsText() function will return the above representations when dealing with 3-D and 4-D geometries.

For well-known binary ([WKB](#)) representation, the format for higher dimensional geometries is given by the ISO SQL/MM specification. The BNF form of the format is available from <http://svn.osgeo.org/postgis/trunk/doc/bnf-wkb.txt>.

In addition to higher-dimensional forms of the standard types, PostGIS includes a few new types that make sense in a 3-D space:

- The TIN type allows you to model triangular meshes as rows in your database.
- The POLYHEDRALSURFACE allows you to model volumetric objects in your database.

Since both these types are for modelling 3-D objects, it only really makes sense to use the Z variants. An example of a POLYHEDRALSURFACE Z would be the 1 unit cube:

```
POLYHEDRALSURFACE Z (
  ((0 0 0, 0 1 0, 1 1 0, 1 0 0, 0 0 0)),
  ((0 0 0, 0 1 0, 0 1 1, 0 0 1, 0 0 0)),
  ((0 0 0, 1 0 0, 1 0 1, 0 0 1, 0 0 0)),
  ((1 1 1, 1 0 1, 0 0 1, 0 1 1, 1 1 1)),
  ((1 1 1, 1 0 1, 1 0 0, 1 1 0, 1 1 1)),
  ((1 1 1, 1 1 0, 0 1 0, 0 1 1, 1 1 1))
)
```

## 5.7.2 3-D Functions

There are a number of functions built to calculate relationships between 3-D objects:

- ST\_3DClosestPoint — Returns the 3-dimensional point on g1 that is closest to g2. This is the first point of the 3D shortest line.
- ST\_3DDistance — For geometry type Returns the 3-dimensional cartesian minimum distance (based on spatial ref) between two geometries in projected units.

- ST\_3DDWithin — For 3d (z) geometry type Returns true if two geometries 3d distance is within number of units.
- ST\_3DDFullyWithin — Returns true if all of the 3D geometries are within the specified distance of one another.
- ST\_3DIIntersects — Returns TRUE if the Geometries “spatially intersect” in 3d - only for points and linestrings
- ST\_3DLongestLine — Returns the 3-dimensional longest line between two geometries
- ST\_3DMaxDistance — For geometry type Returns the 3-dimensional cartesian maximum distance (based on spatial ref) between two geometries in projected units.
- ST\_3DShortestLine — Returns the 3-dimensional shortest line between two geometries

For example, we can calculate the distance between our unit cube and a point using the ST\_3DDistance function:

```
-- This is really the distance between the top corner
-- and the point.
SELECT ST_3DDistance(
    'POLYHEDRALSURFACE Z (
        ((0 0 0, 0 1 0, 1 1 0, 1 0 0, 0 0 0)),
        ((0 0 0, 0 1 0, 0 1 1, 0 0 1, 0 0 0)),
        ((0 0 0, 1 0 0, 1 0 1, 0 0 1, 0 0 0)),
        ((1 1 1, 1 0 1, 0 0 1, 0 1 1, 1 1 1)),
        ((1 1 1, 1 0 1, 1 0 0, 1 1 0, 1 1 1)),
        ((1 1 1, 1 1 0, 0 1 0, 0 1 1, 1 1 1))
    )'::geometry,
    'POINT Z (2 2 2)'::geometry
);

-- So here's a shorter form.
SELECT ST_3DDistance(
    'POINT Z (1 1 1)'::geometry,
    'POINT Z (2 2 2)'::geometry
);

-- Both return 1.73205080756888 == sqrt(3) as expected
```

### 5.7.3 N-D Indexes

Once you have data in higher dimensions it may make sense to index it. However, you should think carefully about the distribution of your data in all dimensions before applying a multi-dimensional index.

Indexes are only useful when they allow the database to drastically reduce the number of return rows as a result of a WHERE condition. For a higher dimension index to be useful, the data must cover a wide range of that dimension, relative to the kinds of queries you are constructing.

- A set of DEM points would probably be a *poor* candidate for a 3-D index, since the queries would usually be extracting a 2-D box of points, and rarely attempting to select a Z-slice of points.
- A set of GPS traces in X/Y/T space might be a *good* candidate for a 3-D index, if the GPS tracks overlapped each other frequently in all dimensions (for example, driving the same route over and over at different times), since there would be large variability in all dimensions of the data set.

You can create a multi-dimensional index on data of any dimensionality (even mixed dimensionality). For example, to create a multi-dimensional index on the `nyc_streets` table,

```
CREATE INDEX sidx_nd_nyc_streets ON nyc_streets
USING GIST (geom gist_geometry_ops_nd);
```

The `gist_geometry_ops_nd` parameter tells PostGIS to use the N-D index instead of the standard 2-D index.

Once you have the index built, you can use it in queries with the `&&&` index operator. `&&&` has the same semantics as `&&`, “bounding boxes interact”, but applies those semantics using all the dimensions of the input geometries. Geometries with mis-matching dimensionality do not interact.

```
-- Returns true (both 3-D on the zero plane)
SELECT 'POINT Z (1 1 0)'::geometry &&&
    'POLYGON ((0 0 0, 0 2 0, 2 2 0, 2 0 0, 0 0 0))'::geometry;

-- Returns false (one 2-D one 3-D)
SELECT 'POINT Z (1 1 1)'::geometry &&&
    'POLYGON ((0 0, 0 2, 2 2, 2 0, 0 0))'::geometry;

-- Returns true (the volume around the linestring interacts with the point)
SELECT 'LINESTRING Z(0 0 0, 1 1 1)'::geometry &&&
    'POINT(0 1 1)'::geometry;
```

To search the `nyc_streets` table using the N-D index, just replace the usual `&&` 2-D index operator with the `&&&` operator.

```
-- N-D index operator
SELECT id, name
FROM nyc_streets
WHERE geom &&&
    ST_SetSRID('LINESTRING(586785 4492901, 587561 4493037)'::geometry,
    ↴26918);

-- 2-D index operator
SELECT id, name
FROM nyc_streets
WHERE geom &&
    ST_SetSRID('LINESTRING(586785 4492901, 587561 4493037)'::geometry,
    ↴26918);
```

The results should be the same. In general the N-D index is very slightly slower than the 2-D index, so only use the N-D index where you are certain that N-D queries will improve the selectivity of your queries.

## 5.8 Nearest-Neighbour Searching

---

**Note:** This section refers to a feature that is only available with PostGIS 2.0 and higher.

---

### 5.8.1 What is a Nearest Neighbour Search?

A frequently posed spatial query is: “what is the nearest <candidate feature> to <query feature>?”

Unlike a distance search, the “nearest neighbour” search doesn’t include any measurement restricting how far away candidate geometries might be, features of any distance away will be accepted, as long as

they are the *nearest*. This poses a problem for traditional index-assisted queries, that require a search box, and therefore need some kind of measurement value to build the box.

The naive way to carry out a nearest neighbour query is to order the candidate table by distance from the query geometry, and then take the record with the smallest distance:

```
-- Closest street to Broad Street station is Wall St
SELECT streets.id, streets.name
FROM
    nyc_streets streets,
    nyc_subway_stations subways
WHERE subways.name = 'Broad St'
ORDER BY ST_Distance(streets.geom, subways.geom) ASC
LIMIT 1;
```

The trouble with this approach is that it forces the database to calculate the distance between the query geometry and *every* feature in the table of candidate features, then sort them all. For a large table of candidate features, it is not a reasonable approach.

One way to improve performance is to add an index constraint to the search. This requires a magic number: what's the smallest box we could search around the query geometry, and still come up with at least one candidate geometry?

If you turn on timing, you can see the performance difference between the box-assisted query below and the simple query above.

```
-- Closest street to Broad Street station is Wall St
SELECT streets.id, streets.name
FROM
    nyc_streets streets,
    nyc_subway_stations subways
WHERE subways.name = 'Broad St'
AND streets.geom && ST_Expand(subways.geom, 200) -- Magic number: 200m
ORDER BY ST_Distance(streets.geom, subways.geom) ASC
LIMIT 1;
```

The problem with this approach is the magic number of 200 meters. What if there had not happened to be any roads within 200m? We would have failed to come up with a result: there is always a nearest neighbour, it just might not be within 200m.

### 5.8.2 Index-based KNN

“KNN” stands for “K nearest neighbours”, where “K” is the number of neighbours you are looking for.

KNN is a pure index based nearest neighbour search. By walking up and down the index, the search can find the nearest candidate geometries without using any magical search radius numbers, so the technique is suitable and high performance even for very large tables with highly variable data densities.

---

**Note:** The KNN feature is only available on PostGIS 2.0 with PostgreSQL 9.1 or greater.

---

The KNN system works by evaluating distances between bounding boxes inside the PostGIS R-Tree index.

Because the index is built using the bounding boxes of geometries, the distances between any geometries that are not points will be inexact: they will be the distances between the bounding boxes of geometries.

The syntax of the index-based KNN query places a special “index-based distance operator” in the ORDER BY clause of the query, in this case “`<->`”. There are two index-based distance operators,

- `<->` means “distance between geometries”
- `<#>` means “distance between bounding boxes”

One side of the index-based distance operator must be a literal geometry value. We can get away with a subquery that returns as single geometry, or we could include a *WKT* geometry instead.

```
-- Closest 10 streets to Broad Street station are ?
SELECT
    streets.id,
    streets.name
FROM
    nyc_streets streets
ORDER BY
    streets.geom <->
        (SELECT geom FROM nyc_subway_stations WHERE name = 'Broad St')
LIMIT 10;

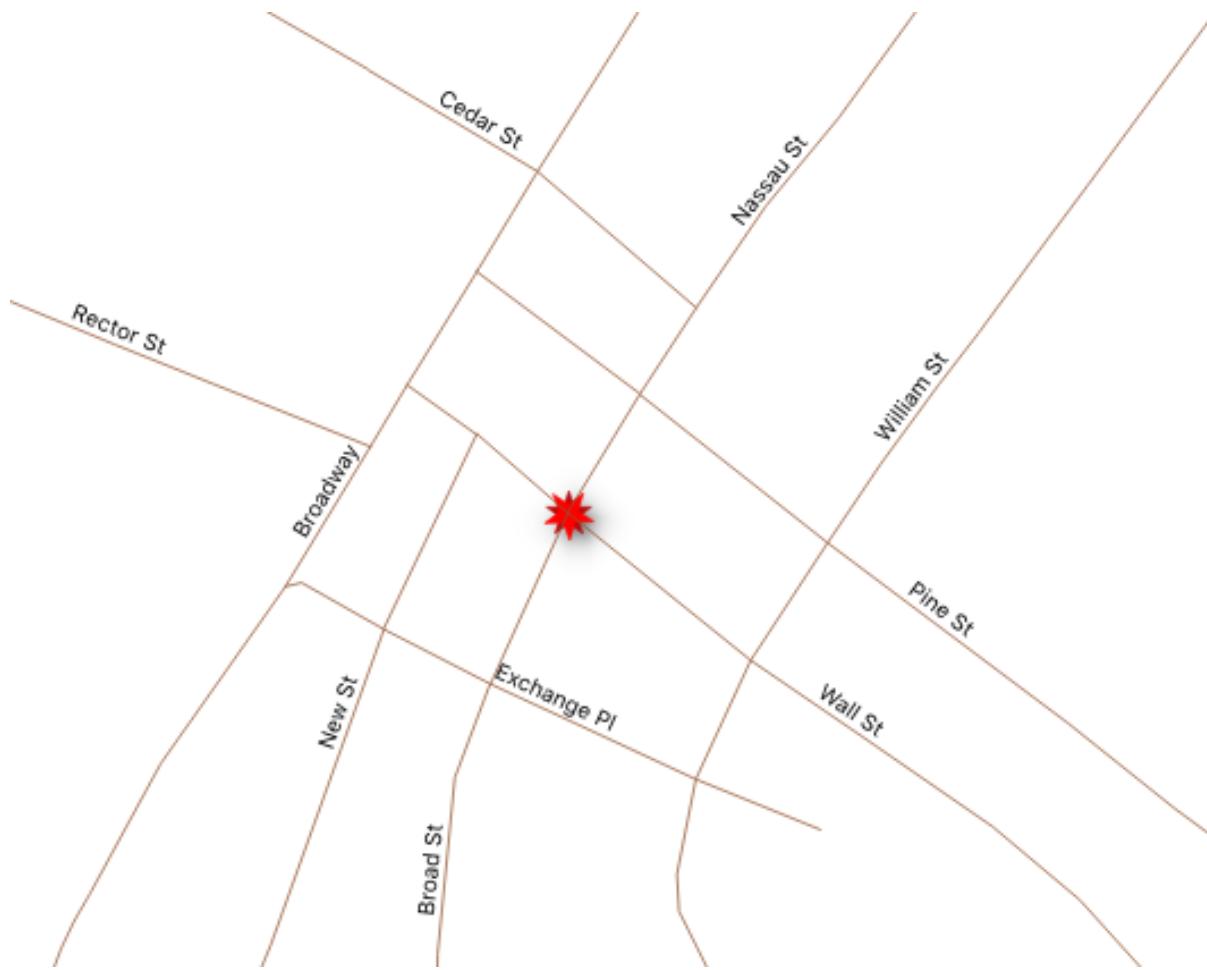
-- Same query using a geometry EWKT literal

SELECT ST_AsEWKT(geom)
FROM nyc_subway_stations
WHERE name = 'Broad St';
-- SRID=26918;POINT(583571 4506714)

SELECT
    streets.id,
    streets.name,
    ST_Distance(
        streets.geom,
        'SRID=26918;POINT(583571.905921312 4506714.34119218)'::geometry
    ) AS distance
FROM
    nyc_streets streets
ORDER BY
    streets.geom <->
        'SRID=26918;POINT(583571.905921312 4506714.34119218)'::geometry
LIMIT 10;
```

<code>id</code>	<code>name</code>	<code>distance</code>
17394	Wall St	0.714202224374917
17399	Broad St	0.872022763400183
17445	Nassau St	1.29928727926582
17357	New St	63.9499165490674
17411	Pine St	75.8461038368021
17367	Exchange Pl	101.6241843136
17322	Broadway	112.049824188021
17296	Rector St	114.442000781044
17478	William St	126.934064759446
17354	Cedar St	133.009278387597

Remember that all the calculations are being done using geometries. Here’s what the map looks like for the results of the query:



We can see that the station falls right on the Wall Street line so the `<->` operator computes the distance between geometries giving the proper answer. Moreover, the `ST_Distance` is in alignment with the order provided by the `<->` operator, proving the knn functionality works!

---

**Note:** For versions of PostgreSQL below 9.5 the `<->` operator would compute distances using the centroid of the bounding boxes of the geometries, producing sometimes approximations that don't match a nearest neighbour search between geometries.

