

Intro to Data Engineering

Paul Singman

Data Engineer, Vroom

LEARNING OBJECTIVES

- ▶ Introduce some common DE technologies/concepts
- ▶ Give example of real world data pipeline
- ▶ Perform coding exercise using some best practices

What is Data Engineering?

- ▶ **Definition:** Moving data from Point A to Point B in a **Robust, Fault-Tolerant**, and **Reliable** manner.



Data Science Unicorn



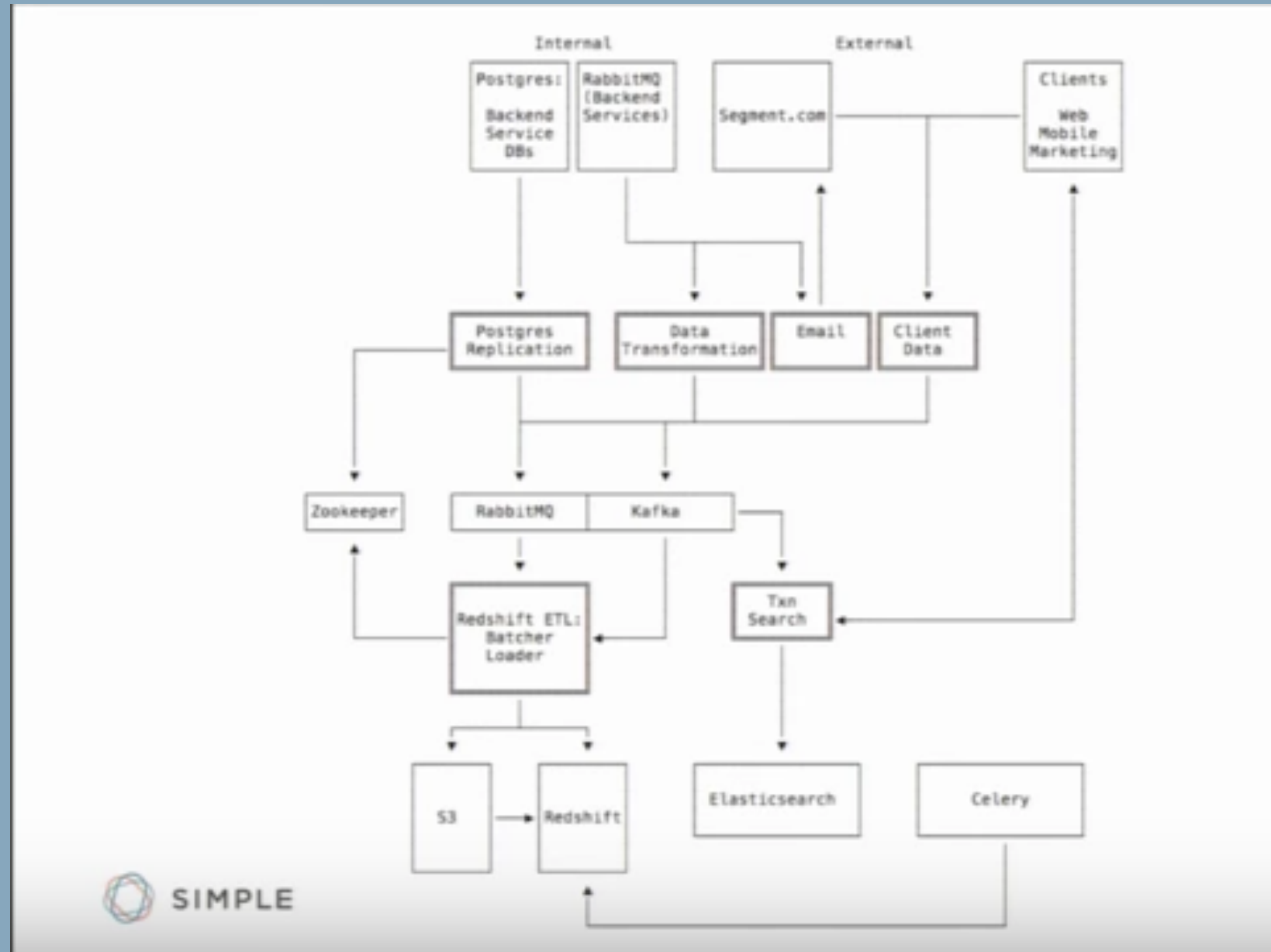
Robustness, Fault-Tolerance & Reliability

