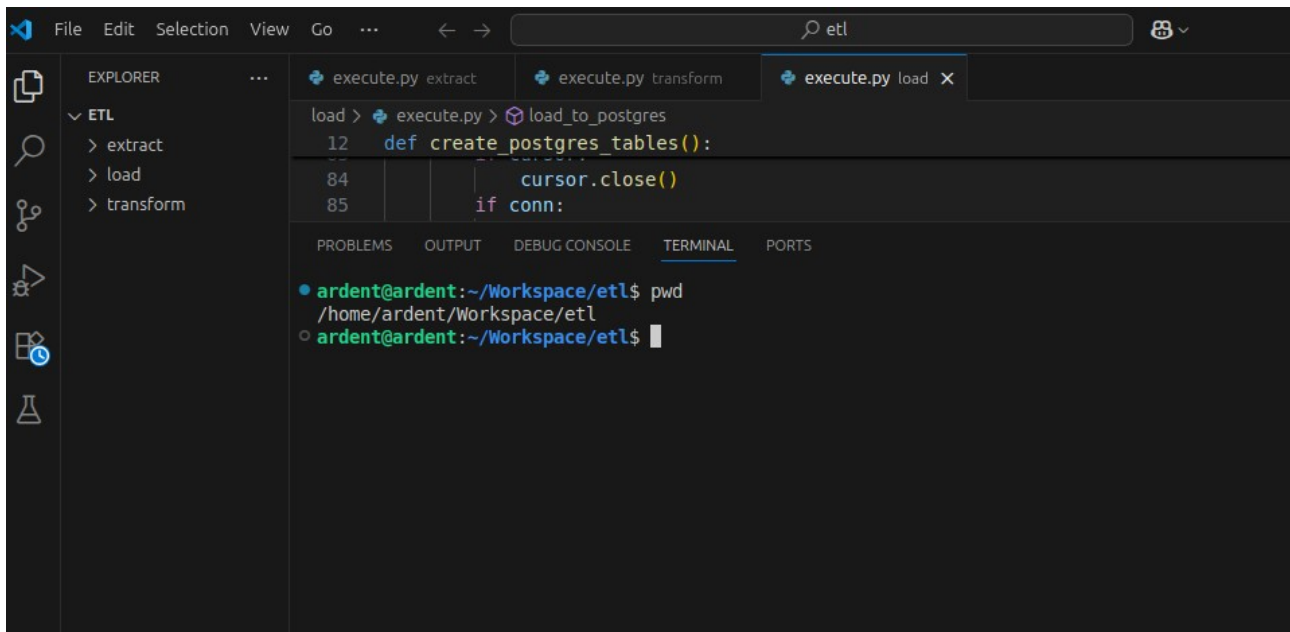
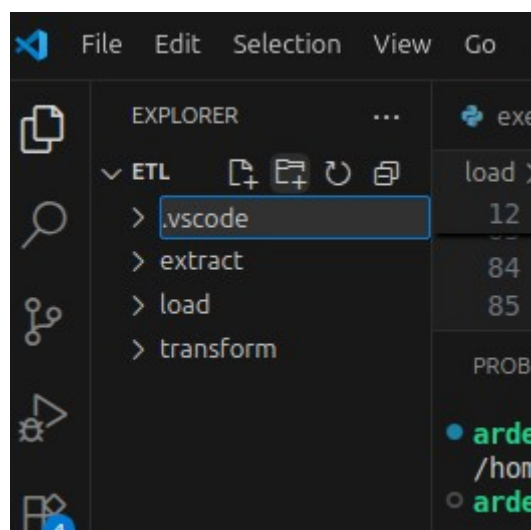


Let us clean our workspace/etl project a little bit. There are several issues with the current code, we will not dive too deep into those. Also, note that this is not the one and only way to do things. Everyone should figure out what works for them and stick with that.