---

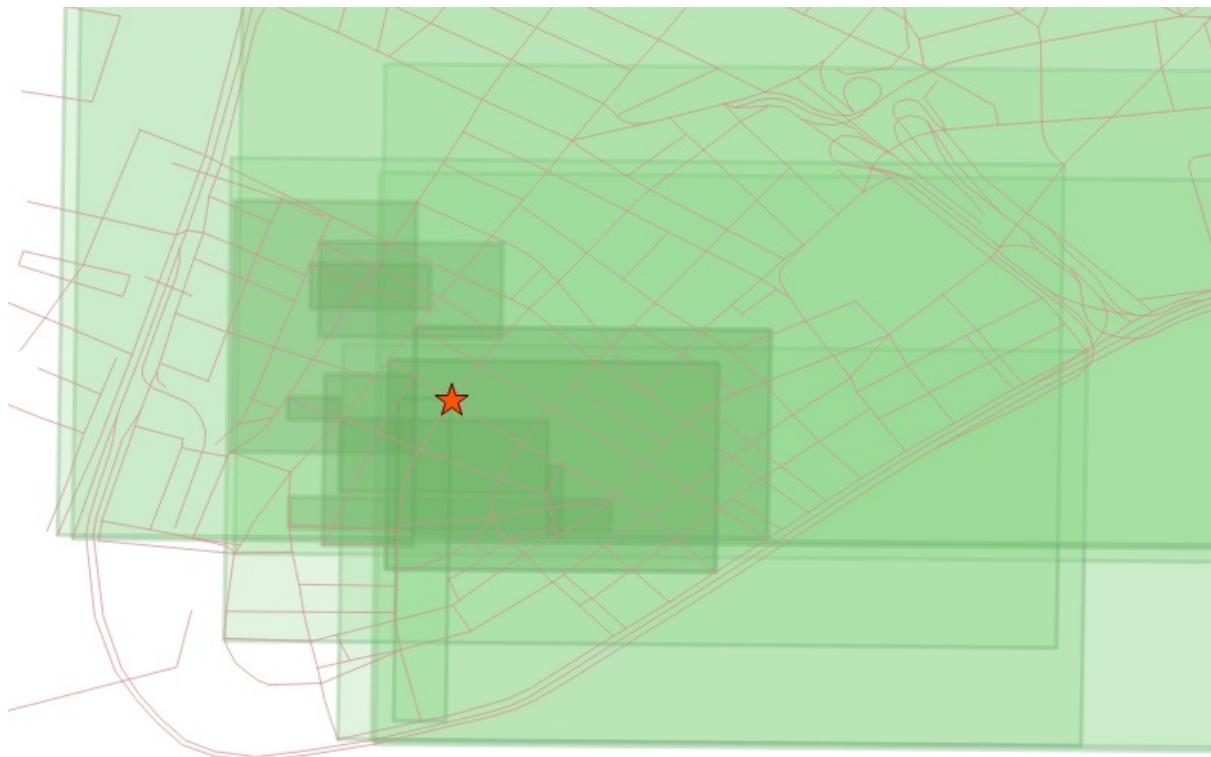
What about the `<#>` operator? If we calculate the distance between box edges, the station would fall **inside** the Wall Street box, giving it a distance of zero and the first entry in the list, right?

```
-- Closest 10 streets to Broad Street station are ?
SELECT
    streets.id,
    streets.name
FROM
    nyc_streets streets
ORDER BY
    streets.geom <#>
    'SRID=26918;POINT(583571.905921312 4506714.34119218)::geometry'
LIMIT 10;
```

Unfortunately, no.

id	name
17315	Pearl St
17364	South St
17394	Wall St
17411	Pine St
17378	FDR Dr
17236	
17241	West Side Highway; West St; West Side Highway; West Side Highway
17322	Broadway
17382	FDR Dr
17399	Broad St

There are a number of large street features with big boxes that **also** overlap the station and yield a box distance of zero.



This may not give the results we were expecting but since it operates on bounding boxes, it provides better performance than the query using `<->`. In case we had a very large table we could enhance the query processing by limiting our search first using the operator `<#>` and then use the operator `<->` in the resulting subset to get the accurate nearest neighbours:

```
-- "Closest" 100 streets to Broad Street station are?
WITH closest_candidates AS (
  SELECT
    streets.id,
    streets.name,
    streets.geom
  FROM
    nyc_streets streets
  ORDER BY
    streets.geom <#>
```

```
'SRID=26918;POINT(583571.905921312 4506714.34119218)::geometry
LIMIT 100
)
SELECT id, name
FROM closest_candidates
ORDER BY
geom <->
'SRID=26918;POINT(583571.905921312 4506714.34119218)::geometry
LIMIT 1;
```

id	name
17394	Wall St

knn  $\leftrightarrow$ : returns the 2D distance between two geometries.

knn  $\#$ : returns the distance between two bounding boxes.

## 5.9 Tracking Edit History using Triggers

A common requirement for production databases is the ability to track history: how has the data changed between two dates, who made the changes, and where did they occur? Some GIS systems track changes by including change management in the client interface, but that adds a lot of complexity to editing tools.

Using the database and the trigger system, it's possible to add history tracking to any table, while maintaining simple "direct edit" access to the primary table.

History tracking works by keeping a history table that records, for every edit:

- If a record was created, when it was added and by whom.
- If a record was deleted, when it was deleted and by whom.
- If a record was updated, adding a deletion record (for the old state) and a creation record (for the new state).

Using this information it is possible to reconstruct the state of the edit table at any point in time. In this example, we will add history tracking to our **nyc\_streets** table.

- First, add a new **nyc\_streets\_history** table. This is the table we will use to store all the historical edit information. In addition to all the fields from **nyc\_streets**, we add five more fields.
  - **hid** the primary key for the history table
  - **created** the date/time the history record was created
  - **created\_by** the database user that caused the record to be created
  - **deleted** the date/time the history record was marked as deleted
  - **deleted\_by** the database user that caused the record to be marked as deleted

Note that we don't actually delete any records in the history table, we just mark the time they ceased to be part of the current state of the edit table.

```
CREATE TABLE nyc_streets_history (
    hid SERIAL PRIMARY KEY,
    id FLOAT8,
    name VARCHAR(200),
```

```

    oneway VARCHAR(10),
    type VARCHAR(50),
    geom GEOMETRY(MultiLinestring,26918),
    created TIMESTAMP,
    created_by VARCHAR(32),
    deleted TIMESTAMP,
    deleted_by VARCHAR(32)
);

```

- Next, we import the current state of the active table, `nyc_streets` into the history table, so we have a starting point to trace history from. Note that we fill in the creation time and creation user, but leave the deletion records NULL.

```

INSERT INTO nyc_streets_history
(id, name, oneway, type, geom, created, created_by)
SELECT id, name, oneway, type, geom, now(), current_user
FROM nyc_streets;

```

- Now we need three triggers on the active table, for INSERT, DELETE and UPDATE actions. First we create the trigger functions, then bind them to the table as triggers.

For an insert, we just add a new record into the history table with the creation time/user:

```

CREATE OR REPLACE FUNCTION nyc_streets_insert() RETURNS trigger AS
$$
BEGIN
    INSERT INTO nyc_streets_history
        (id, name, oneway, type, geom, created, created_by)
    VALUES
        (NEW.id, NEW.name, NEW.oneway, NEW.type, NEW.geom,
         current_timestamp, current_user);
    RETURN NEW;
END;
$$
LANGUAGE plpgsql;

CREATE TRIGGER nyc_streets_insert_trigger
AFTER INSERT ON nyc_streets
    FOR EACH ROW EXECUTE PROCEDURE nyc_streets_insert();

```

For a deletion, we just mark the currently active history record (the one with a NULL deletion time) as deleted:

```

CREATE OR REPLACE FUNCTION nyc_streets_delete() RETURNS trigger AS
$$
BEGIN
    UPDATE nyc_streets_history
        SET deleted = current_timestamp, deleted_by = current_user
        WHERE deleted IS NULL and id = OLD.id;
    RETURN NULL;
END;
$$
LANGUAGE plpgsql;

CREATE TRIGGER nyc_streets_delete_trigger
AFTER DELETE ON nyc_streets
    FOR EACH ROW EXECUTE PROCEDURE nyc_streets_delete();

```

For an update, we first mark the active history record as deleted, then insert a new record for the updated state:

```
CREATE OR REPLACE FUNCTION nyc_streets_update() RETURNS trigger AS
$$
BEGIN

    UPDATE nyc_streets_history
        SET deleted = current_timestamp, deleted_by = current_user
        WHERE deleted IS NULL AND id = OLD.id;

    INSERT INTO nyc_streets_history
        (id, name, oneway, type, geom, created, created_by)
    VALUES
        (NEW.id, NEW.name, NEW.oneway, NEW.type, NEW.geom,
         current_timestamp, current_user);

    RETURN NEW;

END;
$$
LANGUAGE plpgsql;

CREATE TRIGGER nyc_streets_update_trigger
AFTER UPDATE ON nyc_streets
    FOR EACH ROW EXECUTE PROCEDURE nyc_streets_update();
```

Now that the history table is enabled, we can make edits on the main table and watch the log entries appear in the history table.

Note the power of this database-backed approach to history: **no matter what tool is used to make the edits, whether the SQL command line, a web-based JDBC tool, or a desktop tool like QGIS, the history is consistently tracked.**

Let's turn the two streets named "Cumberland Walk" to the more stylish "Cumberland Wynde":

Updating the two streets will cause the original streets to be marked as deleted in the history table, with a deletion time of now, and two new streets with the new name added, with an addition time of now. You can inspect the historical records:

Now that we have a history table, what use is it? It's useful for time travel! To travel to a particular time **T**, you need to construct a query that includes:

- All records created before T, and not yet deleted; and also
- All records created before T, but deleted **after T**.

We can use this logic to create a query, or a view, of the state of the data in the past. Since presumably all your test edits have happened in the past couple minutes, let's create a view of the history table that shows the state of the table 10 minutes ago, **before you started editing** (so, the original data).

```
-- State of history 10 minutes ago
-- Records must have been created at least 10 minute ago and
-- either be visible now (deleted is null) or deleted in the last hour

CREATE OR REPLACE VIEW nyc_streets_ten_min_ago AS
    SELECT * FROM nyc_streets_history
```

```
WHERE created < (now() - '10min'::interval)
AND ( deleted IS NULL OR deleted > (now() - '10min'::interval) );
```

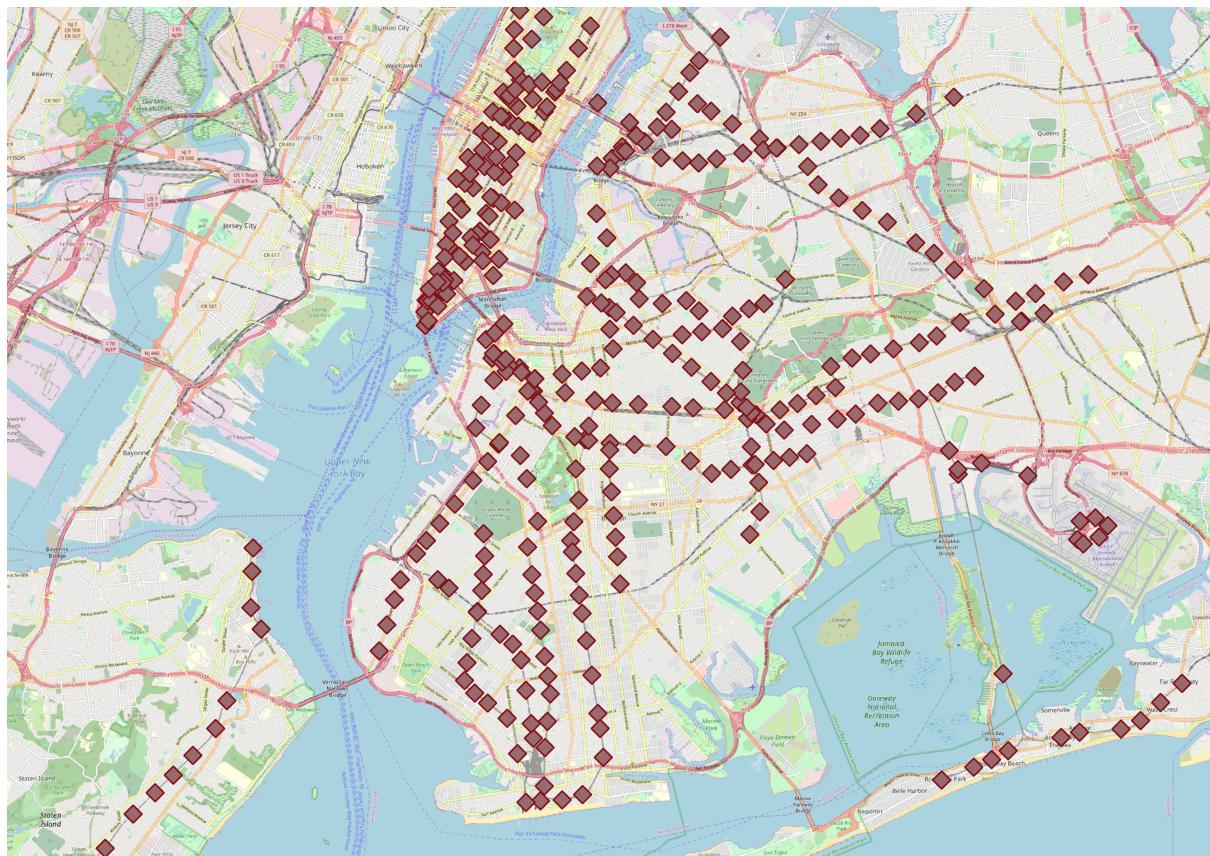
We can also create views that show just what a particular user has added, for example:

```
CREATE OR REPLACE VIEW nyc_streets_postgres AS
SELECT * FROM nyc_streets_history
WHERE created_by = 'postgres';
```

- QGIS open source GIS
- PostgreSQL Triggers

## 5.10 Advanced Geometry Constructions

The `nyc_subway_stations` layer has provided us with lots of interesting examples so far, but there is something striking about it:



Although it is a database of all the stations, it doesn't allow easy visualization of routes! In this chapter we will use advanced features of PostgreSQL and PostGIS to build up a new linear routes layer from the point layer of subway stations.

