# Spatial Cypher Cheat Sheet - Intro To Geospatial Cypher Functions With Neo4j
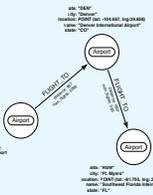
## Intro To Cypher & The Property Graph Data Model

The variable **p** is bound to the path connecting SFO and RSW

Nodes are defined within parentheses **()** properties within curly braces **{}**

The variable length path operator * defines paths of arbitrary length. We specify a maximum length of 2 relationships.

```
MATCH p=(sfo:Airport {iata: "SFO"})-[b:FLIGHT_TO*2]->(rsw:Airport {iata: "RSW"}) RETURN *
```

Search the graph for the following pattern.

This graph pattern represents a node with the label Airport that has a property iata with the value "SFO", connected through one or more outgoing relationships of type FLIGHT_TO to a node with the label Airport and iata value "RSW" . We bind variables sfo, f, and rsw to refer to each part of the graph pattern later in the query if needed.

Return any variables that we've matched on in previous parts of the query

**Neo4j** is a database management system (DBMS) that uses the **property graph data model** which is composed of nodes, relationships, and properties to model, store, and query data as a graph. Nodes can have one or more labels, relationships have a single type and direction. Key-value pair properties can be stored on nodes and relationships.

The **Cypher query language** is used to query data and interact with Neo4j. Cypher is a declarative query language that uses ASCII-art like syntax to define graph patterns that form the basis of most query operations. Nodes are defined with parenthesis, relationships with square brackets, and can be combined to create complex graph patterns. Common Cypher commands are **MATCH** (find where the graph pattern exists), **CREATE** (add data to the database using the specified graph pattern), and **RETURN** (return a subset of the data as a result of a traversal through the graph.

## Spatial Point Type

Neo4j supports 2D or 3D geographic (WGS84) or cartesian coordinate reference system (CRS) points

```
RETURN point( {latitude:37.62245, longitude:-122.383989} )
```
Creating a point by specifying latitude/longitude. WGS84 is assumed when using lat/lon.

Point data can be stored as properties on nodes or relationships. Here we create an Airport node and set its location as a point.

```
CREATE (a:Airport)
SET a.iata = "SFO",
a.location = point( {latitude:37.62245, longitude:-122.383989})
RETURN a
```

```
CREATE POINT INDEX airportIndex
FOR (a:Airport) ON (a.location)
```
Database indexes are used to speed up search performance. Here we create a database index on the location property for Airport nodes. This will help us find airports faster when searching for airports by location (radius distance or bounding box spatial search).

## Spatial Cypher Functions

### Radius Distance

```
MATCH (a:Airport)
WHERE point.distance(
a.location,
point({latitude:37.55948, longitude:-122.32544})) < 20000
RETURN a
```
To find nodes close to other nodes in the graph we can use the **point.distance()** function to perform a radius distance search

### Within Bounding Box

```
MATCH (a:Airport)
WHERE point.withinBBox(
a.location,
point({longitude:-122.325447, latitude: 37.55948 }),
point({longitude:-122.314675 , latitude: 37.563596}))
RETURN a
```
To search for nodes within a bounding box we can use the **point.withinBBox()** function.

### Geocoding

```
CALL apoc.spatial.geocode('SFO Airport') YIELD location
---------------------------------------------------------
{
    "description": "San Francisco International Airport, 780,
South Airport Boulevard, South San Francisco, San Mateo County,
CAL Fire Northern Region, California, 94128, United States",
    "longitude": -122.38398938548363,
    "latitude": 37.622451999999996,
}
```
To geocode a location description into latitude, longitude location we can use the **apoc.spatial.geocode()** procedure. By default this procedure uses the Nominatim geocoding API but can be configured to use other geocoding services, such as Google Cloud.