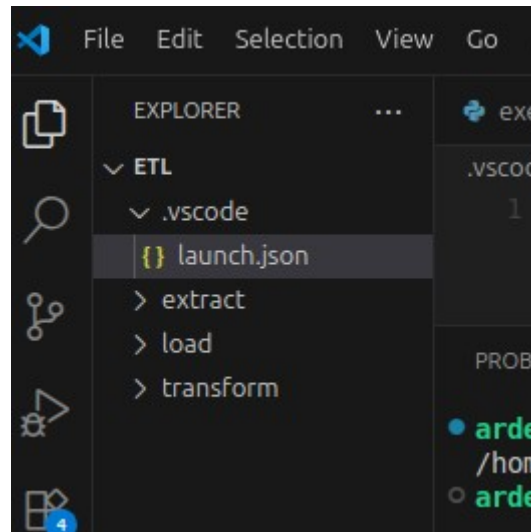
1. Open your ETL project in vscode. It should have three different folders extract,load,transform and an individual execute.py file inside all of it.



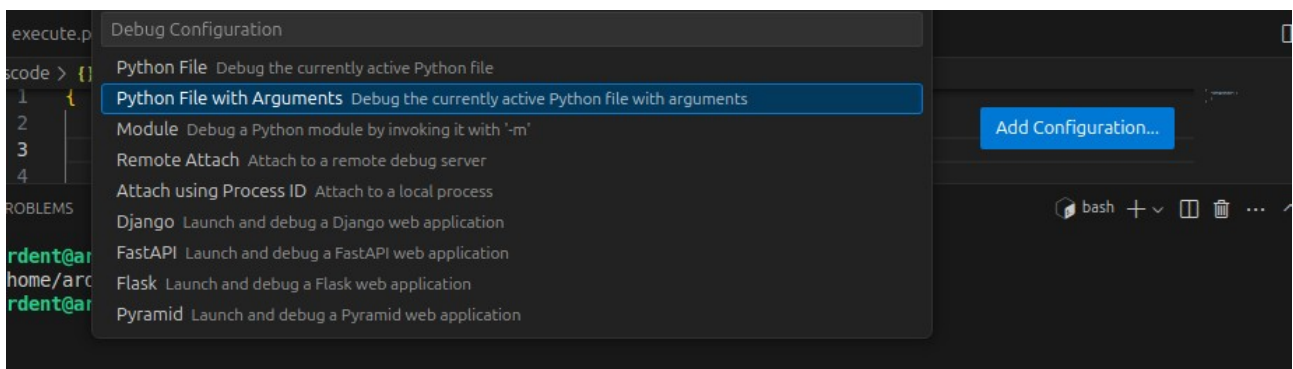
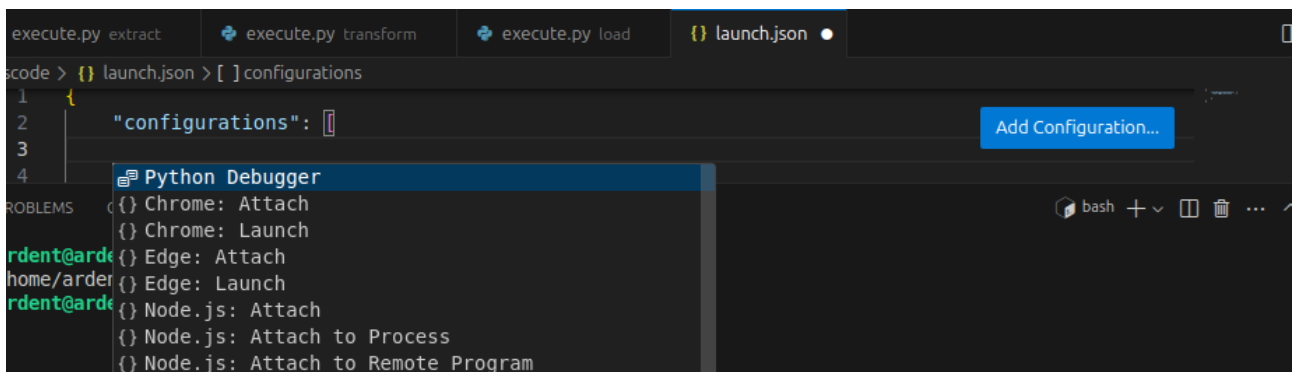
2. Create a new folder titled .vscode in your workspace directory. Make sure the name is exact ".vscode" as it is what the visual studio will look for