Our task is made especially difficult by two issues:

- The `routes` column of `nyc_subway_stations` has multiple route identifiers in each row, so a station that might appear in multiple routes appears only once in the table.
- Related to the previous issue, there is no route ordering information in the stations table, so while it is possible to find all the stations in a particular route, it's not possible using the attributes to

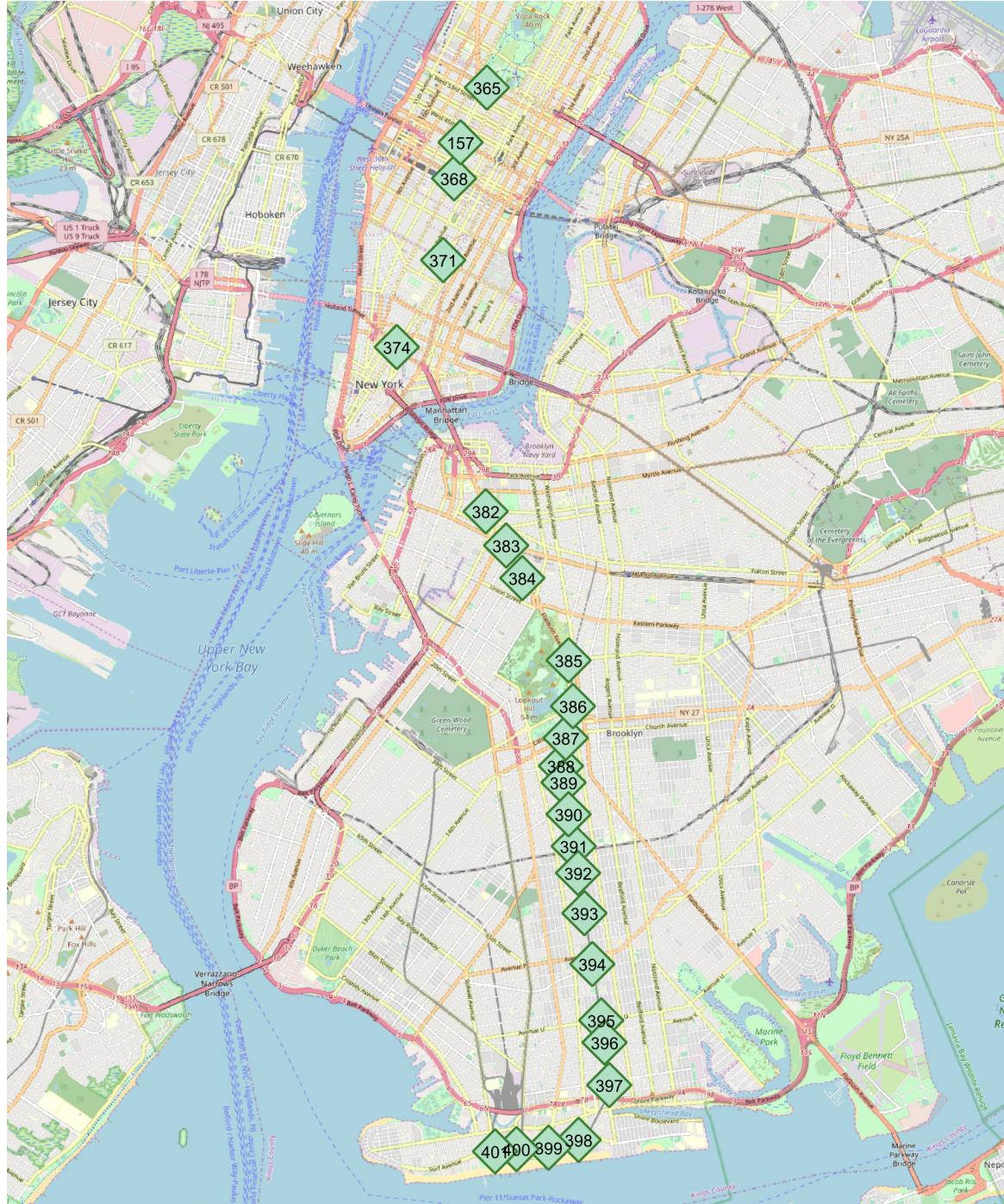
determine what the order in which trains travel through the stations.

The second problem is the harder one: given an unordered set of points in a route, how do we order them to match the actual route.

Here are the stops for the 'Q' train:

```
SELECT s.id, s.geom
FROM nyc_subway_stations s
WHERE (strpos(s.routes, 'Q') <> 0);
```

In this picture, the stops are labelled with their unique id primary key.



If we start at one of the end stations, the next station on the line seems to always be the closest. We can

repeat the process each time as long as we exclude all the previously found stations from our search.

There are two ways to run such an iterative routine in a database:

- Using a procedural language, like PL/PgSQL.
- Using recursive common table expressions.

Common table expressions (CTE) have the virtue of not requiring a function definition to run. Here's the CTE to calculate the route line of the 'Q' train, starting from the northernmost stop (where `id` is 365).

```
WITH RECURSIVE next_stop(geom, idlist) AS (
  SELECT
    geom,
    ARRAY[id] AS idlist
  FROM nyc_subway_stations
  WHERE id = 365)
  UNION ALL
  (SELECT
    s.geom,
    array_append(n.idlist, s.id) AS idlist
  FROM nyc_subway_stations s, next_stop n
  WHERE strpos(s.routes, 'Q') != 0
  AND NOT n.idlist @> ARRAY[s.id]
  ORDER BY ST_Distance(n.geom, s.geom) ASC
  LIMIT 1)
)
SELECT geom, idlist FROM next_stop;
```

The CTE consists of two halves, unioned together:

- The first half establishes a start point for the expression. We get the initial geometry and initialize the array of visited identifiers, using the record of "id" 365 (the end of the line).
- The second half iterates until it finds no further records. At each iteration it takes in the value at the previous iteration via the self-reference to "next\_stop". We search every stop on the Q line (`strpos(s.routes,'Q')`) that we have not already added to our visited list (`NOT n.idlist @> ARRAY[s.id]`) and order them by their distance from the previous point, taking just the first one (the nearest).

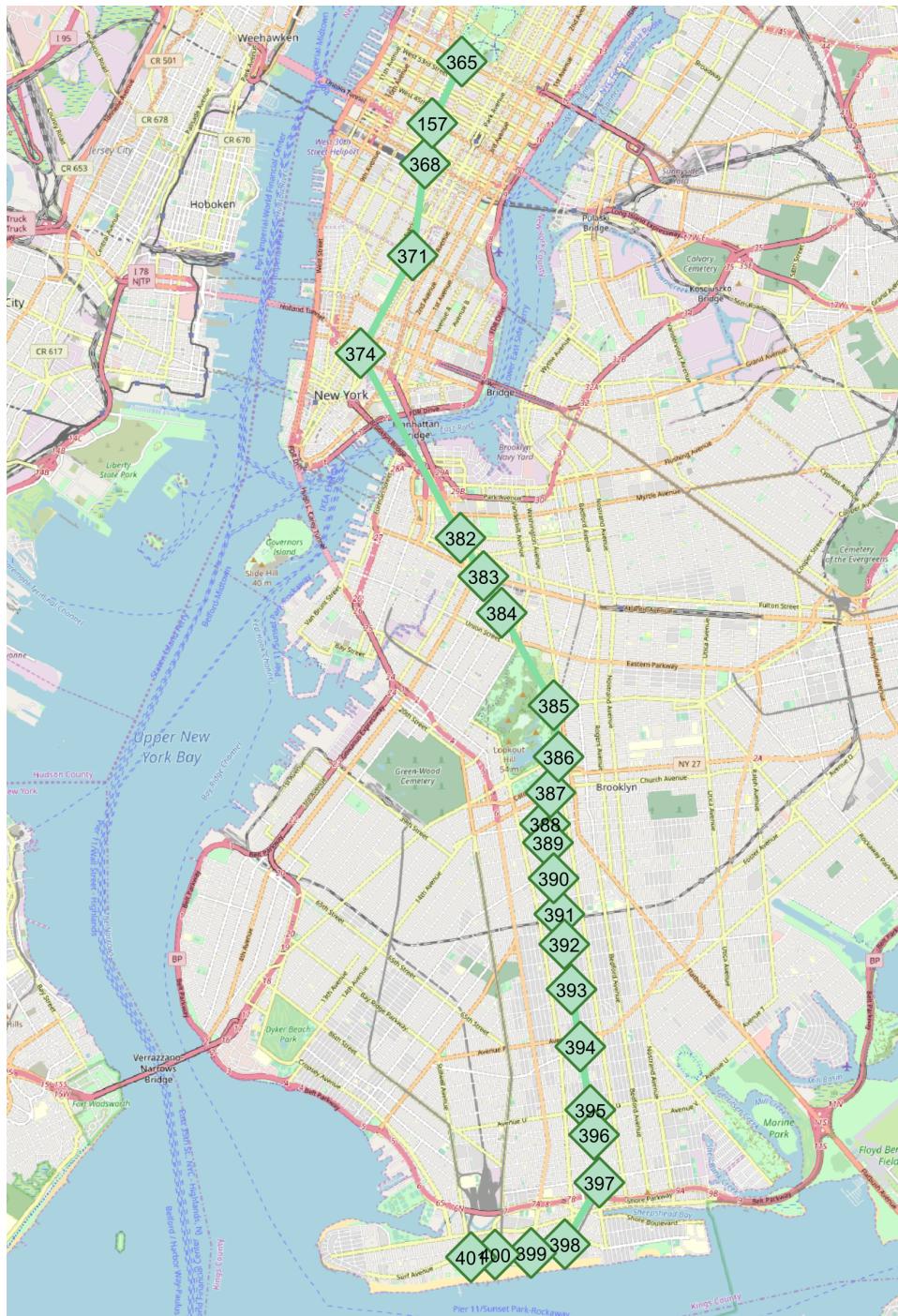
Beyond the recursive CTE itself, there are a number of advanced PostgreSQL array features being used here:

- We are using `ARRAY`! PostgreSQL supports arrays of any type. In this case we have an array of integers, but we could also build an array of geometries, or any other PostgreSQL type.
- We are using `array_append` to build up our array of visited identifiers.
- We are using the `@>` array operator ("array contains") to find which of the Q train stations we have already visited. The `@>` operators requires `ARRAY` values on both sides, so we have to turn the individual "id" numbers into single-entry arrays using the `ARRAY[]` syntax.

When you run the query, you get each geometry in the order it is found (which is the route order), as well as the list of identifiers already visited. Wrapping the geometries into the PostGIS `ST_MakeLine` aggregate function turns the set of geometries into a single linear output, constructed in the provided order.

```
WITH RECURSIVE next_stop(geom, idlist) AS (
  (SELECT
    geom,
    ARRAY[gid] AS idlist
   FROM nyc_subway_stations
  WHERE gid = 304)
 UNION ALL
  (SELECT
    s.geom,
    array_append(n.idlist, s.gid) AS idlist
   FROM nyc_subway_stations s, next_stop n
  WHERE strpos(s.routes, 'Q') != 0
  AND NOT n.idlist @> ARRAY[s.gid]
 ORDER BY ST_Distance(n.geom, s.geom) ASC
 LIMIT 1)
)
SELECT ST_MakeLine(geom) AS geom FROM next_stop;
```

Which looks like this:



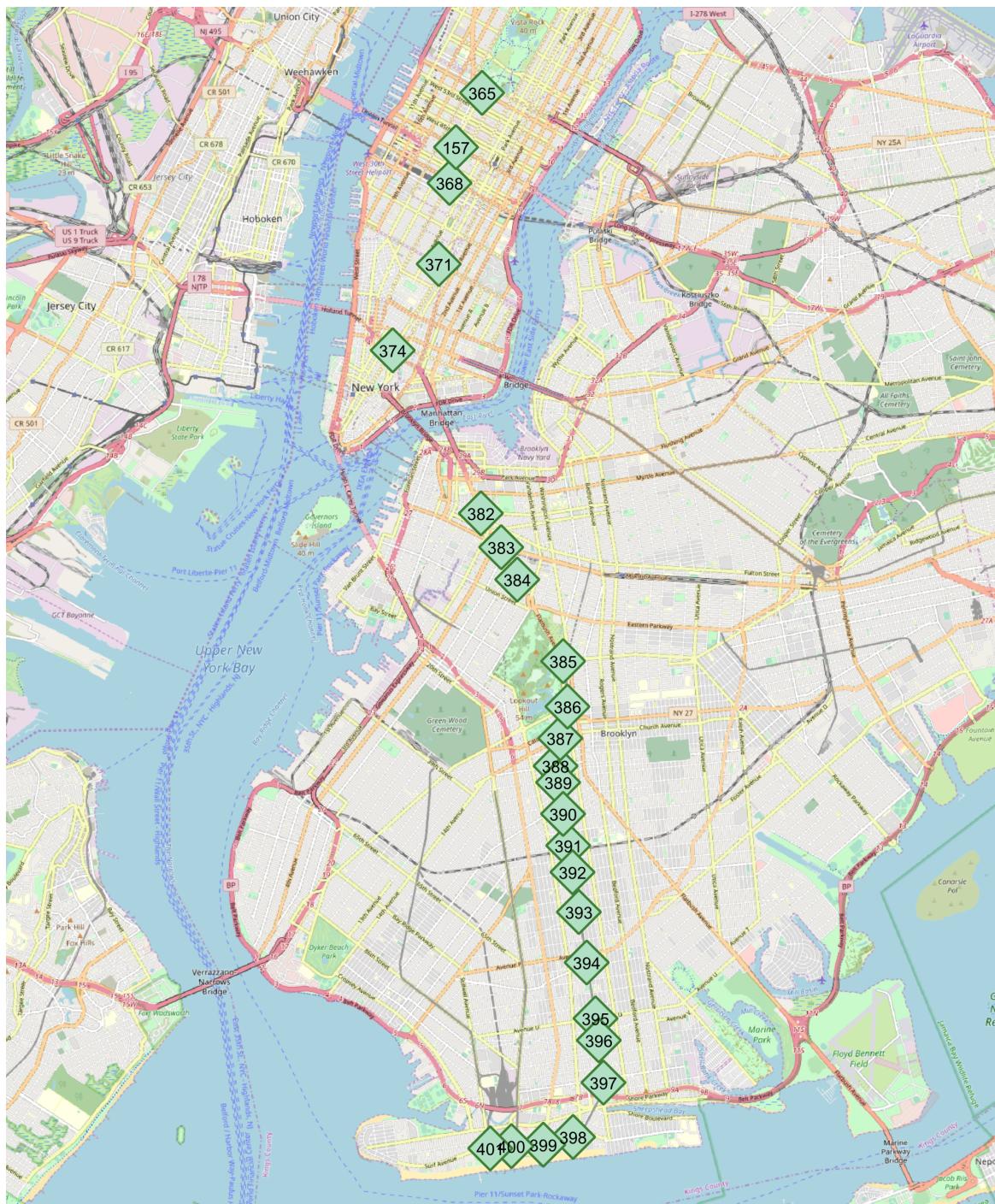
*Success!*

Except, two problems:

- We are only calculating one subway route here, we want to calculate all the routes.
- Our query includes a piece of *a priori* knowledge, the initial station identifier that serves as the seed for the search algorithm that builds the route.

Let's tackle the hard problem first, figuring out the first station on a route without manually eyeballing the set of stations that make up the route.

Our 'Q' train stops can serve as a starting point. What characterizes the end stations of the route?



One answer is “they are the most northerly and southerly stations”. However, imagine if the ‘Q’ train ran from east to west. Would the condition still hold?

A less directional characterization of the end stations is “they are the furthest stations from the middle of the route”. With this characterization it doesn’t matter if the route runs north/south or east/west, just that it run in more-or-less one direction, particularly at the ends.

Since there is no 100% heuristic to figure out the end points, let’s try this second rule out.

---

**Note:** An obvious failure mode of the “furthest from middle” rule is a circular line, like the Circle Line in London, UK. Fortunately, New York doesn’t have any such lines!

---

To work out the end stations of every route, we first have to work out what routes there are! We find the distinct routes.

```
WITH routes AS (
  SELECT DISTINCT unnest(string_to_array(routes, ',')) AS route
  FROM nyc_subway_stations ORDER BY route
)
SELECT * FROM routes;
```

Note the use of two advanced PostgreSQL ARRAY functions:

- **string\_to\_array** takes in a string and splits it into an array using a separator character. PostgreSQL supports arrays of any type, so it's possible to build arrays of strings, as in this case, but also of geometries and geographies as we'll see later in this example.
- **unnest** takes in an array and builds a new row for each entry in the array. The effect is to take a “horizontal” array embedded in a single row and turn it into a “vertical” array with a row for each value.