## Data Import

*We can use Cypher to import data into Neo4j from formats such as CSV and JSON, including GeoJSON.*

**CSV** Using the LOAD CSV Cypher command to create an airport routing graph.

**1 -** Create a constraint on the field that identies uniqueness, in this case Airport IATA code. This ensures we won't create duplicate airports but also creates a database index to improve performance of our data import steps below.

```
CREATE CONSTRAINT FOR (a:Airport) REQUIRE a.iata IS UNIQUE;
```

**2 -** Create Airport nodes, storing their location, name, IATA code, etc as node properties.

```
LOAD CSV WITH HEADERS
FROM "https://cdn.neo4jlabs.com/data/flights/airports.csv"
AS row
MERGE (a:Airport {iata: row.IATA_CODE})
ON CREATE SET a.city = row.CITY,
              a.name = row.AIRPORT,
              a.state = row.STATE,
              a.country = row.country,
              a.location =
                  point({ latitude: toFloat(row.LATITUDE),
                          longitude: toFloat(row.LONGITUDE)
                  });
```

**3 -** Create FLIGHT_TO relationships connecting airports with a connecting flight. Increment the num_flights counter variable to keep track of the number of flights between airports per year.

```
LOAD CSV WITH HEADERS
FROM "https://cdn.neo4jlabs.com/data/flights/flights.csv" AS row
CALL {
    WITH row
    MATCH (origin:Airport {iata: row.ORIGIN_AIRPORT})
    MATCH (dest:Airport {iata: row.DESTINATION_AIRPORT})
    MERGE (origin)-[f:FLIGHT_TO]->(dest)
      ON CREATE SET
        f.num_flights = 0, f.distance = toInteger(row.DISTANCE)
      ON MATCH SET
        f.num_flights = f.num_flights + 1
} IN TRANSACTIONS OF 100000 ROWS;
```

**GeoJSON** We can also store arrays of Points to represent complex geometries like lines and polygons, for example to represent land parcels.

```
CALL apoc.load.json('https://cdn.neo4jlabs.com/data/landgraph/parcels.geojson')
YIELD value
UNWIND value.features AS feature
CREATE (p:Parcel) SET
  p.coordinates = [coord IN feature.geometry.coordinates[0] | point({latitude: coord[1], longitude: coord[0]})]
  p += feature.properties;
```

## Routing With Path Finding Algorithms

### Shortest Path

The shortestPath function performs a **binary breadth-first search** to find the shortest path between nodes in the graph.

```
MATCH p = shortestPath(
    (:Airport {iata: "SFO"})-[:FLIGHT_TO*..10]->(:Airport {iata: "RSW"})
) RETURN p
```

### Shortest Weighted Path

#### Dijkstra's Algorithm

Often we want to consider the shortest weighted path taking into account distance, time or some other cost stored as relationship properties. Dijkstra and A* are two algorithms that take relationship (or edge) weights into account when calculating the shortest path.

```
MATCH (origin:Airport {iata: "SFO"})
MATCH (dest:Airport {iata: "RSW"})
CALL
  apoc.algo.dijkstra(
    origin,
    dest,
    "FLIGHT_TO",
    "distance"
  )
YIELD path, weight
UNWIND nodes(path) AS n
RETURN {
  airport: n.iata,
  lat: n.location.latitude,
  lng: n.location.longitude
} AS route
```

Dijkstra's algorithm is similar to a breadth-first search, but takes into account relationship properties (distance) and prioritizes exploring low-cost routes first using a priority queue.

#### A* Algorithm

```
MATCH (origin:Airport {iata: "SFO"})
MATCH (dest:Airport {iata: "RSW"})
CALL
  apoc.algo.aStarConfig(
    origin,
    dest,
    "FLIGHT_TO",
    {
      pointPropName: "location",
      weight: "distance"
    }
  )
YIELD weight, path
RETURN weight, path
```

The A* algorithm adds a heuristic function to choose which paths to explore. In our case the heuristic is the distance to the final destination.