- 1) Gracefully handle unexpected inputs / schema changes
- 2) Gracefully handle failing or unresponsive nodes
- 3) Provide logging and metrics such as compute time or No. of bytes processed

Engineering Best Practices

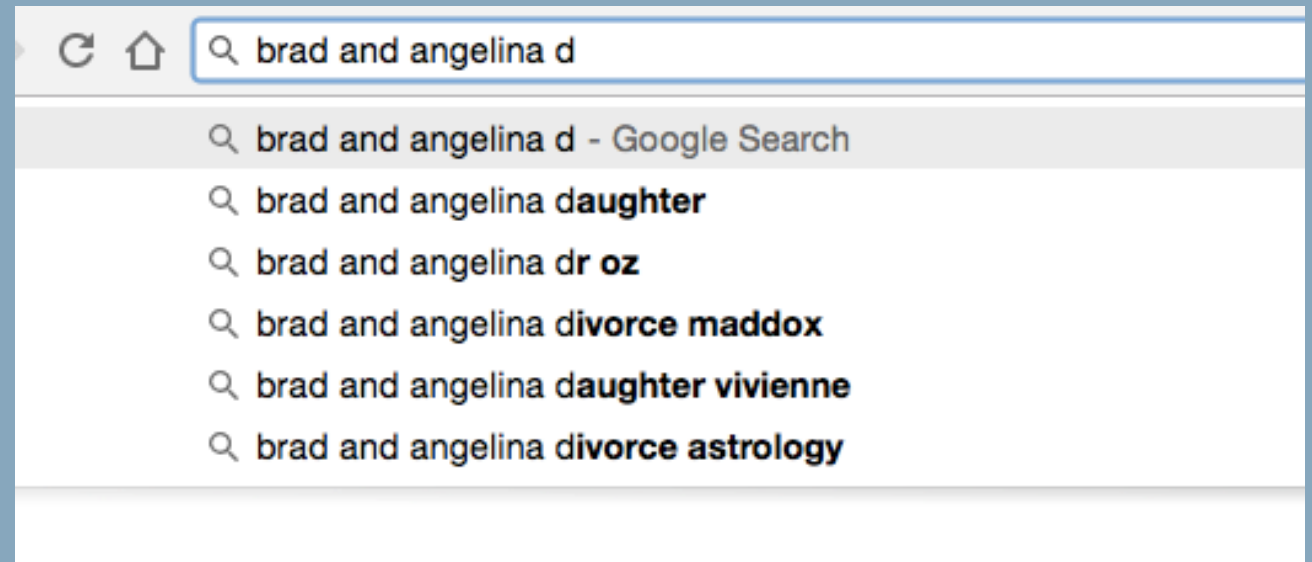
- 1) Operate in isolated development/production environments
- 2) Use version control for code (Git)
- 3) Write unit tests on code and automatically run tests when deploying new code (Continuous Integration)

Source: <https://www.youtube.com/watch?v=9nX35zrN20E>



Batch vs Streaming Pipelines

- 1) Historically companies had batch data pipelines (e.g. job runs once a day at 2:00 AM)
- 2) As technology improves, move towards streaming architectures (real-time processing)



Motivating Example: Google autocomplete suggestions

Hot DE Technologies

- 1) **Spark** (Berkeley) and **Flink**: Distributed Stream Computing
- 2) **Kafka** (LinkedIn) and **Kinesis** (AWS): Distributed Queue
- 3) **Airflow** (Airbnb) and **Luigi** (Spotify): Task Scheduling and Monitoring
- 4) Cloud Computing: Google Cloud Compute (AWS replacement)
- 5) Environment creation via containers: **Docker**

Lab Time!

Storing Spotify Song Data in a
Sqlite DB

Step 1: ???