The result is a list of all the unique subway route identifiers.

route
-----
1
2
3
4
5
6
7
A
B
C
D
E
F
G
J
L
M
N
Q
R
S
V
W
Z
(24 rows)

We can build on this result by joining it back to the `nyc_subway_stations` table to create a new table that has, for each route, a row for every station on that route.

```
WITH routes AS (
  SELECT DISTINCT unnest(string_to_array(routes, ',')) AS route
  FROM nyc_subway_stations ORDER BY route
),
stops AS (
  SELECT s.id, s.geom, r.route
```

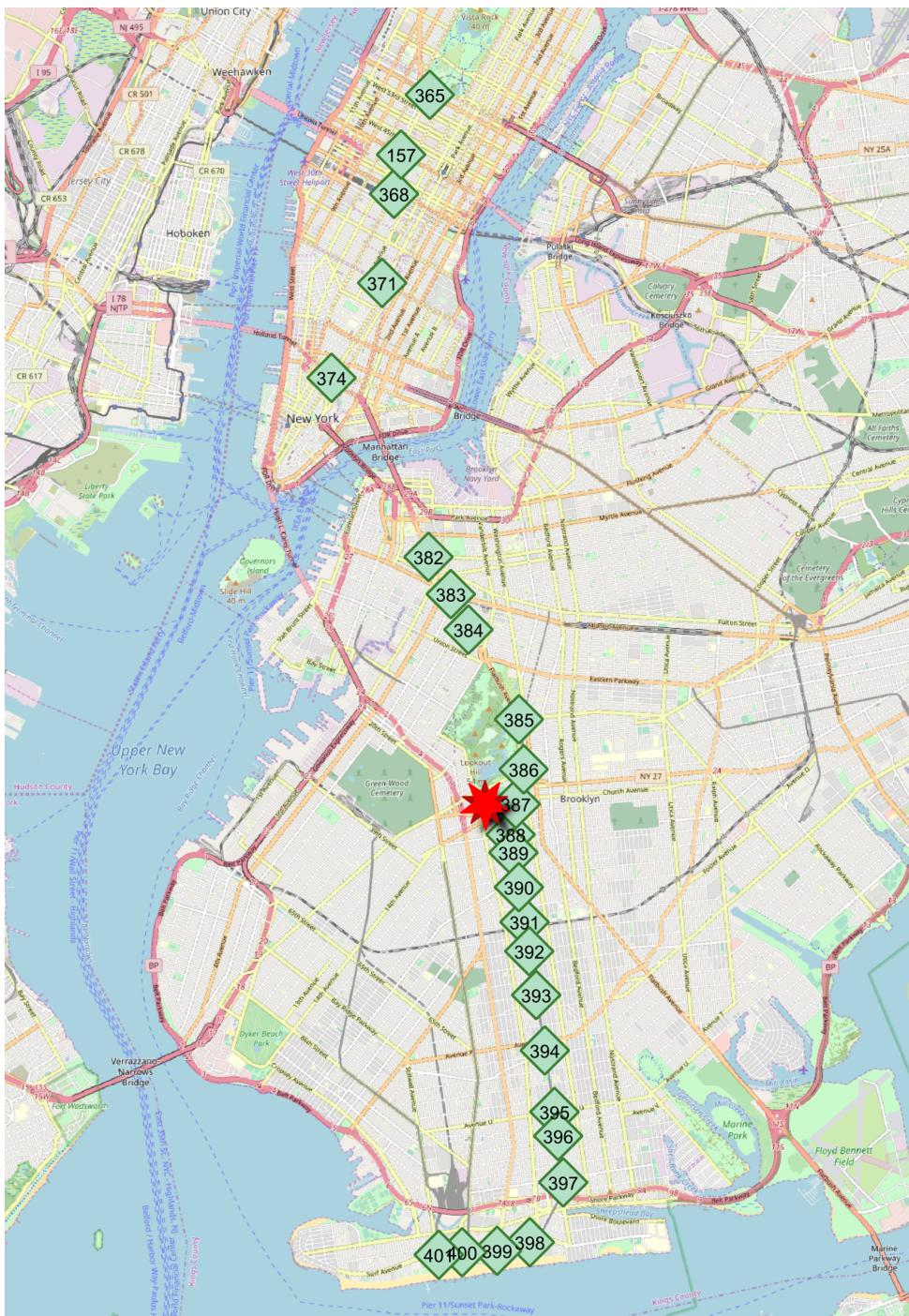
```
FROM routes r
JOIN nyc_subway_stations s
ON (strpos(s.routes, r.route) <> 0)
)
SELECT * FROM stops;
```

id	geom	route
2	010100002026690000CBE327F938CD21415EDBE1572D315141	1
1	010100002026690000C676635D10CD2141A0ECDB6975305141	1
36	010100002026690000AE59A3F82C132241D835BA14D1435141	1
37	0101000020266900003495A303D615224116DA56527D445141	1
...etc...		

Now we can find the center point by collecting all the stations for each route into a single multi-point, and calculating the centroid of that multi-point.

```
WITH routes AS (
  SELECT DISTINCT unnest(string_to_array(routes, ',')) AS route
  FROM nyc_subway_stations ORDER BY route
),
stops AS (
  SELECT s.id, s.geom, r.route
  FROM routes r
  JOIN nyc_subway_stations s
  ON (strpos(s.routes, r.route) <> 0)
),
centers AS (
  SELECT ST_Centroid(ST_Collect(geom)) AS geom, route
  FROM stops
  GROUP BY route
)
SELECT * FROM centers;
```

The center point of the collection of ‘Q’ train stops looks like this:



So the northern most stop, the end point, appears to also be the stop furthest from the center. Let's calculate the furthest point for every route.

```
WITH routes AS (
  SELECT DISTINCT unnest(string_to_array(routes, ',')) AS route
  FROM nyc_subway_stations ORDER BY route
),
stops AS (
  SELECT s.id, s.geom, r.route
  FROM routes r
  JOIN nyc_subway_stations s
  ON (strpos(s.routes, r.route) <> 0)
),
```

```
centers AS (
    SELECT ST_Centroid(ST_Collect(geom)) AS geom, route
    FROM stops
    GROUP BY route
),
stops_distance AS (
    SELECT s.*, ST_Distance(s.geom, c.geom) AS distance
    FROM stops s JOIN centers c
    ON (s.route = c.route)
    ORDER BY route, distance DESC
),
first_stops AS (
    SELECT DISTINCT ON (route) stops_distance.*
    FROM stops_distance
)
SELECT * FROM first_stops;
```

We've added two sub-queries this time:

- **stops\_distance** joins the centers points back to the stations table and calculates the distance between the stations and center for each route. The result is ordered such that the records come out in batches for each route, with the furthest station as the first record of the batch.
- **first\_stops** filters the **stops\_distance** output by only taking the first record for each distinct group. Because of the way we ordered **stops\_distance** the first record is the furthest record, which means it's the station we want to use as our starting seed to build each subway route.

Now we know every route, and we know (approximately) what station each route starts from: we're ready to generate the route lines!

But first, we need to turn our recursive CTE expression into a function we can call with parameters:

```
CREATE OR REPLACE function walk_subway(integer, text) returns geometry AS
$$
WITH RECURSIVE next_stop(geom, idlist) AS (
    (SELECT
        geom AS geom,
        ARRAY[id] AS idlist
    FROM nyc_subway_stations
    WHERE id = $1)
    UNION ALL
    (SELECT
        s.geom AS geom,
        array_append(n.idlist, s.id) AS idlist
    FROM nyc_subway_stations s, next_stop n
    WHERE strpos(s.routes, $2) != 0
    AND NOT n.idlist @> ARRAY[s.id]
    ORDER BY ST_Distance(n.geom, s.geom) ASC
    LIMIT 1)
)
SELECT ST_MakeLine(geom) AS geom
FROM next_stop;
$$
language 'sql';
```

And now we are ready to go!

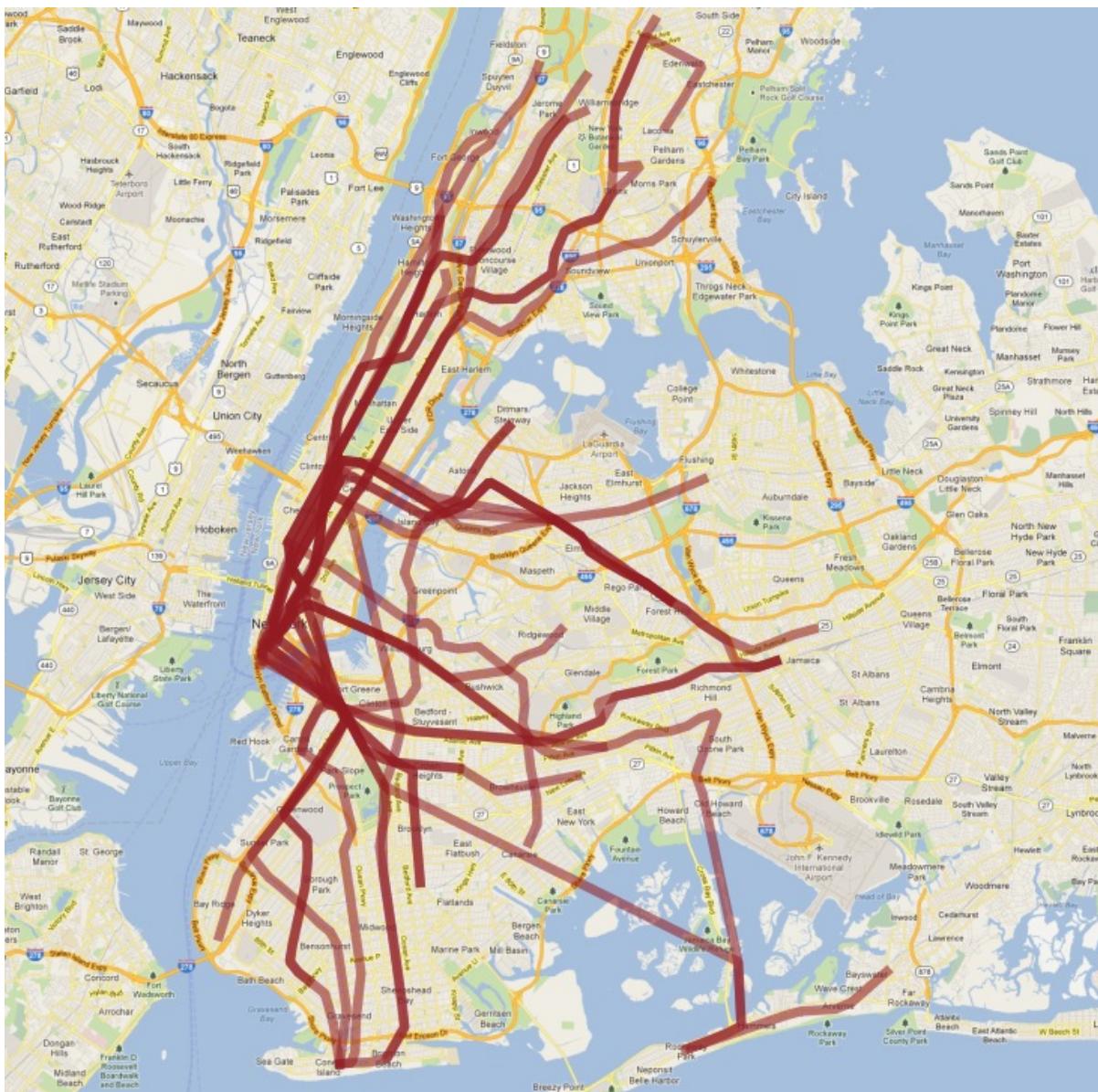
```

CREATE TABLE nyc_subway_lines AS
-- Distinct route identifiers!
WITH routes AS (
  SELECT DISTINCT unnest(string_to_array(routes,',')) AS route
  FROM nyc_subway_stations ORDER BY route
),
-- Joined back to stops! Every route has all its stops!
stops AS (
  SELECT s.id, s.geom, r.route
  FROM routes r
  JOIN nyc_subway_stations s
  ON (strpos(s.routes, r.route) <> 0)
),
-- Collects stops by routes and calculate centroid!
centers AS (
  SELECT ST_Centroid(ST_Collect(geom)) AS geom, route
  FROM stops
  GROUP BY route
),
-- Calculate stop/center distance for each stop in each route.
stops_distance AS (
  SELECT s.*, ST_Distance(s.geom, c.geom) AS distance
  FROM stops s JOIN centers c
  ON (s.route = c.route)
  ORDER BY route, distance DESC
),
-- Filter out just the furthest stop/center pairs.
first_stops AS (
  SELECT DISTINCT ON (route) stops_distance.*
  FROM stops_distance
)
-- Pass the route/stop information into the linear route generation
-- function!
SELECT
  ascii(route) AS id, -- QGIS likes numeric primary keys
  route,
  walk_subway(id::Integer, route) AS geom
FROM first_stops;

-- Do some housekeeping too
ALTER TABLE nyc_subway_lines ADD PRIMARY KEY (id);

```

Here's what our final table looks like visualized in QGIS:



As usual, there are some problems with our simple understanding of the data:

- there are actually two ‘S’ (short distance “shuttle”) trains, one in Manhattan and one in the Rockaways, and we join them together because they are both called ‘S’;
- the ‘4’ train (and a few others) splits at the end of one line into two terminuses, so the “follow one line” assumption breaks and the result has a funny hook on the end.

Hopefully this example has provided a taste of some of the complex data manipulations that are possible combining the advanced features of PostgreSQL and PostGIS.

### 5.10.1 See Also

- PostgreSQL Arrays
- PostgreSQL Array Functions
- PostgreSQL Recursive Common TABLE Expressions
- PostGIS ST\_MakeLine

## MAINTENANCE

### 6.1 PostgreSQL Security

PostgreSQL has a rich and flexible permissions system, with the ability to parcel out particular privileges to particular `roles`, and provide users with the powers of one or more of those `roles`.

In addition, the PostgreSQL server can use multiple different systems to authenticate users. This means that the database can use the same authentication infrastructure as other architecture components, simplifying password management.

#### 6.1.1 Users and Roles

In this chapter we will create two useful production users:

- A read-only user for use in a publishing application.
- A read/write user for use by a developer in building a software or analyzing data.

Rather than creating users and granting them the necessary powers, we will create two roles with the right powers and then create two users and add them to the appropriate roles. That way we can easily reuse the roles when we create further users.

#### Creating Roles

A role is a user and a user is a role. The only difference is that a “user” can be said to be a role with the “login” privilege.

So functionally, the two SQL statements below are the same, they both create a “role with the login privilege”, which is to say, a “user”.

```
CREATE ROLE mrbean LOGIN;  
CREATE USER mrbean;
```

#### Read-only Users

Our read-only user will be for a web application to use to query the `nyc_streets` table.

