# Explain me in detailed way why we are using each line what will happen how functions calls how the code related just explain me in detaild way. Think that I am preparing for the interview I am early learner of data engineer.

Certainly! Here is a detailed explanation of the entire code, line by line and function by function, with the flow and role of each part emphasized. The explanation is tailored for an early learner preparing for a data engineering interview focused on Apache Beam pipelines with BigQuery integration.

## Detailed Explanation of the Apache Beam Dataflow Pipeline Code

### Import Statements

```
import argparse
import json
import apache_beam as beam
from apache_beam.options.pipeline_options import PipelineOptions
import csv
```

- `argparse`: Used to read command-line arguments passed when running the Python script. This makes the script reusable for different inputs without hardcoding.

- `json`: To load JSON files such as config files, schema definitions, and pipeline options from disk into Python dicts for use in the pipeline.

- `apache_beam`: The core library to build batch or streaming data pipelines that can run on different runners like Google Dataflow.

- `PipelineOptions`: A Beam utility class that helps manage pipeline configuration. It converts basic Python dict options into a Beam-compatible options object.

- `csv`: CSV parsing is done via this Python standard library module to process raw CSV lines into separate columns, handling comma separation and quoting correctly.

**Function:** `load_json(path)`

```python
def load_json(path):
    with open(path) as f:
        return json.load(f)
```

- This function loads and parses JSON data from the file at the given `path`.
- The JSON could represent the pipeline's configuration, BigQuery schema, or pipeline options.
- Using this function improves code reusability and abstraction — you can load any JSON by just passing the path.

**Function:** `schema_to_bq_string(schema_list)`

```python
def schema_to_bq_string(schema_list):
    return ",".join(f"{field.get('name', '')}:{field.get('type', '')}" for field in schem
```

- BigQuery expects schemas as comma-separated `"field_name:field_type"` strings.
- This function converts the schema (a list of dicts each describing a field with `name` and `type`) into the string format BigQuery expects.
- Example: `[{"name":"id", "type":"STRING"}, {"name":"age", "type":"INTEGER"}]` → `"id:STRING,age:INTEGER"`.
- It calls `.get()` with defaults so it handles missing keys gracefully.

**Function:** `extract_field_properties(schema_list)`

```python
def extract_field_properties(schema_list):
    field_names = []
    required_fields = set()
    nullable_fields = set()
    for field in schema_list:
        name = field.get('name')
        mode = field.get('mode', 'NULLABLE').upper()
        field_names.append(name)
        if mode == 'REQUIRED':
            required_fields.add(name)
        elif mode == 'NULLABLE':
            nullable_fields.add(name)
    return field_names, required_fields, nullable_fields
```

- This function parses the BigQuery schema to prepare for data validation.
- It collects:
  - `field_names`: All fields in the schema, preserving order, to know column names.

- - `required_fields`: Fields that must have non-empty values (marked as `REQUIRED`).
    - `nullable_fields`: Fields that can legitimately be empty (`NULLABLE`).
- This helps when cleaning and validating CSV data to check existence and nullability constraints before loading.

## Function: `parse_csv_line(line, field_names)`

```
def parse_csv_line(line, field_names):
    values = next(csv.reader([line]))
    if len(values) < len(field_names):
        values.extend([None] * (len(field_names) - len(values)))
    return dict(zip(field_names, values))
```

- This takes one raw CSV line (a string), parses it into a list of values using `csv.reader` (which handles commas, quotes, escaping).
- It pads missing trailing fields with None if the line has fewer fields than expected.
- Then creates a dictionary mapping from each field name to its corresponding value (`field_names` come from schema).
- Result is a structured row dict for further processing.
- This transform is crucial to convert raw unstructured text data into structured format usable downstream.

## Function: `clean_row(row, nullable_fields)`

```
def clean_row(row, nullable_fields):
    if row is None:
        return None
    return {k: (None if (v in ("NULL", "")) and k in nullable_fields else v) for k, v in
```

- This cleans the row dictionary by handling values marked as `"NULL"` or empty string `""` in CSV data.
- For fields flagged as nullable, it replaces `"NULL"` or `""` with Python `None` to represent null.
- For other fields or non-nullable fields, keeps original value.
- Returns a cleaned row or None if input was None.
- Ensures proper null representation when loading into BigQuery to avoid errors or type mismatches.

**Function:** `validate_row(row, required_fields, field_names, table_name)`

```python
def validate_row(row, required_fields, field_names, table_name):
    if row is None:
        print(f"{table_name}: Dropping None row")
        return False
    if set(row.keys()) != set(field_names):
        print(f"{table_name}: Dropping row with unexpected fields: {row}")
        return False
    for field in required_fields:
        val = row.get(field)
        if val is None or val == '':
            print(f"{table_name}: Dropping row missing required '{field}': {row}")
            return False
    return True
```

- This function is a filter to validate rows before loading.
- It checks:
    - If the row is None (discard).
    - If any unexpected columns/fields are present or missing (discard).
    - If any required field is missing or empty after cleaning (discard).
- It logs the reason for discarding to console (useful for debugging).
- Returns True only when row fully passes validation, allowing it downstream.
- Avoids loading corrupted or incomplete data into BigQuery.