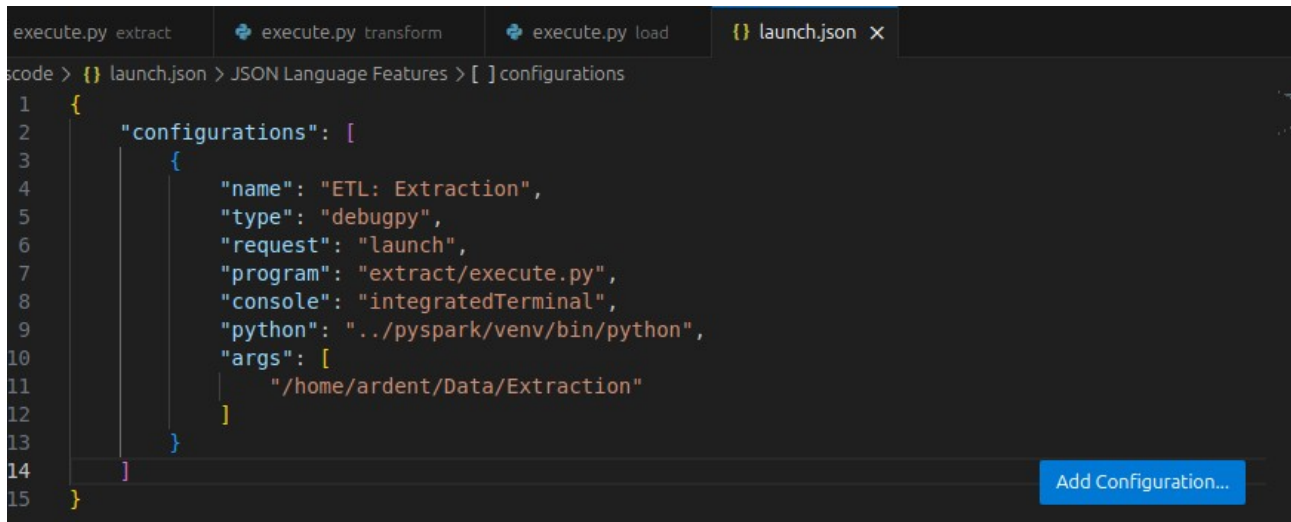
3. Inside the newly created .vscode directory create a file named launch.json



4. Click on Add Configuration button and then python debugger. Then select python file with arguments. If the add configuration or other options after that did not appear. You should finally type the config parameters shown in the next section



5. Modify the config parameter to point to your file and accommodate your needs



```
execute.py extract | execute.py transform | execute.py load | {} launch.json x
code > {} launch.json > JSON Language Features > [ ] configurations
1 {
2   "configurations": [
3     {
4       "name": "ETL: Extraction",
5       "type": "debugpy",
6       "request": "launch",
7       "program": "extract/execute.py",
8       "console": "integratedTerminal",
9       "python": "../pyspark/venv/bin/python",
10      "args": [
11        "/home/ardent/Data/Extraction"
12      ]
13    }
14  ]
15 }
```

The screenshot shows the VS Code editor with the `launch.json` file open. The file contains a JSON configuration for a debug session. The configuration is named "ETL: Extraction" and is of type "debugpy". It is a "launch" request, meaning it starts the program. The program to run is "extract/execute.py". The console is set to "integratedTerminal". The python executable is specified as "../pyspark/venv/bin/python", which is a relative path from the file's location. The arguments to pass are [" /home/ardent/Data/Extraction"]. A button "Add Configuration..." is visible in the bottom right corner of the editor window.