The application will have specific access to the `nyc_streets` table, but will also inherit the necessary system access for PostGIS operations from the `postgis_reader` role.

```
-- A user account for the web app  
CREATE USER app1;  
-- Web app needs access to specific data tables  
GRANT SELECT ON nyc_streets TO app1;
```

```
-- A generic role for access to PostGIS functionality
CREATE ROLE postgis_reader INHERIT;
-- Give that role to the web app
GRANT postgis_reader TO app1;
```

Now, use the psql command line to login as as app1:

```
psql -d nyc -U app1
```

Try the following queries:

```
-- This works!
SELECT * FROM nyc_streets LIMIT 1;

-- This works too!
SELECT ST_AsText(ST_Transform(geom, 4326))
  FROM nyc_streets LIMIT 1;
```

The `INHERIT` clause granted `postgis_reader` with the privileges that all the roles in the database have. Now we have a nice generic `postgis_reader` role we can apply to any user that need to read from PostGIS tables.

### Read/write Users

There are two kinds of read/write scenarios we need to consider:

- Web applications and others that need to write to existing data tables.
- Developers or analysts that need to create new tables and geometry columns as part of their work.

For web applications that require write access to data tables, we just need to grant extra permissions to the tables themselves, and we can continue to use the `postgis_reader` role.

```
-- Add insert/update/delete abilities to our web application
GRANT INSERT, UPDATE, DELETE ON nyc_streets TO app1;
```

These kinds of permissions would be required for a read/write WFS service, for example.

For developers and analysts, a little more access is needed to the main PostGIS metadata tables. We will need a `postgis_writer` role that can edit the PostGIS metadata tables!

```
-- Make a postgis writer role
CREATE ROLE postgis_writer;

-- Start by giving it the postgis_reader powers
GRANT postgis_reader TO postgis_writer;

-- Add insert/update/delete powers for the PostGIS tables
GRANT INSERT, UPDATE, DELETE ON spatial_ref_sys TO postgis_writer;

-- Make app1 a PostGIS writer to see if it works!
GRANT postgis_writer TO app1;
```

Now try the table creation SQL above as the `app1` user and see how it goes!

### 6.1.2 Encryption

PostgreSQL provides a lot of [encryption facilities](#), many of them optional, some of them on by default.

- By default, all passwords are MD5 encrypted. The client/server handshake double encrypts the MD5 password to prevent re-use of the hash by anyone who intercepts the password. It is also possible to use SCRAM which is an Internet standard.
- [SSL connections](#) are optionally available between the client and server, to encrypt all data and login information. SSL certificate authentication is also available when SSL connections are used.
- Columns inside the database can be encrypted using the [pgcrypto](#) module, which includes hashing algorithms, direct ciphers (blowfish, aes) and both public key and symmetric PGP encryption.

### SSL Connections

In order to use SSL connections, both your client and server must support SSL. Your version of PostgreSQL may have SSL support built, but not enabled. If so, we have to carry out a few steps to turn it on first.

- First, turn shut down PostgreSQL, since activating SSL will require a restart.
- Next, we acquire or generate an SSL certificate and key. The certificate will need to have no passphrase on it, or the database server won't be able to start up. You can generate a self-signed key as follows:

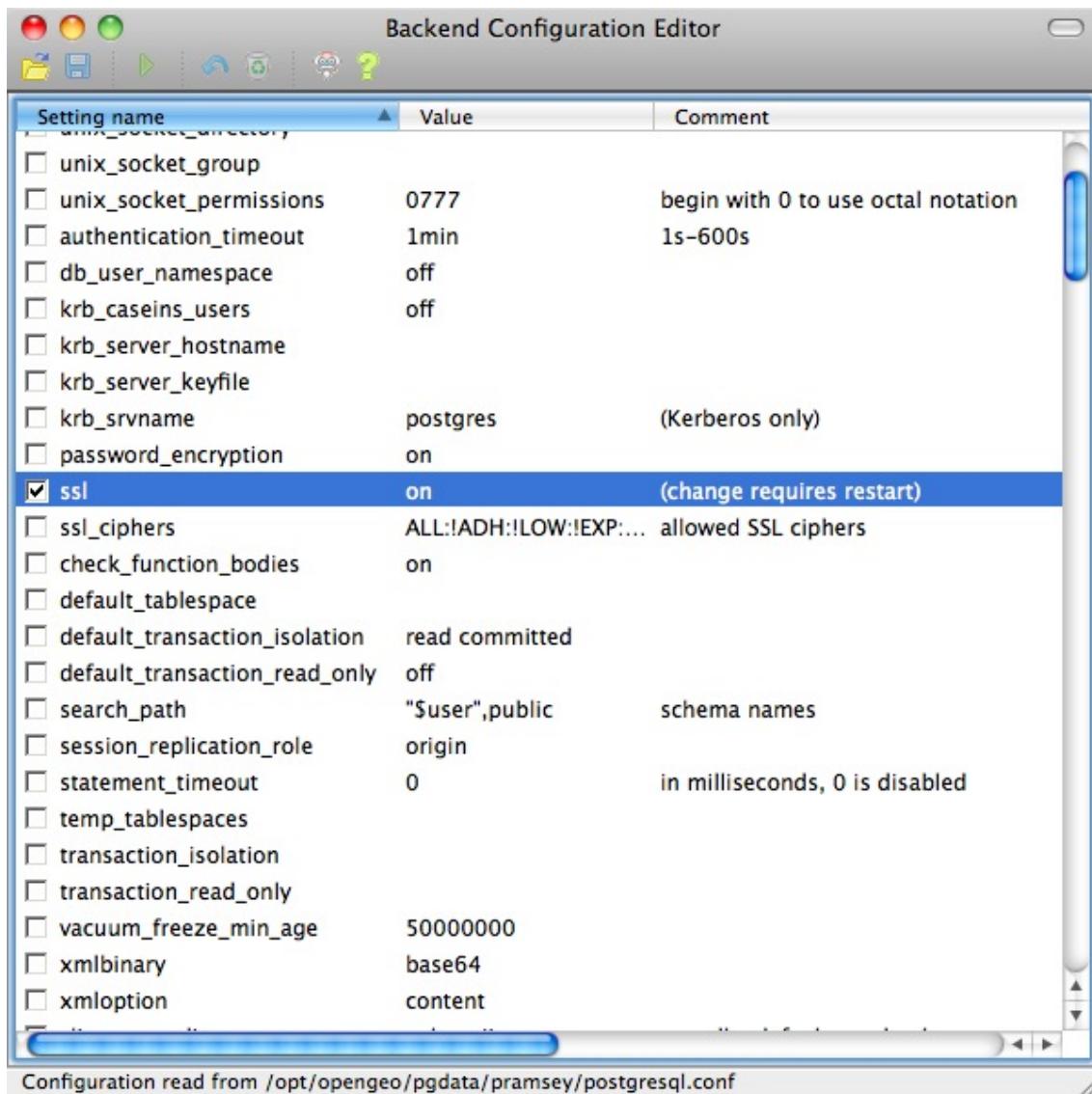
```
# Create a new certificate, filling out the certification info as prompted
openssl req -new -text -out server.req

# Strip the passphrase from the certificate
openssl rsa -in privkey.pem -out server.key

# Convert the certificate into a self-signed cert
openssl req -x509 -in server.req -text -key server.key -out server.crt

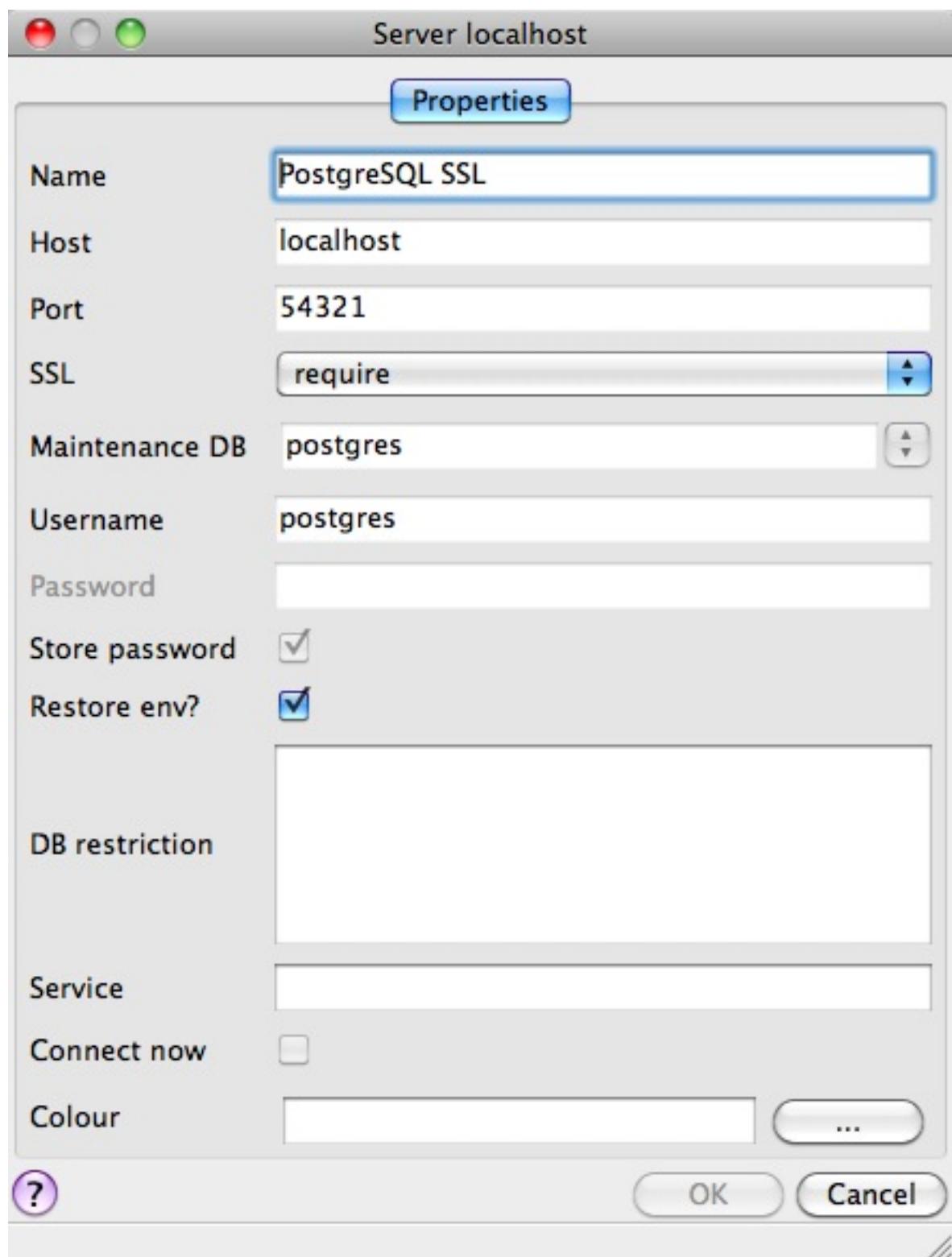
# Set the permission of the key to private read/write
chmod og-rwx server.key
```

- Copy the server.crt and server.key into the PostgreSQL data directory.
- Enable SSL support in the postgresql.conf file by turning the “ssl” parameter to “on”.



- Now re-start PostgreSQL. The server is now ready for SSL operation.

With the server enabled for SSL, creating an encrypted connection is easy. In PgAdmin, create a new server connection (File > Add Server...), and set the SSL parameter to “require”.



Once you connect with the new connection, you can see in its properties that it is using an SSL connection.

Properties		Statistics	Dependencies	Dependents
<b>Property</b>				<b>Value</b>
Description	PostgreSQL SSL			
Hostname	localhost			
Port	54321			
Encryption	SSL encrypted			
Maintenance database	postgres			
Username	postgres			
Store password?	Yes			
Restore environment?	Yes			
Version string	PostgreSQL 8.4.9 on i386-apple-darwin, compiled by GCC i686-apple-darwin			
Version number	8.4			
Last system OID	11563			
Connected?	Yes			
Up since	08/08/2012 11:51:37			
Autovacuum	running			

Since the default SSL connection mode is “prefer”, you don’t even need to specify an SSL preference when connecting. A connection with the command line `psql` terminal will pick up the SSL option and use it by default:

```
psql (8.4.9)
SSL connection (cipher: DHE-RSA-AES256-SHA, bits: 256)
Type "help" for help.

postgres=#
```

Note how the terminal reports the SSL status of the connection.

### Data Encryption

The `pgcrypto` module has a huge range of encryption options, so we will only demonstrate the simplest use case: encrypting a column of data using a symmetric cipher.

- First, enable `pgcrypto` by creating the EXTENSION, either in PgAdmin or `psql`.

```
CREATE EXTENSION IF NOT EXISTS pgcrypto;
```

- Then, test the encryption function.

```
-- encrypt a string using blowfish (bf)
SELECT encrypt('this is a test phrase', 'mykey', 'bf');
```

- And make sure it’s reversible too!

```
-- round-trip a string using blowfish (bf)
SELECT decrypt(encrypt('this is a test phrase', 'mykey', 'bf'), 'mykey
→ ', 'bf');
```

### 6.1.3 Authentication

PostgreSQL supports 14 different authentication methods, to allow easy integration into existing enterprise architectures. For production purposes, the following methods are commonly used:

- **Password** is the basic system where the passwords are stored by the database, with MD5 encryption.
- **Kerberos** is a standard enterprise authentication method, which is used by both the **GSSAPI** and **SSPI** schemes in PostgreSQL. Using **SSPI**, PostgreSQL can authenticate against Windows servers.
- **LDAP** is another common enterprise authentication method. The **OpenLDAP** server bundled with most Linux distributions provides an open source implementation of **LDAP**.
- **Certificate** authentication is an option if you expect all client connections to be via SSL and are able to manage the distribution of keys.
- **PAM** authentication is an option if you are on Linux or Solaris and use the **PAM** scheme for transparent authentication provision.

Authentication methods are controlled by the `pg_hba.conf` file. The “HBA” in the file name stands for “host based access”, because in addition to allowing you to specify the authentication method to use for each database, it allows you to limit host access using network addresses.

Here is an example `pg_hba.conf` file:

<code># TYPE</code>	<code>DATABASE</code>	<code>USER</code>	<code>CIDR-ADDRESS</code>	<code>METHOD</code>
<i># "local" is for Unix domain socket connections only</i>				
local	all	all		trust
<i># IPv4 local connections:</i>				
host	all	all	127.0.0.1/32	trust
<i># IPv6 local connections:</i>				
host	all	all	::1/128	trust
<i># remote connections for nyc database only</i>				
host	nyc	all	192.168.1.0/2	ldap

The file consists of five columns

- **TYPE** determines the kind of access, either “local” for connections from the same server or “host” for remote connections.
- **DATABASE** specifies what database the configuration line refers to or “all” for all databases
- **USER** specifies what users the line refers to or “all” for all users
- **CIDR-ADDRESS** specifies the network limitations for remote connections, using network/netmask syntax
- **METHOD** specifies the authentication protocol to use. “trust” skips authentication entirely and simply accepts any valid username without challenge.

It’s common for local connections to be trusted, since access to the server itself is usually privileged. Remote connections are disabled by default when PostgreSQL is installed: if you want to connect from remote machines, you’ll have to add an entry.

The line for `nyc` in the example above is an example of a remote access entry. The `nyc` example allows LDAP authenticated access only to machines on the local network (in this case the 192.168.1. network) and only to the `nyc` database. Depending on the security of your network, you will use more or less strict versions of these rules in your production set-up.