# Spatial Cypher Cheat Sheet - Using Neo4j With Python

## The Neo4j Python Driver

### Creating A GeoDataFrame From Data Stored In Neo4j

```
import neo4j

NEO4J_URI        = "neo4j://localhost:7689"
NEO4J_USER       = "neo4j"
NEO4J_PASSWORD   = "letmeinnow"
NEO4J_DATABASE   = "neo4j"
```
*Connection credentials for our Neo4j database*

```
driver = neo4j.GraphDatabase.driver(
  NEO4J_URI, auth=(NEO4J_USER, NEO4J_PASSWORD)
)
```
*Create a connection to the database*

*Define a Cypher query to fetch data from Neo4j*

```
AIRPORT_QUERY = """
  MATCH (origin:Airport)-[f:FLIGHT_TO]->(dest:Airport)
  CALL {
    WITH origin
    MATCH (origin)-[f:FLIGHT_TO]-()
    RETURN sum(f.num_flights) AS origin_centrality
  }
  CALL {
    WITH dest
    MATCH (dest)-[f:FLIGHT_TO]-()
    RETURN sum(f.num_flights) AS dest_centrality
  }
  RETURN {
    origin_wkt:
      "POINT (" + origin.location.longitude + " " + origin.location.latitude + ")",
    origin_iata: origin.iata,
    origin_city: origin.city,
    origin_centrality: origin_centrality,
    dest_centrality: dest_centrality,
    dest_wkt:
      "POINT (" + dest.location.longitude + " " + dest.location.latitude + ")",
    dest_iata: dest.iata,
    dest_city: dest.city,
    distance: f.distance,
    num_flights: f.num_flights,
    geometry:
      "LINESTRING (" + origin.location.longitude + " " + origin.location.latitude + ","
      + dest.location.longitude + " " + dest.location.latitude + ")"
  }
```

*Weighted degree centrality is a measure of a node's importance in the network and is the sum of all relationship weights connected to a given node.*

*We return the geometry of our origin and destination airports, and the flight route as Well Known Text (WKT). POINT for the airports and LINESTRING for the flight route. We'll parse this WKT when defining the geometry in our GeoDataFrame.*

```
def get_airport(tx):
    results = tx.run(AIRPORT_QUERY)
    df = results.to_df(expand=True)
    df.columns =
      ['origin_city','origin_wkt', 'dest_city', 'dest_wkt', 'origin_centrality', 'distance', 'origin_iata',
      'geometry','num_flights', 'dest_centrality', 'dest_iata']
    df['geometry'] = geopandas.GeoSeries.from_wkt(df['geometry'])
    df['origin_wkt'] = geopandas.GeoSeries.from_wkt(df['origin_wkt'])
    df['dest_wkt'] = geopandas.GeoSeries.from_wkt(df['dest_wkt'])
    gdf = geopandas.GeoDataFrame(df, geometry='geometry')
    return gdf
```

*The Neo4j Python Driver has a .to_df() function which will convert a Neo4j result set to a pandas DataFrame*

*Here we parse the WKT columns into GeoSeries and convert our pandas DataFrame into a GeoPandas GeoDataFrame*

```
with driver.session(database=NEO4J_DATABASE) as session:
    airport_df = session.execute_read(get_airport)
```
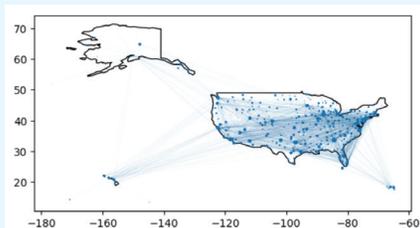
*We now have a GeoDataFrame where each row is a flight route between two airports. We can plot the airport and routes, using the centrality metric to size airport nodes: more important airports should be larger.*

```
world = geopandas.read_file(
  geopandas.datasets.get_path('naturalearth_lowres')
)

ax = world[world.continent == 'North America']
  .plot(color='white', edgecolor='black')

flights_gdf = flights_gdf.set_geometry("origin_wkt")
flights_gdf.plot(ax=ax, markersize='origin_centrality')

flights_gdf = flights_gdf.set_geometry("geometry")
flights_gdf.plot(ax=ax, markersize=0.1, linewidth=0.01)
```