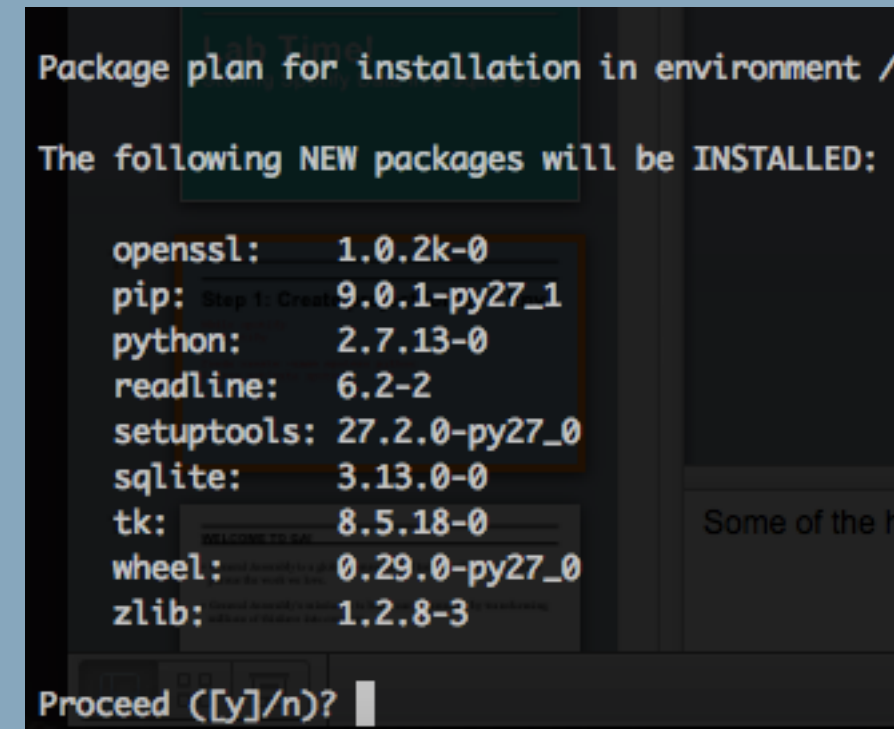
Step 1: Create project folder & Env

`Mkdir spotify`

`Cd spotify`

`Conda create -name spotenv python`

`Source activate spotenv`



```
Package plan for installation in environment /home/.../spotenv

The following NEW packages will be INSTALLED:

openssl: 1.0.2k-0
pip: 9.0.1-py27_1
python: 2.7.13-0
readline: 6.2-2
setuptools: 27.2.0-py27_0
sqlite: 3.13.0-0
tk: 8.5.18-0
wheel: 0.29.0-py27_0
zlib: 1.2.8-3

Proceed ([y]/n)?
```

Step 2: Open Sublime and create python file get_data.py

```
touch get_data.py
```

```
vim get_data.py (to open file and edit in terminal)
```

For fun, type: `import pandas as pd`

Save and run: `python get_data.py`

```
(spotenv) Pauls-MacBook-Air:spotfiy paulsingman$ python get_data.py
Traceback (most recent call last):
  File "get_data.py", line 1, in <module>
    import pandas as pd
ImportError: No module named pandas
```

Step 3: Add to get_data.py

```
get_data.py  create_db.py x
1  import json
2  import requests
3  import sqlite3
4
5
6  def fetch_album_songs(album_id):
7      ...
8      pass
9
10
11 |
12 def store_in_database(songs):
13     ...
14     pass
15
16
17 if __name__ == '__main__':
18     ...
19     songs = fetch_album_songs()
20
21     store_in_database(songs)
22
23
24
```

Step 4: Plan fetch_album_songs

- 1) Build api endpoint url with album_id
- 2) Use requests to get API response
- 3) Convert response into json object
- 4) Parse response for relevant fields
- 5) Return parsed data

<https://open.spotify.com/special/thebeatles>

<https://developer.spotify.com/web-api/get-albums-tracks/>

Step 6: code to fetch_album_songs

- 1) Build custom url with album_id
- 2) Use requests to get API response
- 3) Convert response into json object
- 4) Parse response for relevant fields
- 5) Return parsed data

```
8  def fetch_album_songs(album_id):
9
10     url = URL_BASE + 'albums/' + album_id + '/tracks'  # 1
11
12     response = requests.get(url)  # 2
13
14     json_response = json.loads(response.text)  # 3
15
16     songs = []
17
18     for item in json_response['items']:  # 4
19
20         track_num = item['track_number']
21         name = item['name']
22         duration = item['duration_ms']
23
24         songs.append((track_num, name, duration))
25
26     return songs  # 5
```

```
{u'items': [{u'name': u"Sgt. Pepper's Lonely Hearts Club Band - Remastered", u'external_urls': {u'spotify': u'https://open.spotify.com/track/4fUKE8EULjQdHF
```