#### 6.1.4 Links

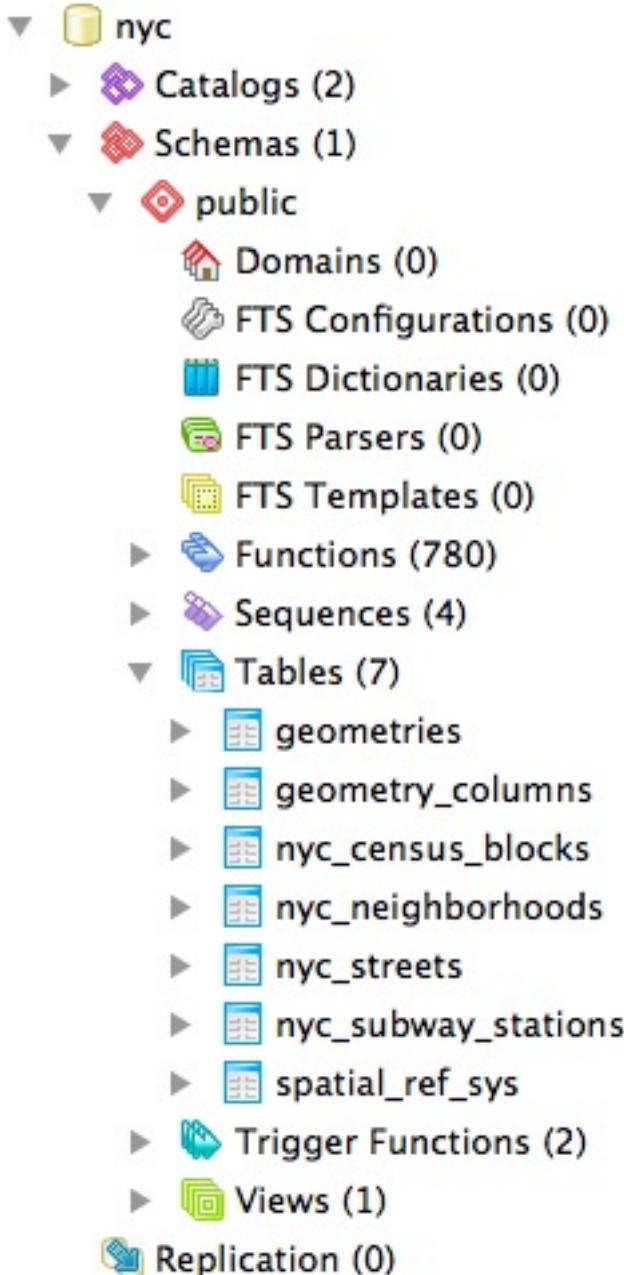
- PostgreSQL Authentication

- PostgreSQL Encryption
- PostgreSQL SSL Support

## 6.2 PostgreSQL Schemas

Production databases inevitably have a large number of tables and views, and managing them all in one schema can become unwieldy quickly. Fortunately, PostgreSQL includes the concept of a “\_Schema”.

Schemas are like folders, and can hold tables, views, functions, sequences and other relations. Every database starts out with one schema, the public schema.



Inside that schema, the default install of PostGIS creates the geometry\_columns, geography\_columns and spatial\_ref\_sys metadata relations, as well as all the types and functions used by PostGIS. So users of PostGIS always need access to the public schema.

In the public schema you can also see all the tables we have created so far in the workshop.

### 6.2.1 Why use Schemas?

There are two very good reasons for using schemas:

- Data that is managed in a schema is easier to apply bulk actions to.
  - It's easier to back-up data that's in a separate schema: so volatile data can have a different back-up schedule from non-volatile data.
  - It's easier to restore data that's in a separate schema: so application-oriented schemas can be separately restored and backed up for time travel and recovery.
  - It's easier to manage application differences when the application data is in a schema: so a new version of software can work off a table structure in a new schema, and cut-over involves a simple change to the schema name.
- Users can be restricted in their work to single schemas to allow isolation of analytical and test tables from production tables.

So for production purposes, keeping your application data separate in schemas improves management; and for user purposes, keeping your users in separate schemas keeps them from treading on each other.

### 6.2.2 Creating a Data Schema

Let's create a new schema and move a table into it. First, create a new schema in the database:

```
CREATE SCHEMA census;
```

Next, we will move the `nyc_census_blocks` table to the `census` schema:

```
ALTER TABLE nyc_census_blocks SET SCHEMA census;
```

If you're using the `psql` command-line program, you'll notice that `nyc_census_blocks` has disappeared from your table listing now! If you're using PgAdmin, you might have to refresh your view to see the new schema and the table inside it.

You can access tables inside schemas in two ways:

- by referencing them using `schema.table` notation
- by adding the schema to your `search_path`

Explicit referencing is easy, but it gets tiring to type after a while:

```
SELECT * FROM census.nyc_census_blocks LIMIT 1;
```

Manipulating the `search_path` is a nice way to provide access to tables in multiple schemas without lots of extra typing.

You can set the `search_path` at run time using the `SET` command:

```
SET search_path = census, public;
```

This ensures that all references to relations and functions are searched in both the `census` and the `public` schemas. Remember that all the PostGIS functions and types are in `public` so we don't want to drop that from the search path.

Setting the search path every time you connect can get tiring too, but fortunately it's possible to permanently set the search path for a user:

```
ALTER USER postgres SET search_path = census, public;
```

Now the `postgres` user will always have the `census` schema in their search path.

### 6.2.3 Creating a User Schema

Users like to create tables, and PostGIS users do so particularly: analysis operations with SQL demand temporary tables for visualization or interim results, so spatial SQL tends to require that users have `CREATE` privileges more than ordinary database workloads.

By default, every role in Oracle is given a personal schema. This is a nice practice to use for PostgreSQL users too, and is easy to replicate using PostgreSQL roles, schemas, and search paths.

Create a new user with table creation privileges (see security for information about the `postgis_writer` role), then create a schema with that user as the authorization:

```
CREATE USER myuser WITH ROLE postgis_writer;
CREATE SCHEMA myuser AUTHORIZATION myuser;
```

If you log in as that user, you'll find the default `search_path` for PostgreSQL is actually this:

```
show search_path;
```

```
search_path
```

```
-----
```

```
"$user",public
```

The first schema on the search path is the user's named schema! So now the following conditions exist:

- The user exists, with the ability to create spatial tables.
- The user's named schema exists, and the user owns it.
- The user's search path has the user schema first, so new tables are automatically created there, and queries automatically search there first.

That's all there is to it, the user's default work area is now nicely separated from any tables in other schemas.

## 6.3 PostgreSQL Backup and Restore

There are lots of ways to backup a PostgreSQL database, and the one you choose will depend a great deal on how you are using the database.

- For relatively static databases, the basic `pg_dump/pg_restore` tools can be used to take periodic snapshots of the data.
- For frequently changing data, using an “online backup” scheme allows continuous archiving of updates to a secure location.

Online backup is the basis for replication and stand-by systems for [high availability](#), particularly for versions of PostgreSQL  $\geq 9.0$ .

### 6.3.1 Laying Out your Data

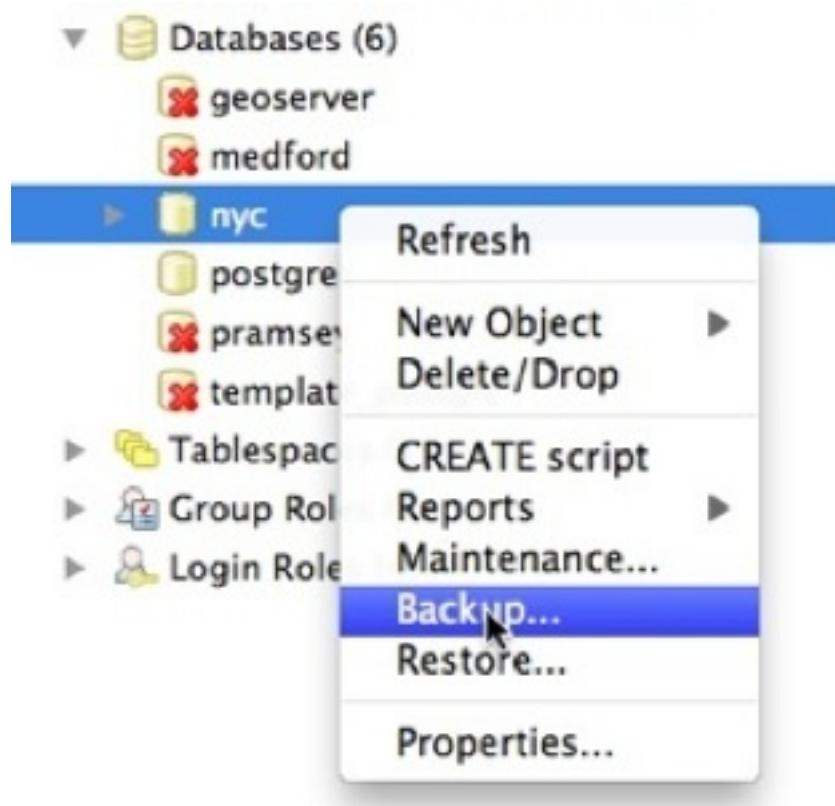
As discussed in schemas, ensuring that production data is always stored in separate schemas is a very important **best practice** in managing data. There are two reasons:

- Backing up and restoring data in schemas is much simpler than managing lists of tables to be backed up individually.
- Keeping data tables out of the “public” schema allows far easier upgrades, as discussed in upgrades.

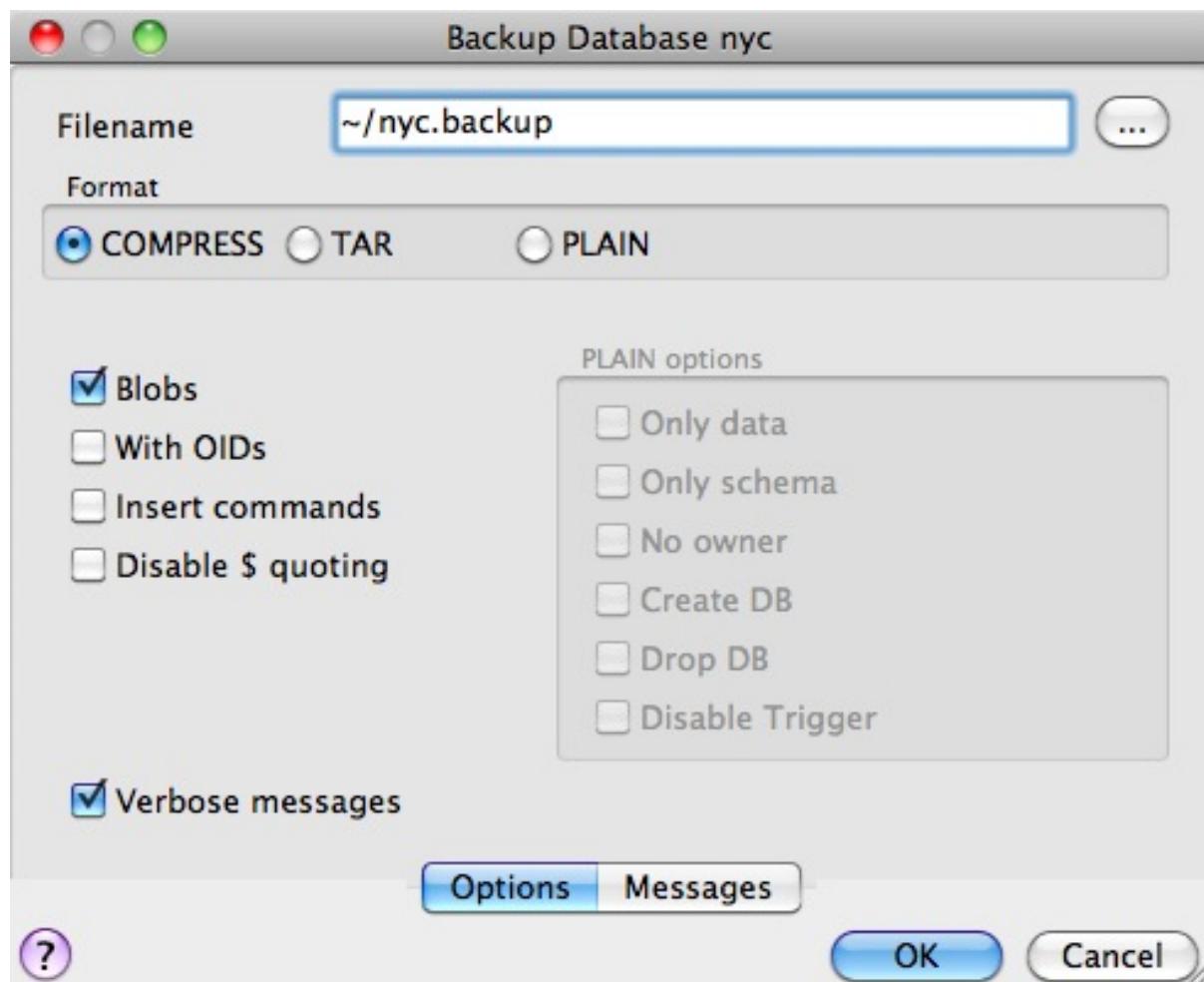
### 6.3.2 Basic Backup and Restore

Backing up a full database is easy using the `pg_dump` utility. The utility is a command-line tool, which makes it easy to automate with scripting, and it can also be invoke via a GUI in the PgAdmin utility.

To backup our `nyc` database, we can use the GUI, just right-click the database you want to backup:



Enter the name of the backup file you want to create.



Note that there are three backup format options: compress, tar and plain.