In this section we'll use the Neo4j Python Driver to **create a GeoData-Frame** of our flight data. We'll also compute weighted degree centrality so we can plot airport size relative to their "importance" in the US airline network. The **Neo4j Python Driver** can be installed with:

```
pip install neo4j
```
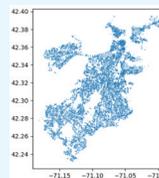
## Working With OpenStreetMap Data

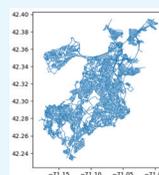### Loading A Road Network With OSMNx

```
pip install osmnx
```

```
import osmnx as ox

G = ox.graph_from_place("Boston, MA, USA", network_type="drive")
fig, ax = ox.plot_graph(G)

gdf_nodes, gdf_relationships = ox.graph_to_gdfs(G)
gdf_nodes.reset_index(inplace=True)
gdf_relationships.reset_index(inplace=True)
```

In this section we will import data from **OpenStreetMap** into Neo4j using the **OSMNx Python package**. Below is the property graph data model we will use to model the road network of Boston.

**gdf_nodes**

*Here is our nodes GeoDataFrame. Each row represents an intersection in the Boston road network.*

**gdf_relationships**

*Here is our relationships GeoDataFrame. Each row represents a road segment connecting two intersec-tions.*

```
node_query = '''
  UNWIND $rows AS row
  WITH row WHERE row.osmid IS NOT NULL
  MERGE (i:Intersection {osmid: row.osmid})
    SET i.location =
      point({latitude: row.y, longitude: row.x }),
      i.ref = row.ref,
      i.highway = row.highway,
      i.street_count = toInteger(row.street_count)
  RETURN COUNT(*) as total
'''
```
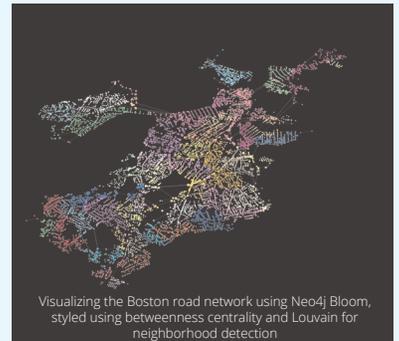*Define a Cypher query to add intersection nodes from the nodes GeoDataFrame*

```
rels_query = '''
  UNWIND $rows AS road
  MATCH (u:Intersection {osmid: road.u})
  MATCH (v:Intersection {osmid: road.v})
  MERGE (u)-[r:ROAD_SEGMENT {osmid: road.osmid}]->(v)
    SET r.oneway = road.oneway,
        r.lanes = road.lanes,
        r.ref = road.ref,
        r.name = road.name,
        r.highway = road.highway,
        r.max_speed = road.maxspeed,
        r.length = toFloat(road.length)
  RETURN COUNT(*) AS total
'''
```
*Adding road segments from the relationships GeoDataFrame connecting intersection nodes*

```
def insert_data(tx, query, rows, batch_size=1000):
    total = 0
    batch = 0

    while batch * batch_size < len(rows):
        results = tx.run(query, parameters = {
          'rows':
            rows[batch * batch_size: (batch + 1) * batch_size]
            .to_dict('records')
        }).data()
        print(results)
        total += results[0]['total']
        batch += 1

with driver.session() as session:
    session.execute_write(insert_data, node_query, gdf_nodes.drop(columns=['geometry']))
    session.execute_write(insert_data, rels_query, gdf_relationships.drop(columns=['geometry']))
```
*Because our GeoDataFrames can be very large we break them up into batches to avoid sending too much data to the database at once.*

Visualizing the Boston road network using Neo4j Bloom, styled using betweenness centrality and Louvain for neighborhood detection