A look At the API Response

```
(spotenv) Pauls-MacBook-Air:spotfiy paulsingman$ python get_data.py
{'u'name': u"Sgt. Pepper's Lonely Hearts Club Band - Remastered", u'external_urls': {'u'spotify': u'https://open.spotify.com/track/4fUKE8EULjQdHF4zb0M8F0'}, u'uri': u'spotify:track:4fUKE8EULjQdHF4zb0M8F0', u'explicit': False, u'preview_url': u'https://p.scdn.co/mp3-preview/7ae81e104c9b55dfd0c203678d29a264801711c6?cid=null', u'track_number': 1, u'disc_number': 1, u'href': u'https://api.spotify.com/v1/tracks/4fUKE8EULjQdHF4zb0M8F0', u'artists': [{u'name': u'The Beatles', u'external_urls': {'u'spotify': u'https://open.spotify.com/artist/3WrFJ7ztbogyGnTHbHJfL2'}, u'uri': u'spotify:artist:3WrFJ7ztbogyGnTHbHJfL2', u'href': u'https://api.spotify.com/v1/artists/3WrFJ7ztbogyGnTHbHJfL2', u'type': u'artist', u'id': u'3WrFJ7ztbogyGnTHbHJfL2'}], u'duration_ms': 122893, u'type': u'track', u'id': u'4fUKE8EULjQdHF4zb0M8F0', u'available_markets': [u'CA', u'MX', u'US']}
```

```
{u'name': u'With A Little Help From My Friends - Remastered', u'external_urls': {'u'spotify': u'https://open.spotify.com/track/2RnPATK99oG0ZygnD2GT06'}, u'uri': u'spotify:track:2RnPATK99oG0ZygnD2GT06', u'explicit': False, u'preview_url': u'https://p.scdn.co/mp3-preview/2574cc919f6f70013598de262fc8c3d39b55fbbc?cid=null', u'track_number': 2, u'disc_number': 1, u'href': u'https://api.spotify.com/v1/tracks/2RnPATK99oG0ZygnD2GT06', u'artists': [{u'name': u'The Beatles', u'external_urls': {'u'spotify': u'https://open.spotify.com/artist/3WrFJ7ztbogyGnTHbHJfL2'}, u'uri': u'spotify:artist:3WrFJ7ztbogyGnTHbHJfL2', u'href': u'https://api.spotify.com/v1/artists/3WrFJ7ztbogyGnTHbHJfL2', u'type': u'artist', u'id': u'3WrFJ7ztbogyGnTHbHJfL2'}], u'duration_ms': 164106, u'type': u'track', u'id': u'2RnPATK99oG0ZygnD2GT06', u'available_markets': [u'CA', u'MX', u'US']}
```

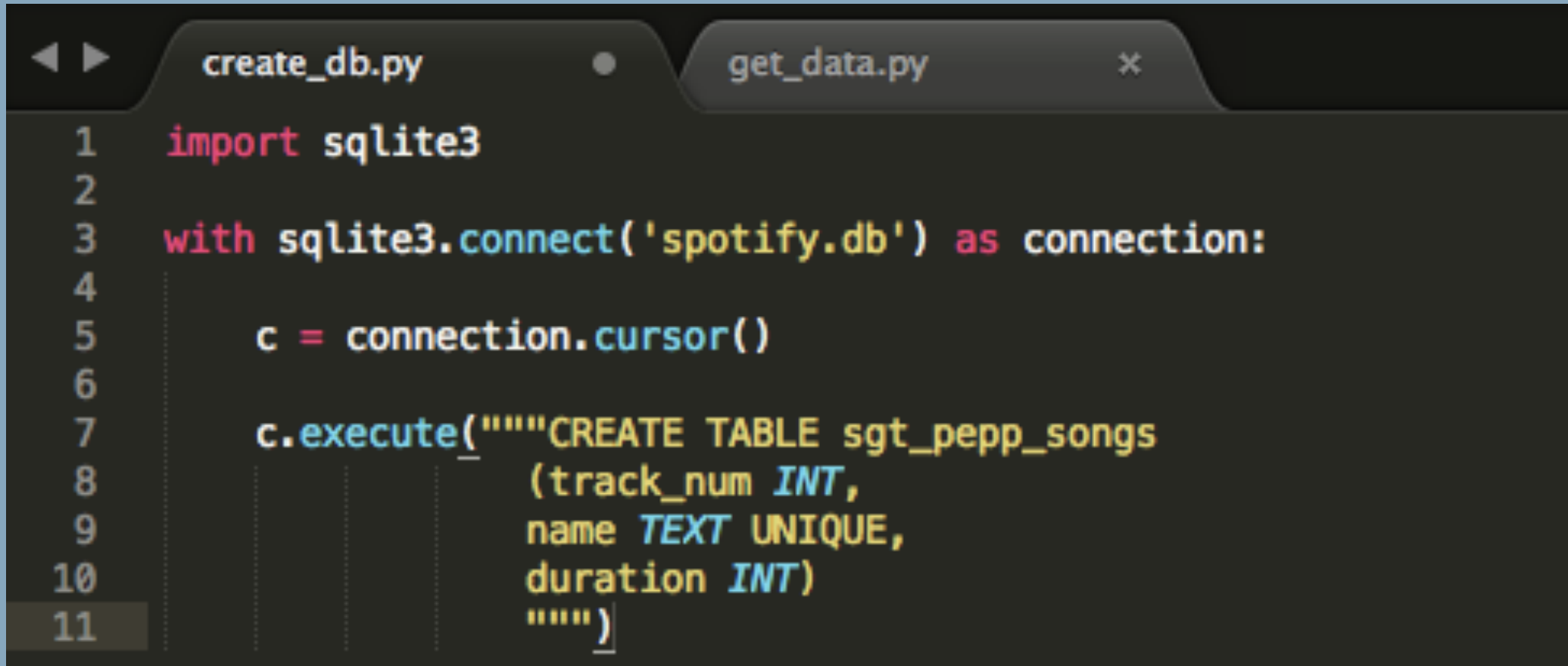
Step 7: code to store_in_database

```
30 def store_in_database(songs):
31
32     connection = sqlite3.connect(DATABASE_NAME)
33
34     cursor = connection.cursor()
35
36     cursor.executemany("""INSERT INTO sgt_pepp_songs VALUES(?,?,?)""", songs)
37
38     connection.commit()
39
40     connection.close()
41
```