Here,

type parameter is set to debugpy to use the python debugger. Make sure you install it from extensions marketplace.

request parameter is set to launch to launch the program.

program parameter takes the path to the python file to run.

console is the terminal you want run the code on. We have selected integratedTerminal.

python parameter is the python executable to use to run the program. Note that I have used the relative path to provide the python from the pyspark virtual environment. Depending on use case using absolute path might be a better option

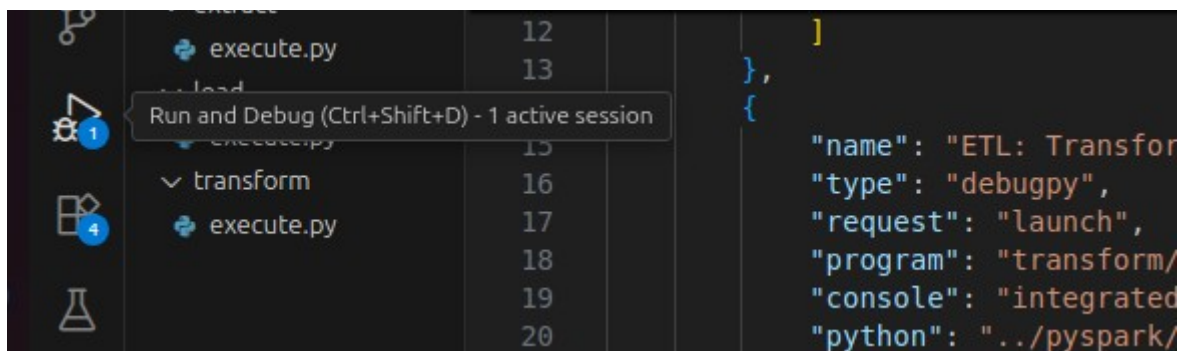
args parameter is the array of arguments to be passed to the program during execution

name is the name given to the specific configuration. You can add multiple configuration to the configurations list. **Try to do this for transform and load** before looking at below configuration. It is pretty straightforward

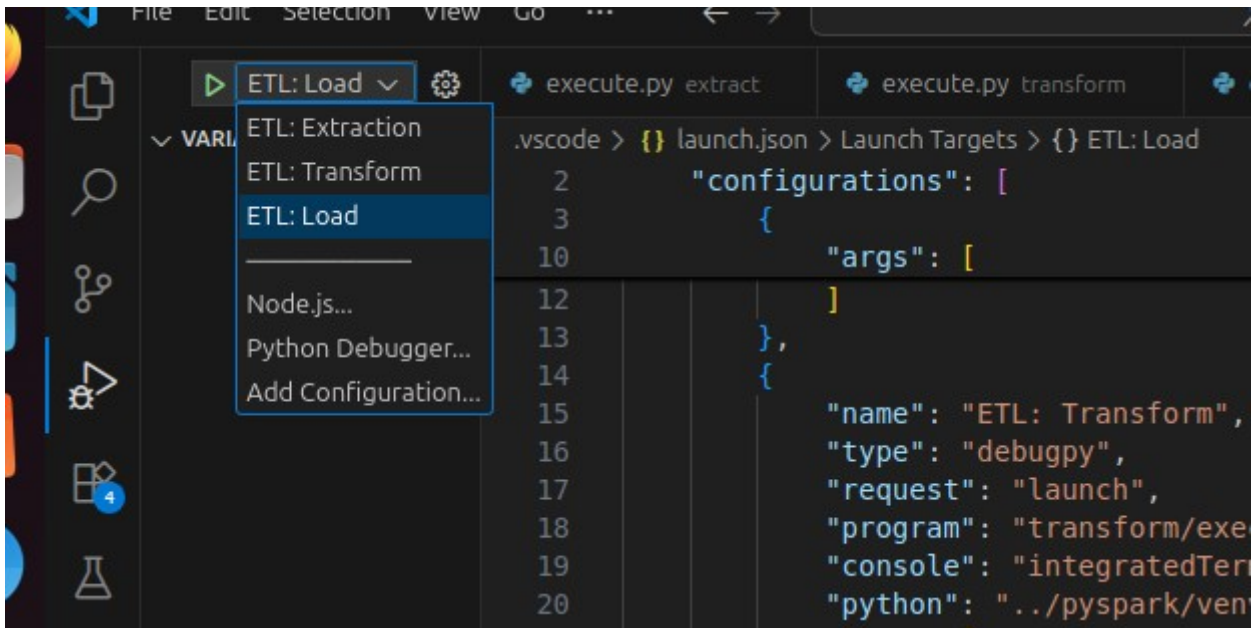
```
{
  "args": [
  ],
},
{
  "name": "ETL: Transform",
  "type": "debugpy",
  "request": "launch",
  "program": "transform/execute.py",
  "console": "integratedTerminal",
  "python": "../pyspark/venv/bin/python",
  "args": [
    "/home/ardent/Data/Extraction",
    "/home/ardent/Data/Transform"
  ]
},
{
  "name": "ETL: Load",
  "type": "debugpy",
  "request": "launch",
  "program": "load/execute.py",
  "console": "integratedTerminal",
  "python": "../pyspark/venv/bin/python",
  "args": [
    "/home/ardent/Data/transform"
  ]
}
}
```