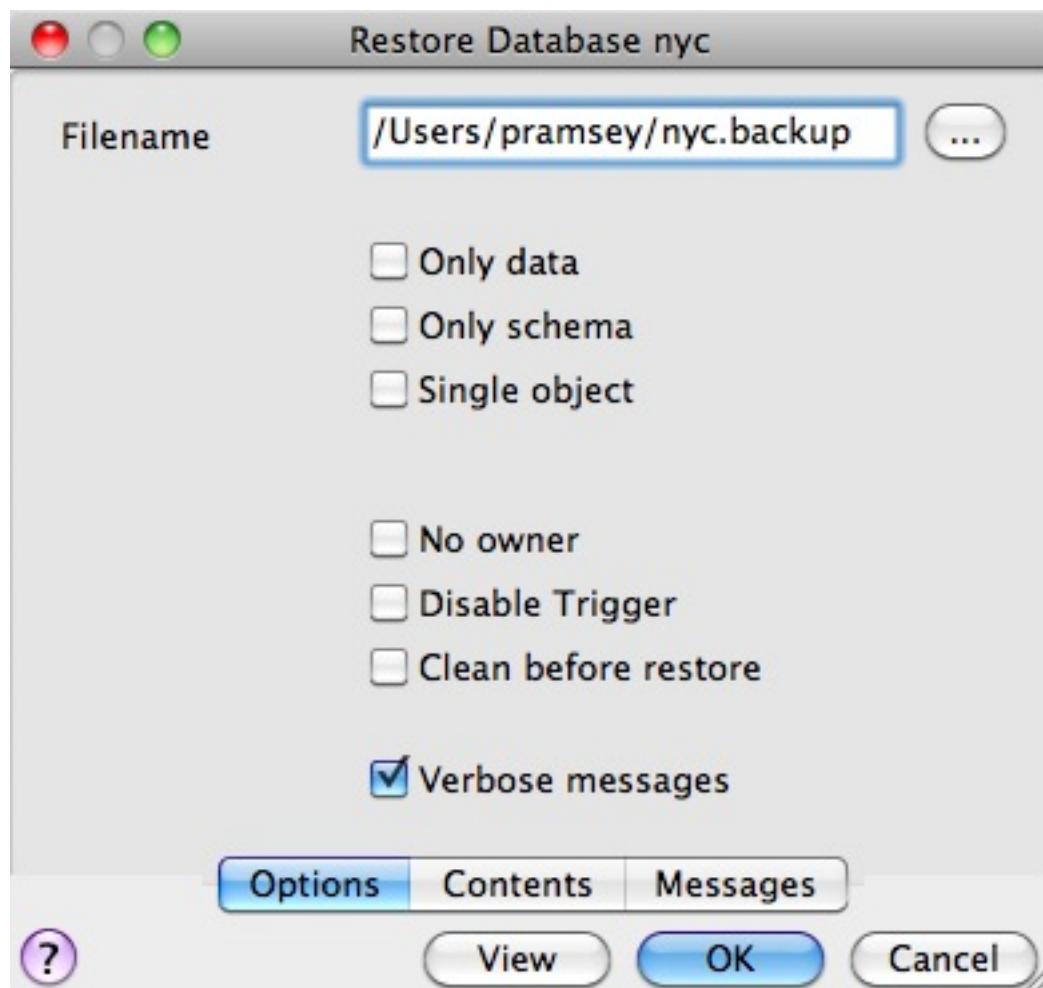
- **Plain** is just a textual SQL file. This is the simplest format and in many ways the most flexible, since it can be editing or altered easily and then loaded back into a database, allowing offline changes to things like ownership or other global information.
- **Tar** using a UNIX archive format to hold components of the dump in separate files. Using the tar format allows the `pg_restore` utility to selectively restore parts of the dump.
- **Compress** is like the Tar format, but compresses the internal components individually, allowing them to be selectively restored without decompressing the entire archive.

We'll check the Compress option and go, saving out a backup file.

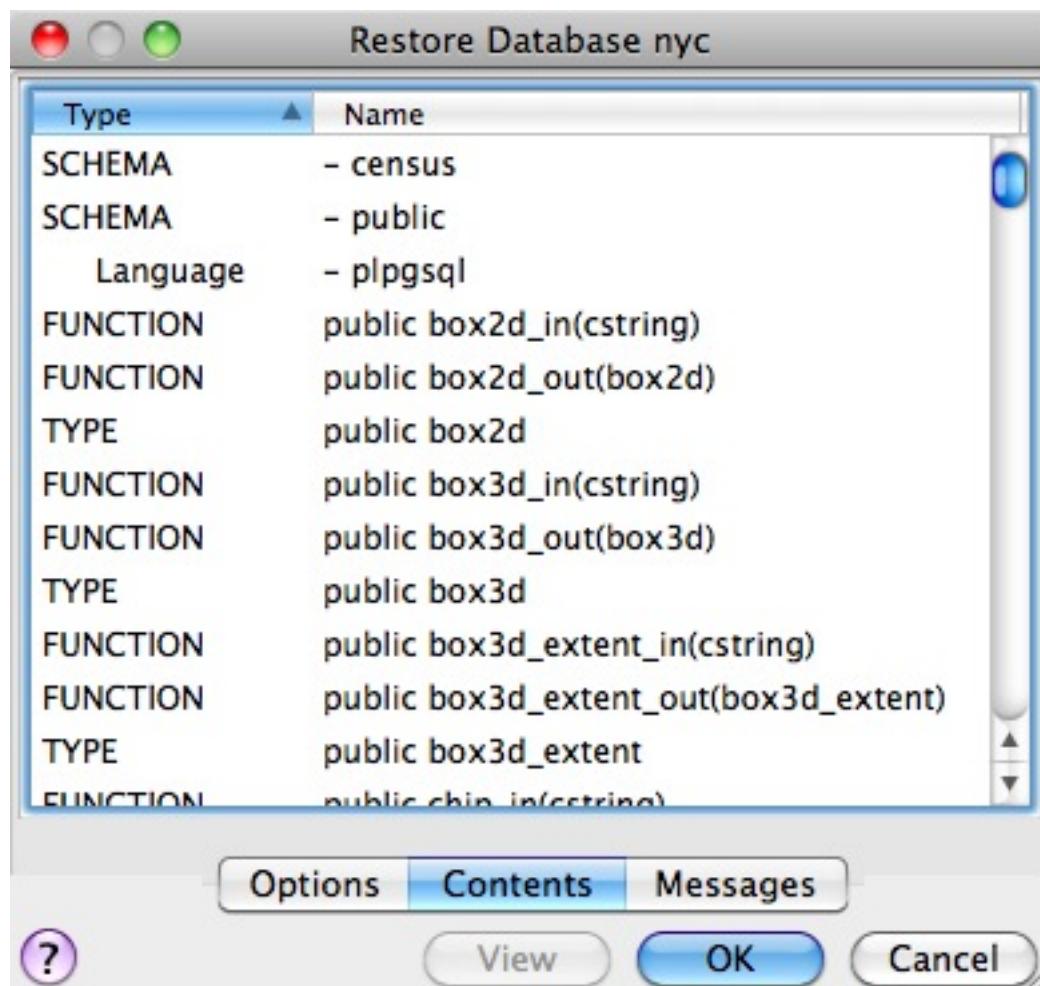
The same operation can be done with the command line like this:

```
pg_dump --file=nyc.backup --format=c --port=54321 --username=postgres nyc
```

Because the backup file is in Compress format, we can view the contents using the `pg_restore` command to list the manifest. In the PgAdmin GUI, “View” is an option in the panel.



When you look at the manifest, one of the things you might notice is that there are a lot of “FUNCTION” signatures in there.



That's because the `pg_dump` utility dumps **every** non-system object in the database, and that includes the PostGIS function definitions.

---

**Note:** PostgreSQL 9.1+ includes an “EXTENSION” feature that allows add-on packages like PostGIS to be installed as registered system components and therefore excluded from `pg_dump` output. PostGIS 2.0 and higher support installation using this extension system.

---

We can see the same manifest from the command-line using `pg_restore` directly:

```
pg_restore --list nyc.backup
```

The problem with a dump file full of PostGIS function signatures is that we really wanted a dump of our data, not our system functions.

Since every object is in the dump file, we can restore to a blank database and get full functionality. In doing so, we are expecting that system we are restoring to has exactly the same version of PostGIS as the one we dumped from (since the function signature definitions reference a particular version of the PostGIS shared library).

From the command-line the restore looks like this:

```
createdb --port 54321 nyc2
pg_restore --dbname=nyc2 --port 54321 --username=postgres nyc.backup
```

Dumping just data, without function signatures, is where having data in schemas is handy, because there is a command-line flag to only dump a particular schema:

```
pg_dump --port=54321 --format=c --schema=census --file=census.backup
```

Now when we list the contents of the dump, we see just the data tables we wanted:

```
pg_restore --list census.backup

;
; Archive created at Thu Aug  9 11:02:49 2012
;     dbname: nyc
;     TOC Entries: 11
;     Compression: -1
;     Dump Version: 1.11-0
;     Format: CUSTOM
;     Integer: 4 bytes
;     Offset: 8 bytes
;     Dumped from database version: 8.4.9
;     Dumped by pg_dump version: 8.4.9
;
;
; Selected TOC Entries:
;
6; 2615 20091 SCHEMA - census postgres
146; 1259 19845 TABLE census nyc_census_blocks postgres
145; 1259 19843 SEQUENCE census nyc_census_blocks_gid_seq postgres
2691; 0 0 SEQUENCE OWNED BY census nyc_census_blocks_gid_seq postgres
2692; 0 0 SEQUENCE SET census nyc_census_blocks_gid_seq postgres
2681; 2604 19848 DEFAULT census gid postgres
2688; 0 19845 TABLE DATA census nyc_census_blocks postgres
2686; 2606 19853 CONSTRAINT census nyc_census_blocks_pkey postgres
2687; 1259 20078 INDEX census nyc_census_blocks_geom_gist postgres
```

Having just the data tables is handy, because it means we can store to a database with any version of PostGIS installed, as we talk about in upgrades.

## Backing Up Users

The `pg_dump` utility operates a database at a time (or a schema or table at a time, if you restrict it). However, information about users is stored across an entire cluster, it's not stored in any one database!

To backup your user information, use the `pg_dumpall` utility, with the “`--globals-only`” flag.

```
pg_dumpall --globals-only --port 54321
```

You can also use `pg_dumpall` in its default mode to backup an entire cluster, but be aware that, as with `pg_dump`, you will end up backing up the PostGIS function signatures, so the dump will have to be restored against an identical software installation, it can't be used as part of an upgrade process.

### 6.3.3 Online Backup and Restore

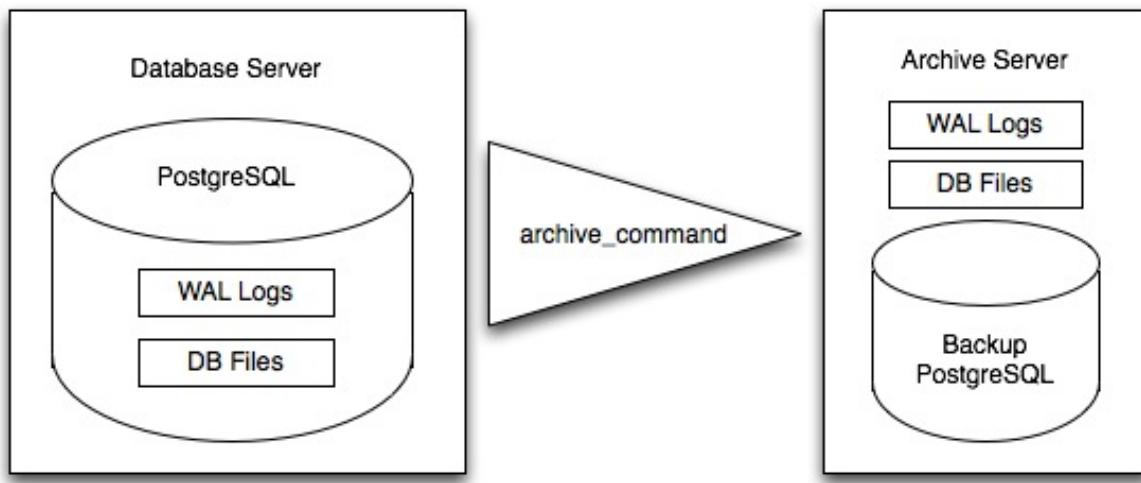
Online backup and restore allows an administrator to keep an extremely up-to-date set of backup files without the overhead of repeatedly dumping the entire database. If the database is under frequent insert and update load, then online backup might be preferable to basic backup.

**Note:** The best way to learn about online backup is to read the relevant sections of the PostgreSQL manual on [continuous archiving and point-in-time recovery](#). This section of the PostGIS workshop will just provide a brief snapshot of online backup set-up.

---

### How it Works

Rather than continually write to the main data tables, PostgreSQL stores changes initially in “write-ahead logs” (WAL). Taken together, these logs are a complete record of all changes made to a database. Online backup consists of taking a copy of the database main data table, then taking a copy of each WAL that is generated from then on.



When it is time to recover to a new database, the system starts on the main data copy, then replays all the WAL files into the database. The end result is a restored database in the same state as the original at the time of the last WAL received.

Because WAL are being written anyways, and transferring copies to an archive server is computationally cheap, online backup is an effective means of keeping a very up-to-date backup of a system without resorting to intensive regular full dumps.

### Archiving the WAL Files

The first thing to do in setting up online backup is to create an archiving method. PostgreSQL archiving methods are the ultimate in flexibility: the PostgreSQL backend simply calls a script specified in the `archive_command` configuration parameter.

That means archiving can be as simple as copying the file to a network-mounted drive, and as complex as encrypting and emailing the files to the remote archive. Any process you can script you can use to archive the files.

To turn on archiving we will edit `postgresql.conf`, first turning on WAL archiving:

```
wal_level = archive  
archive_mode = on
```

And then setting the `archive_command` to copy our archive files to a safe location (changing the destination paths as appropriate):

```
# Unix
archive_command = 'test ! -f /archivedir/%f && cp %p /archivedir/%f'

# Windows
archive_command = 'copy "%p" "C:\\\\archivedir\\\\%f"'
```

It is important that the archive command not over-write existing files, so the unix command includes an initial test to ensure that the files aren't already there. It is also important that the command returns a non-zero status if the copy process fails.

Once the changes are made you can re-start PostgreSQL to make them effective.

## Taking the Base Backup

Once the archiving process is in place, you need to take a base back-up.

Put the database into backup mode (this doesn't do anything to alter operation of queries or data updates, it just forces a checkpoint and writes a label file indicating when the backup was taken).

```
SELECT pg_start_backup('/archivedir/basebackup.tgz');
```

For the label, using the path to the backup file is a good practice, as it helps you track down where the backup was stored.

Copy the database to an archival location:

```
# Unix
tar cvfz /archivedir/basebackup.tgz ${PGDATA}
```

Then tell the database the backup process is complete.

```
SELECT pg_stop_backup();
```

All these steps can of course be scripted for regular base backups.

## Restoring from the Archive

These steps are taking from the PostgreSQL manual on [continuous archiving and point-in-time recovery](#).

- Stop the server, if it's running.
- If you have the space to do so, copy the whole cluster data directory and any tablespaces to a temporary location in case you need them later. Note that this precaution will require that you have enough free space on your system to hold two copies of your existing database. If you do not have enough space, you should at least save the contents of the cluster's pg\_xlog subdirectory, as it might contain logs which were not archived before the system went down.
- Remove all existing files and subdirectories under the cluster data directory and under the root directories of any tablespaces you are using.
- Restore the database files from your file system backup. Be sure that they are restored with the right ownership (the database system user, not root!) and with the right permissions. If you are using tablespaces, you should verify that the symbolic links in pg\_tblspc/ were correctly restored.
- Remove any files present in pg\_xlog/; these came from the file system backup and are therefore probably obsolete rather than current. If you didn't archive pg\_xlog/ at all, then recreate it with

proper permissions, being careful to ensure that you re-establish it as a symbolic link if you had it set up that way before.

- If you have unarchived WAL segment files that you saved in step 2, copy them into pg\_xlog/. (It is best to copy them, not move them, so you still have the unmodified files if a problem occurs and you have to start over.)
- Create a recovery command file recovery.conf in the cluster data directory (see Chapter 26). You might also want to temporarily modify pg\_hba.conf to prevent ordinary users from connecting until you are sure the recovery was successful.
- Start the server. The server will go into recovery mode and proceed to read through the archived WAL files it needs. Should the recovery be terminated because of an external error, the server can simply be restarted and it will continue recovery. Upon completion of the recovery process, the server will rename recovery.conf to recovery.done (to prevent accidentally re-entering recovery mode later) and then commence normal database operations.
- Inspect the contents of the database to ensure you have recovered to the desired state. If not, return to step 1. If all is well, allow your users to connect by restoring pg\_hba.conf to normal.