Full Script

```
get_data.py x create_db.py x
1 import json
2 import requests
3 import sqlite3
4
5 URL_BASE = 'https://api.spotify.com/v1/'
6 DATABASE_NAME = 'spotify.db'
7
8 def fetch_album_songs(album_id):
9
10     url = URL_BASE + 'albums/' + album_id + '/tracks' # 1
11
12     response = requests.get(url) # 2
13
14     json_response = json.loads(response.text) # 3
15
16     songs = []
17
18     for item in json_response['items']: # 4
19
20         track_num = item['track_number']
21         name = item['name']
22         duration = item['duration_ms']
23
24         songs.append((track_num, name, duration))
25
26     return songs # 5
27
28
29 def store_in_database(songs):
30
31     connection = sqlite3.connect(DATABASE_NAME)
32
33     cursor = connection.cursor()
34
35     cursor.executemany("""INSERT INTO sgt_pepp_songs VALUES(?,?,?)""", songs)
36
37     connection.commit()
38
39     connection.close()
40
41
42 if __name__ == '__main__':
43
44     songs = fetch_album_songs('6QaVfG1pHYl1z15ZxkvVDW')
45
46     store_in_database(songs)
47
```

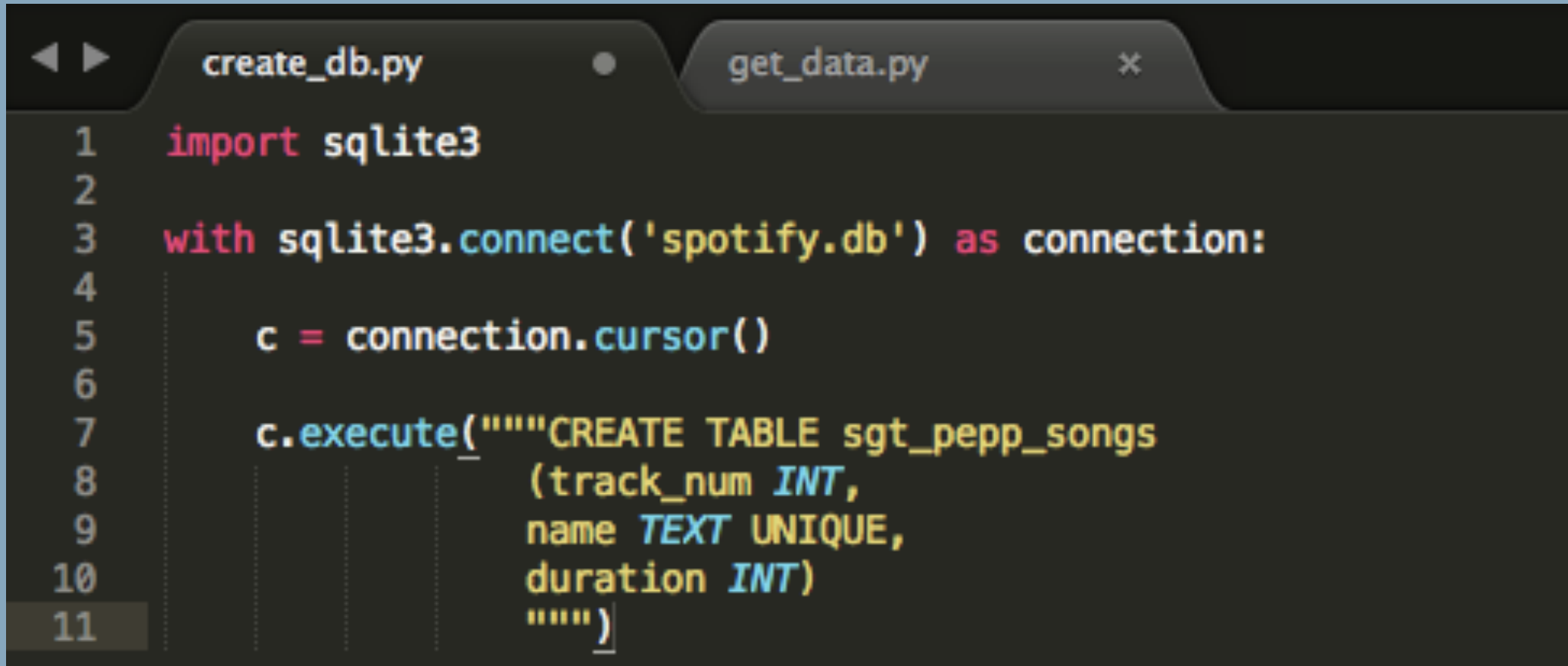
Step 8: Create Database in create_db.py

A screenshot of a code editor with two tabs: 'create_db.py' and 'get_data.py'. The 'create_db.py' tab is active, showing Python code that uses the sqlite3 module to create a database named 'spotify.db' and a table named 'sgt_pepp_songs'. The code includes imports, a with-statement for the database connection, a cursor object, and an execute statement with a multi-line SQL CREATE TABLE command. Line numbers 1 through 11 are visible on the left side of the editor.

```
1 import sqlite3
2
3 with sqlite3.connect('spotify.db') as connection:
4
5     c = connection.cursor()
6
7     c.execute("""CREATE TABLE sgt_pepp_songs
8                 (track_num INT,