Add Configuration

6. Switch to debugger tab by clicking on the debug icon at the left hand side



You should be able to see the name of the configurations you have set in launch.json file here. Now you can simply run the file by selecting the desire configuration and clicking the green button. This will run our code in debugger mode allow us to set breakpoint and access to watch window among many other features. Explore this on your own. This will be a valuable addition to your toolset



7. We need to modify our load module. Currently our load module contains sensitive information such as password in plain text. The code will eventually be pushed to github and it is very bad practice to expose password in plain text to public repositories.

There are several ways of solving this issue. The required parameters can be passed to the program during execution as an argument, or these can be set as environment variable, or several other password protector tools can be utilized. We will use the first approach as it is the easiest to implement in our current scenario.

7.1 Modify the entry point to take in username and password as input arguments:

```
if __name__ == "__main__":
    if len(sys.argv) != 4:
        print("Usage: python load/execute.py <input_dir> <pg_un> <pg_pw>")
        sys.exit(1)

    input_dir = sys.argv[1]
    pg_un = sys.argv[2]
    pg_pw = sys.argv[3]

    if not os.path.exists(input_dir):
```

Notice the len has been changed from 2 to 4 and two more variables are being stored

7.2 Modify the two function to take username and password as input

```
input_dir = sys.argv[1]
pg_un = sys.argv[2]
pg_pw = sys.argv[3]

if not os.path.exists(input_dir):
    print(f"Error: Input directory {input_dir} does not exist")
    sys.exit(1)

spark = create_spark_session()
create_postgres_tables(pg_un, pg_pw)
load_to_postgres(spark, input_dir, pg_un, pg_pw)

print("Load stage completed")
```

7.3 Modify the actual function definition to take it as input as well

```
.getOrCreate()

def create_postgres_tables(pg_un, pg_pw):
    """Create PostgreSQL tables if they don't exist using psycopg2."""
    conn = None
    cursor = None

    try:
        conn = psycopg2.connect(
            dbname="postgres",
            user=pg_un,
            password=pg_pw,
            host="localhost",
            port="5432"
        )
        cursor = conn.cursor()

        create_table_queries = [
```



```

8         conn.close()
9
10    def load_to_postgres(spark, input_dir, pg_un, pg_pw):
11        """Load Parquet files to PostgreSQL."""
12        jdbc_url = "jdbc:postgresql://localhost:5432/postgres"
13        connection_properties = {
14            "user": pg_un,
15            "password": pg_pw,
16            "driver": "org.postgresql.Driver"
17        }
18
19        tables = [
20            ("stage2/master_table", "master_table"),
21            ("stage3/recommendations_exploded", "recommendations_exploded"),
22            ("stage3/artist_track", "artist_track"),
23            ("stage3/track_metadata", "track_metadata"),
24            ("stage3/artist_metadata", "artist_metadata")
25        ]