### 6.3.4 Links

- [pg\\_dump](#)
- [pg\\_dumpall](#)
- [pg\\_restore](#)
- [PostgreSQL High Availability](#)
- [PostgreSQL High Availability Continuous Archiving and PITR](#)

## 6.4 Software Upgrades

Because PostGIS resides within PostgreSQL every PostGIS installation actually consists of two versions of software: the PostgreSQL version and the PostGIS version. The OpenGeo Suite only ships one combination at a time, so the number of potential combinations is reduced, but as a general principle, each version of PostGIS can be theoretically run within a number of versions of PostgreSQL, and vice versa.

So, upgrades need to be considered in terms of upgrading each component.

### 6.4.1 Upgrading PostgreSQL

There are two kinds of PostgreSQL upgrade scenarios:

- A “minor upgrade” when the software version increases at the “patch” level. For example, from 8.4.3 to 8.4.4, or from 9.0.1 to 9.0.3. Increases of more than one patch version are just fine. Minor upgrades fix bugs but do not add any new features or change behaviour.
- A “major upgrade” when the “major” or “minor” versions increase. For example, from 8.4.5 to 9.0.0, or from 9.0.5 to 9.1.1. Major upgrades add new features and change behavior.

#### Minor PostgreSQL Upgrades

For “minor upgrades”, no special process is necessary. Simply install the new software, and re-start the server.

## Major PostgreSQL Upgrades

For “major upgrades” there are two ways to carry out the upgrade.

### Dump/Restore

Dumping and restoring involves converting all the data to a platform neutral format (text representations) on dump, and back to native representations on restore, so it can be time consuming and CPU intensive. However, if you are migrating to a new architecture or operating system, it’s a required process. It’s also a time-tested and well-understood upgrade path, so if your database is not too big, there’s no reason not to stick with it.

- Dump your data `pg_dumpall` from the old database.
- Install the new version of PostgreSQL and the same version of PostGIS you are using in your old database. You need to match the PostGIS version so that the dump file function definitions reference an expected version of the PostGIS library.
- Initialize the new data area using the `initdb` program from the new software.
- Start the new server on the new data area.
- Restore the dump file using `pg_restore`.

### `pg_upgrade`

The `pg_upgrade` utility allows PostgreSQL data directories to be upgraded without the requirement for a dump/restore step. The utility cannot handle changes to the data files themselves, but handles the more common and frequent changes to system tables that occur in PostgreSQL major upgrades.

---

**Note:** The full instructions for running the upgrade process are in the `pg_upgrade` web page at the PostgreSQL site.

---

The `pg_upgrade` program expects to have access to both versions of PostgreSQL it is working with, the old and the new version, so you will have to install them both.

- Install the new version of PostgreSQL you will be using.
- Install the same version of PostGIS you are using in the old PostgreSQL into the new PostgreSQL.
- Initialize the new PostgreSQL data area with the new copy of `initdb`.
- Ensure both the old and new PostgreSQL servers are turned off.
- Run `pg_upgrade`, making sure to use the binary from the new software installation.

```
pg_upgrade
--old-datadir "C:/Program Files/PostgreSQL/8.4/data"
--new-datadir "C:/Program Files/PostgreSQL/9.0/data"
--old-bindir "C:/Program Files/PostgreSQL/8.4/bin"
--new-bindir "C:/Program Files/PostgreSQL/9.0/bin"
```

- If `pg_upgrade` generated any `.sql` files, run them now.
- Start the new server.

### 6.4.2 Upgrading PostGIS

There are two upgrade scenarios for PostGIS too, but they are slightly different from the PostgreSQL scheme:

There are two kinds of PostGIS upgrade scenarios:

- A “minor or patch upgrade” is when the software version increases at the “patch” or “minor” level. For example, from 2.0.1 to 2.0.2, or from 2.0.2 to 2.1.0. Patch upgrades fix bugs only and do not add new features. Minor upgrades fix add new features or change behaviour.
- A “major upgrade” is when the “major” version increases. This is extremely rare. For example, from 0.9.4 to 1.0.0, or from 1.5.4 to 2.0.0. Major upgrades change the on-disk storage format for geometries and require a full database dump and restore.

#### Minor/Patch PostGIS Upgrades

PostGIS deals with minor and upgrades through the EXTENSION mechanism. If you spatially-enabled your database using `CREATE EXTENSION postgis`, you can update your database using the same functionality.

First, install the new software so it is available to the database.

Then, run the SQL to upgrade your PostGIS extension.

```
-- To upgrade to 2.1.2
ALTER EXTENSION postgis UPDATE TO '2.1.2';
```

#### Major PostGIS Upgrades

Major upgrades involve changes to the actual data format for the on-disk storage of geometry and geography data. As such, the data tables need to be re-written. The only way to achieve this is to dump (creating a neutral text-based output) and restore (writing the new table format to disk).

To upgrade, you will have to dump your data first, as discussed in *PostgreSQL Backup and Restore*.

#### With Data in Schemas

- Dump your data by schema.

```
pg_dump
--port=54321
--type=compressed
--file=yourschema.backup
--schema=yourschema
yourdatabase
```

- Install the new version of the PostGIS software.
- Create a new blank database, and enable PostGIS in it.
- Load your data using pg\_restore.

```
pg_restore
--port=54321
--type=compressed
--dbname=yournewdatabase
yourschema.backup
```

## Without Data in Schemas

In this case you have to dump the whole database, which means the dump file will contain PostGIS function and type signatures, and old ones at that. Before loading that file back into the new database, we strip out all the PostGIS-specific bits using a magic script from the PostGIS distribution.

- Dump your whole database, using the “compressed” backup format.

```
pg_dump
--port=54321
--type=compressed
--file=yourdatabase.backup yourdatabase
```

- Install the new version of the PostGIS software.
- Filter your database backup using the ./utils/postgis\_restore.pl script from the new version of PostGIS.

```
postgis_restore.pl yourdatabase.backup > yourdatabase.sql
```

- Create a new blank database, and enable PostGIS in it.

```
-- New in PostGIS 2+ / PostgreSQL 9.1+
-- Formal extensions replace hand loading sql files!
CREATE EXTENSION postgis;
```

- Load the filtered data back into the new database

```
psql
--port=54321
--file=yourdatabase.sql
--dbname=yournewdatabase
```

You should now have an upgraded database ready to use.

### 6.4.3 PostGIS 2.0 Upgrade Issues

In addition to being a major upgrade, and therefore requiring a dump and restore, PostGIS 2.0 made some major changes to behaviour and functionality, some of which are user facing. If you are upgrading from 1.X to 2.X you'll want to remain aware of these.

#### Function Signatures

Not only are there new functions, to support new features, but PostGIS 2.0 **removed a large number of old function signatures**. In particular, most of the function variants that are not prefixed by “ST\_” have been removed.

For example, ST\_Intersects() exists, but Intersects() no longer exists in PostGIS 2.0.

For some client applications, upgrading the software to no longer use the old signatures is not an option. For those users, **it is possible to restore the old signatures**, by loading the legacy.sql file into your database.

Loading legacy.sql reestablishes all the old function signatures as aliases to the new signatures.

### Default WKT and WKB

Prior to PostGIS 2.0, the default forms of the ST\_AsBinary() and ST\_AsText() functions were the *OGC SFSQL* defined versions, which only supported two dimensions. Running ST\_AsBinary() and ST\_AsText() on 3-D and 4-D features just caused the extra dimensions to be stripped, and the returns were 2-D.

For PostGIS 2.0, the ISO SQL/MM definitions of ST\_AsBinary() and ST\_AsText() are used. For 2-D features, the representations are the same, so no changes will be noticed. For 3-D and 4-D features, however, legal representations in ISO SQL/MM exist, so the dimensions will no longer be stripped, and ISO text and binary results will be returned.

For well-known text, that means that the type string will include dimensionality information, and there will be extra ordinates, eg:

```
POINT Z (0 0 0)
LINESTRING ZM (0 0 0 0, 1 1 1 1)
```

For well-known binaries, that means that the type number will be promoted by a multiple of 1000 to indicate the dimensionality.

- 0 implies 2D
- 1000 implied 3D with a Z
- 2000 implies 3D with an M
- 3000 implies 4D

So, for example

- A POINT has type number 1, a POINT ZM has type number 3001
- A LINESTRING has type number 2, LINESTRING M has type number 2002.

### Default SRID

For PostGIS 0.X and 1.X, the SRID assigned to geometries created without specifying an SRID was -1.

For PostGIS 2.X, the SRID assigned to geometries created without specifying an SRID is 0.

This is only important to client applications calling the ST\_SRID() function and testing the result.

### SRID Range Limits

In order to fit the SRID number into a limited address range in the PostgreSQL system tables, the range of values PostGIS 2.X supports for SRID numbers is actually smaller than the range supported in 1.X.

Legal user-defined SRIDs in PostGIS 2.X are from 1 to 998999. The top 10000 SRIDs are retained by PostGIS for internal use.

## APPENDIX A: POSTGIS FUNCTIONS

### 7.1 Constructors

**ST\_MakePoint (Longitude, Latitude)** Returns a new point. Note the order of the coordinates (longitude then latitude).

**ST\_GeomFromText (WellKnownText, srid)** Returns a new geometry from a standard WKT string and srid.

**ST\_SetSRID (geometry, srid)** Updates the srid on a geometry. Returns the same geometry. This does not alter the coordinates of the geometry, it just updates the srid. This function is useful for conditioning geometries created without an srid.

**ST\_Expand (geometry, Radius)** Returns a new geometry that is an expanded bounding box of the input geometry. This function is useful for creating envelopes for use in indexed searches.

### 7.2 Outputs

**ST\_AsText (geometry)** Returns a geometry in a human-readable text format.

**ST\_AsGML (geometry)** Returns a geometry in standard OGC *GML* format.

**ST\_AsGeoJSON (geometry)** Returns a geometry to a standard *GeoJSON* format.

### 7.3 Measurements

**ST\_Area (geometry)** Returns the area of the geometry in the units of the spatial reference system.

**ST\_Length (geometry)** Returns the length of the geometry in the units of the spatial reference system.

**ST\_Perimeter (geometry)** Returns the perimeter of the geometry in the units of the spatial reference system.

**ST\_NumPoints (linestring)** Returns the number of vertices in a linestring.

**ST\_NumRings (polygon)** Returns the number of rings in a polygon.

**ST\_NumGeometries (geometry)** Returns the number of geometries in a geometry collection.

### 7.4 Relationships

**ST\_Distance (geometry, geometry)** Returns the distance between two geometries in the units of the spatial reference system.

**ST\_DWithin(geometry, geometry, radius)** Returns true if the geometries are within the radius distance of one another, otherwise false.

**ST\_Intersects(geometry, geometry)** Returns true if the geometries are not disjoint, otherwise false.

**ST\_Contains(geometry, geometry)** Returns true if the first geometry fully contains the second geometry, otherwise false.

**ST\_Crosses(geometry, geometry)** Returns true if a line or polygon boundary crosses another line or polygon boundary, otherwise false.

## APPENDIX B: GLOSSARY

**CRS** A “coordinate reference system”. The combination of a geographic coordinate system and a projected coordinate system.

**GDAL** [Geospatial Data Abstraction Library](#), pronounced “GOO-duhl”, an open source raster access library with support for a large number of formats, used widely in both open source and proprietary software.

**GeoJSON** “Javascript Object Notation”, a text format that is very fast to parse in Javascript virtual machines. In spatial, the extended specification for [GeoJSON](#) is commonly used.

**GIS** [Geographic information system](#) or geographical information system captures, stores, analyzes, manages, and presents data that is linked to location.

**GML** [Geography Markup Language](#). GML is the [OGC](#) standard XML format for representing spatial feature information.

**JSON** “[Javascript Object Notation](#)”, a text format that is very fast to parse in Javascript virtual machines. In spatial, the extended specification for [GeoJSON](#) is commonly used.

**JSTL** “JavaServer Page Template Library”, a tag library for [JSP](#) that encapsulates many of the standard functions handled in JSP (database queries, iteration, conditionals) into a terse syntax.

**JSP** “JavaServer Pages” a scripting system for Java server applications that allows the interleaving of markup and Java procedural code.

**KML** “Keyhole Markup Language”, the spatial XML format used by Google Earth. Google Earth was originally written by a company named “Keyhole”, hence the (now obscure) reference in the name.

**OGC** The [Open Geospatial Consortium](#) (OGC) is a standards organization that develops specifications for geospatial services.

**OSGeo** The [Open Source Geospatial Foundation](#) (OSGeo) is a non-profit foundation dedicated to the promotion and support of open source geospatial software.

**SFSQL** The [Simple Features for SQL](#) (SFSQL) specification from the [OGC](#) defines the types and functions that make up a standard spatial database.

**SLD** The [Styled Layer Descriptor](#) (SLD) specification from the [OGC](#) defines a format for describing cartographic rendering of vector features.

**SRID** “Spatial reference ID” a unique number assigned to a particular “coordinate reference system”. The PostGIS table `spatial_ref_sys` contains a large collection of well-known srid values and text representations of the coordinate reference systems.

**SQL** “[Structured query language](#)” is the standard means for querying relational databases.

**SQL/MM** SQL Multimedia; includes several sections on extended types, including a substantial section on spatial types.

**SVG** “Scalable vector graphics” is a family of specifications of an XML-based file format for describing two-dimensional vector graphics, both static and dynamic (i.e. interactive or animated).

**WFS** The Web Feature Service (WFS) specification from the [OGC](#) defines an interface for reading and writing geographic features across the web.

**WMS** The Web Map Service (WMS) specification from the [OGC](#) defines an interface for requesting rendered map images across the web.

**WKB** “Well-known binary”. Refers to the binary representation of geometries described in the Simple Features for SQL specification ([SFSQL](#)).

**WKT** “Well-known text”. Can refer either to the text representation of geometries, with strings starting “POINT”, “LINESTRING”, “POLYGON”, etc. Or can refer to the text representation of a [CRS](#), with strings starting “PROJCS”, “GEOGCS”, etc. Well-known text representations are [OGC](#) standards, but do not have their own specification documents. The first descriptions of WKT (for geometries and for CRS) appeared in the [SFSQL](#) 1.0 specification.

---

**CHAPTER  
NINE**

---

**APPENDIX C: LICENSE**

This work is licensed under the Creative Commons Attribution-Share Alike, United States License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

Our attribution requirement is that you retain the visible copyright notices in all materials.



## INDEX

### C

CRS, **179**

### G

GDAL, **179**

GeoJSON, **179**

GIS, **179**

GML, **179**

### J

JSON, **179**

JSP, **179**

JSTL, **179**

### K

KML, **179**

### O

OGC, **179**

OSGeo, **179**

### S

SFSQL, **179**

SLD, **179**

SQL, **179**

SQL/MM, **180**

SRID, **179**

SVG, **180**

### W

WFS, **180**

WKB, **180**

WKT, **180**

WMS, **180**