9                  name TEXT UNIQUE,
10                 duration INT)
11                """)
```

Run: `python create_db.py`

Step 8: Create Database in create_db.py

A screenshot of a code editor with two tabs: 'create_db.py' and 'get_data.py'. The 'create_db.py' tab is active, showing Python code that uses the sqlite3 module to create a database named 'spotify.db' and a table named 'sgt_pepp_songs'. The code includes imports, a with-statement for the database connection, a cursor object, and an execute statement with a multi-line SQL string. Line numbers 1 through 11 are visible on the left side of the editor.

```
1 import sqlite3
2
3 with sqlite3.connect('spotify.db') as connection:
4
5     c = connection.cursor()
6
7     c.execute("""CREATE TABLE sgt_pepp_songs
8                 (track_num INT,
9                  name TEXT UNIQUE,
10                 duration INT)
11                """)
```

Run: `python create_db.py`

Step 9: Query Away!

Q: Can you write a query that returns the name and duration of the second longest song?

Step 10 (Optional): Adding Command Line Args

```
43 if __name__ == '__main__':
44
45     parser = argparse.ArgumentParser(description='Specify an album id')
46     parser.add_argument('-a', '--album_id', required=True, help='Spotify Album id')
47     args = parser.parse_args()
48
49     songs = fetch_album_songs(args.album_id)
50
```

Run: `python create_db.py -a 6QaVfG1pHYI1z15ZxkvVDW`

Where to go from here

- 1) Check out Spotipy on Github for a generalized python wrapper for Spotify API**
- 2) Use credentials to access all API endpoints**
- 3) Add documentation and logging to the code**
- 4) Get more Spotify data into a database!**

WELCOME TO DATA Engineering

EXIT TICKET

**DON'T FORGET TO FILL OUT YOUR EXIT
TICKET**