```

7.4 Finally, modify the launch.json file we created earlier to pass these two new input arguments

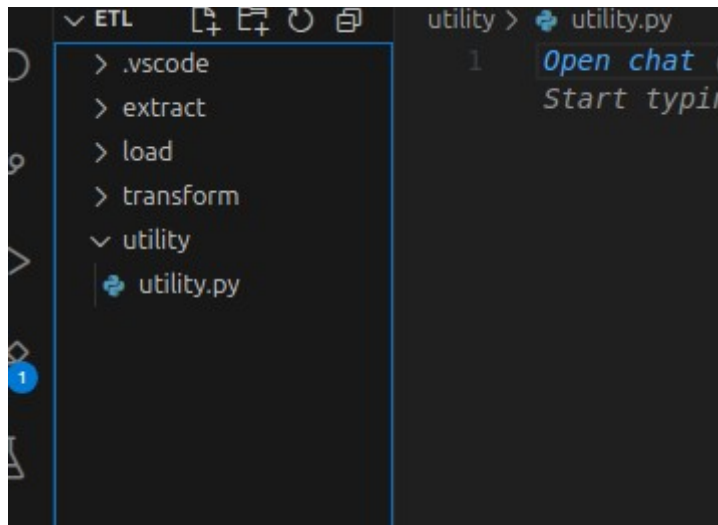
```

8    "python": "../pyspark/venv/bin/python",
9    "args": [
10        "/home/ardent/Data/Extraction", //input directory
11        "/home/ardent/Data/Transform" // output directory
12    ]
13 },
14 {
15     "name": "ETL: Load",
16     "type": "debugpy",
17     "request": "launch",
18     "program": "load/execute.py",
19     "console": "integratedTerminal",
20     "python": "../pyspark/venv/bin/python",
21     "args": [
22         "/home/ardent/Data/transform", //input directory
23         "postgres", //username
24         "postgres" //password
25     ]
26 }
27 ]
28 }
29 }
30 }

```

7.5 Pass the host, database, and port as an argument as well. Do this on your own

8. In the ETL directory. Create a new folder titled utility and create an new python file inside it titled utility.py



8.1 Let's setup logging functionality first. Create a function named `setup_logger` which will take the name of the log file and create instance of root logger. Then add two handlers. Handlers are responsible for handling the output of logs. We will setup one handler to output to stdout or terminal and the next handler to write to a log file. Logging is not limited to this. Learn about it in detail on your own

Do this in the newly create utility.py file

```
import logging

def setup_logging(log_file_name):

    logFormatter = logging.Formatter("%(asctime)s [%(threadName)-12.12s] [%(levelname)-5.5s] %(message)s")
    rootLogger = logging.getLogger()
    rootLogger.setLevel(logging.INFO)

    fileHandler = logging.FileHandler(log_file_name)
    fileHandler.setFormatter(logFormatter)
    rootLogger.addHandler(fileHandler)

    consoleHandler = logging.StreamHandler()
    consoleHandler.setFormatter(logFormatter)
    rootLogger.addHandler(consoleHandler)
    return rootLogger
```


8.3 Import the logger in the execute.py file of load directory

Notice that the absolute path to the etl folder has been added to the system path before importing from the utility directory. This is required because the python has no idea where to search for the utility file from

```
import os
import psycopg2
from pyspark.sql import SparkSession
sys.path.append(os.path.abspath(os.path.join(os.path.dirname(__file__), '..')))
from utility.utility import setup_logging

def create_spark_session(logger):
    """Initialize Spark session with PostgreSQL JDBC driver."""
    logger.debug("Initializing Spark Session with default parameters")
    return SparkSession.builder \
```

8.4 Let us modify the execute.py of load module to include the logger.

```
load > execute.py > ...
91 def load_to_postgres(spark, input_dir, pg_un, pg_pw):
115     print(f"Loaded {table_name} to PostgreSQL")
116     except Exception as e:
117         print(f"Error loading {table_name}: {e}")
118
119 if __name__ == "__main__":
120
121     logger = setup_logging("load.log")
122
123     if len(sys.argv) != 4:
124         logger.error("Usage: python load/execute.py <input_dir> <pg_un> <pg_pw>")
125         sys.exit(1)
126
127     input_dir = sys.argv[1]
128     pg_un = sys.argv[2]
129     pg_pw = sys.argv[3]
```

8.5 Replace all instances of print using logger. For function make sure logger is passed as first argument and later used throughout it instead of print. There are 4 different methods of logger that we will use mostly. logger.info, logger.debug, logger.warning, logger.error Use these accordingly based on the type of message you are trying to log

```

load > execute.py > ...
129     pg_pw = sys.argv[2]
130
131     if not os.path.exists(input_dir):
132         logger.error(f"Error: Input directory {input_dir} does not exist")
133         sys.exit(1)
134
135
136     spark = create_spark_session(logger)
137     create_postgres_tables(logger, pg_un, pg_pw)
138     load_to_postgres(logger, spark, input_dir, pg_un, pg_pw)
139
140     logger.info("Load stage completed")
141
142

```

Do the same for other functions and print statement as well.

```

load > execute.py > create_postgres_tables
7  def create_spark_session(logger):
12      .getOrCreate()
13
14  def create_postgres_tables(logger, pg_un, pg_pw):
15      """Create PostgreSQL tables if they don't exist using psycopg2."""
16      conn = None
17      cursor = None
18
19      try:
20          conn = psycopg2.connect(
21              dbname="postgres",
22              user=pg_un,
23              password=pg_pw,
24              host="localhost",
25              port="5432"
26          )
27          cursor = conn.cursor()
28
29          logger.debug("Successfully connected to postgres database")
30
31          create_table_queries = []
32

```

```

load > execute.py > ...
14  def create_postgres_tables(logger, pg_un, pg_pw):
75      followers FLOAT,
76      popularity INTEGER
77      );
78      """
79  ]
80
81      for query in create_table_queries:
82          cursor.execute(query)
83      conn.commit()
84      logger.info("PostgreSQL tables created successfully")
85
86  except Exception as e:
87      logger.warning(f"Error creating tables: {e}")
88  finally:
89      logger.debug("Closing connection and cursor to postgres db")
90      if cursor:
91          cursor.close()
92      if conn:
93          conn.close()
94

```

```

93     conn.close()
94
95 def load_to_postgres(logger, spark, input_dir, pg_un, pg_pw):
96     """Load Parquet files to PostgreSQL."""
97     jdbc_url = "jdbc:postgresql://localhost:5432/postgres"
98     connection_properties = {
99         "user": pg_un,
100         "password": pg_pw,
101     }
102
103     tables = [
104         ("stage3/artist_metadata", "artist_metadata")
105     ]
106
107     for parquet_path, table_name in tables:
108         try:
109             df = spark.read.parquet(os.path.join(input_dir, parquet_path))
110             mode = "append" if 'master' in parquet_path else "overwrite"
111             df.write \
112                 .mode(mode) \
113                 .jdbc(url=jdbc_url, table=table_name, properties=connection_properties)
114             logger.info(f"Loaded {table_name} to PostgreSQL")
115         except Exception as e:
116             logger.warning(f"Error loading {table_name}: {e}")
117
118

```

9. Make similar changes for execute and transform module on your own

10. In the utility.py file add new function that will take seconds as an input and return formatted string with hours, minutes and seconds

```

utility > utility.py > format_time
3 def setup_logging(log_file_name):
15     consoleHandler.setFormatter(LogFormatter)
16     rootLogger.addHandler(consoleHandler)
17     return rootLogger
18
19
20 def format_time(seconds):
21     hours, remainder = divmod(seconds, 3600)
22     minutes, seconds = divmod(remainder, 60)
23     return f"{int(hours)} hours, {int(minutes)} minutes, {int(seconds)} seconds"

```

10.1 Import this newly created function in the execute.py of load module Also import the time module

```

load > execute.py > ...
1 import time
2 import sys
3 import os
4 import psycopg2
5 from pyspark.sql import SparkSession
6 sys.path.append(os.path.abspath(os.path.join(os.path.dirname(__file__), '..')))
7 from utility.utility import setup_logging, format_time
8
9 def create_spark_session(logger):
10     """Initialize Spark session with PostgreSQL JDBC driver."""
11     logger.debug("Initializing Spark Session with default parameters")
12     return SparkSession.builder \

```

10.2 In the entry point start recording time before execute the main logic and record the time after the code has been executed

```
139         logger.error(f"Error: input directory {input_dir} does not exist")
140         sys.exit(1)
141     logger.info("Load stage started")
142     start = time.time()
143
144     spark = create_spark_session(logger)
145     create_postgres_tables(logger, pg_un, pg_pw)
146     load_to_postgres(logger, spark, input_dir, pg_un, pg_pw)
147
148     end = time.time()
149     logger.info("Load stage completed")
150
151
```

10.3 Finally pass the difference between these time to newly created function and log the output of that response

```
141     logger.info("Load stage started")
142     start = time.time()
143
144     spark = create_spark_session(logger)
145     create_postgres_tables(logger, pg_un, pg_pw)
146     load_to_postgres(logger, spark, input_dir, pg_un, pg_pw)
147
148     end = time.time()
149     logger.info("Load stage completed")
150     logger.info(f"Total time taken {format_time(end-start)}")
151
152
```

11. Do this for execute and transform module as well. Do this on your own. The process is exactly same as shown above

The below is the sample output after everything above has been implemented. Also a load.log file will be created in the etl directory. Similarly for execute and transform implement and verify it on your own

```
2025-07-27 17:50:05,252 [MainThread ] [INFO ] PostgreSQL tables created successfully
2025-07-27 17:50:22,582 [MainThread ] [INFO ] Loaded master_table to PostgreSQL
2025-07-27 17:50:56,700 [MainThread ] [INFO ] Loaded recommendations_explored to PostgreSQL
2025-07-27 17:50:59,704 [MainThread ] [INFO ] Loaded artist_track to PostgreSQL
2025-07-27 17:51:02,709 [MainThread ] [INFO ] Loaded track_metadata to PostgreSQL
2025-07-27 17:51:06,818 [MainThread ] [INFO ] Loaded artist_metadata to PostgreSQL
2025-07-27 17:51:06,819 [MainThread ] [INFO ] Load stage completed
2025-07-27 17:51:06,819 [MainThread ] [INFO ] Total time taken 0 hours, 1 minutes, 4 seconds
2025-07-27 17:51:07,495 [MainThread ] [INFO ] Closing down clientserver connection
(venv) ardent@ardent:~/Workspace/etl$
```

```
1 2025-07-27 17:41:02,863 [MainThread ] [INFO ] Load stage started
2 2025-07-27 17:41:06,372 [MainThread ] [INFO ] PostgreSQL tables created successfully
3 2025-07-27 17:50:01,833 [MainThread ] [INFO ] Load stage started
4 2025-07-27 17:50:05,252 [MainThread ] [INFO ] PostgreSQL tables created successfully
5 2025-07-27 17:50:22,582 [MainThread ] [INFO ] Loaded master_table to PostgreSQL
6 2025-07-27 17:50:56,700 [MainThread ] [INFO ] Loaded recommendations_exploded to PostgreSQL
7 2025-07-27 17:50:59,704 [MainThread ] [INFO ] Loaded artist_track to PostgreSQL
8 2025-07-27 17:51:02,709 [MainThread ] [INFO ] Loaded track_metadata to PostgreSQL
9 2025-07-27 17:51:06,818 [MainThread ] [INFO ] Loaded artist_metadata to PostgreSQL
10 2025-07-27 17:51:06,819 [MainThread ] [INFO ] Load stage completed
11 2025-07-27 17:51:06,819 [MainThread ] [INFO ] Total time taken 0 hours, 1 minutes, 4 seconds
12 2025-07-27 17:51:07,495 [MainThread ] [INFO ] Closing down clientserver connection
13
```

Git related:

Gitignore Documentation: <https://git-scm.com/docs/gitignore>

Commit message convention:

<https://gist.github.com/qoomon/5dfcdf8eec66a051ecd85625518cfd13>

Example of .gitignore file for python project:

<https://github.com/github/gitignore/blob/main/Python.gitignore>