## Main function: `run(...)`

```python
def run(project_id, dataset_id, table_name, input_path, config_path, schema_path, options
```

- This is the core function executing the pipeline end-to-end, parameterized for flexibility.
- Inputs:
    - GCP project and dataset where BigQuery tables live.
    - The specific table to process.
    - File paths for input CSV, and JSON config/schema/options files.

Inside `run`:

```python
print(f"Running pipeline with table_name={table_name}, input_path={input_path}")
```

- Logs to identify pipeline start and parameters used.

Load external JSON configs:

```
config = load_json(config_path)
schemas = load_json(schema_path)
options_dict = load_json(options_path)
```

- Load table-specific config (optional, not deeply used here),
- Load BigQuery schemas dict (mapping table names to schema lists),
- Load Dataflow pipeline options (runner settings, temp location, etc.).

Create Beam pipeline options:

```
pipeline_options = PipelineOptions(**options_dict)
```

- Converts options dict to Beam compatible options object.

Extract required info for the specified table:

```
schema = schemas[table_name]
field_names, required_fields, nullable_fields = extract_field_properties(schema)
bq_schema_str = schema_to_bq_string(schema)
table_spec = f"{project_id}:{dataset_id}.{table_name}"
```

- Gets schema fields for the table.
- Prepares lists for validating data.
- Formats schema string for BigQuery.
- Creates fully qualified BigQuery table spec string.

Set up metric counter for processed rows — useful for monitoring:

```
rows_processed_counter = beam.metrics.Metrics.counter('pipeline', f'{table_name}_rows_pro
def count_row(row):
    rows_processed_counter.inc()
    return row
```

## Building the Pipeline

```
with beam.Pipeline(options=pipeline_options) as p:
```

- Context manager creates and runs the pipeline with provided options.

Pipeline steps (PTransforms):

```
p
| f"Read_{table_name}" >> beam.io.ReadFromText(input_path, skip_header_lines=1)
```

- Reads CSV file line by line from the input path.
- Skips header row.

```
| f"Parse_{table_name}" >> beam.Map(lambda r: parse_csv_line(r, field_names))
```

- Parses each CSV line into structured dictionaries keyed by field names.

```
| f"Clean_{table_name}" >> beam.Map(lambda r: clean_row(r, nullable_fields))
```

- Cleans rows by converting `"NULL"` or `""` in nullable fields to `None`.

```
| f"FilterNone_{table_name}" >> beam.Filter(lambda r: r is not None)
```

- Removes rows filtered out (None).

```
| f"Validate_{table_name}" >> beam.Filter(lambda r: validate_row(r, required_fields, fiel
```

- Keeps only rows which have valid field names and values per schema constraints.

```
| f"CountRows_{table_name}" >> beam.Map(count_row)
```

- Increments metrics counter for reporting.

```
| f"WriteToBigQuery_{table_name}" >> beam.io.WriteToBigQuery(
        table_spec,
        schema=bq_schema_str,
        write_disposition=beam.io.BigQueryDisposition.WRITE_TRUNCATE,
        create_disposition=beam.io.BigQueryDisposition.CREATE_IF_NEEDED,
        custom_gcs_temp_location=options_dict.get("temp_location"),
    )
```

- Writes processed, validated rows into the BigQuery target table.
- `WRITE_TRUNCATE` overwrites existing table data.
- `CREATE_IF_NEEDED` will create table if missing.
- Temporary files during load are stored in GCS temp location from options.

## Script Entry Point and Argument Parsing

```
if __name__ == "__main__":
    parser = argparse.ArgumentParser(...)
    ...
    args = parser.parse_args()

    run(
        args.project_id,
        args.dataset_id,
```

```
        args.table_name,
        args.input_path,
        args.config_path,
        args.schema_path,
        args.options_path,
    )
```

- This block makes the script runnable from the command line.

- Arguments are required for project, dataset, table, input file path.

- Optional config files have defaults.

- Parses arguments and calls `run()` with these to start the pipeline.

# Summary: How it all fits together in the pipeline

1. **Configuration loading:** JSON files provide schema and options externally so pipeline can be reused for different tables and settings.

2. **Schema parsing:** Before reading any data, schema is parsed to understand column names and validation rules.

3. **Reading raw input:** Text CSV files are read line by line.

4. **Parsing to dict:** CSV string rows parsed into structured dict rows.

5. **Cleaning data:** Replace "NULL" and "" with Python `None` for nullable fields, improving data quality.

6. **Validation:** Drop incomplete or malformed rows to avoid ingesting bad data.

7. **Metrics:** Count successfully processed rows to monitor data flow.

8. **Write to BigQuery:** Store processed data into Google's managed data warehouse in a clean, typed schema.

9. **Parameterization:** Everything is modular and configurable, usable for any table or dataset by command line inputs and JSON configs.


### Interview Tip

Explain your design focus on **modularity, data validation, and configurability**—key for production data pipelines. Emphasize how reusable functions make maintenance easier and how metrics help operational reliability.

Understand Beam concepts like **PCollections**, **PTransforms** (`Map`, `Filter`), and **PipelineOptions** are crucial to explain your pipeline well.

Being able to walk through the flow: **config → read → parse → clean → validate → count → write** clearly shows solid grasp on ETL/data pipeline patterns.

Feel free to ask if a particular function or concept needs even deeper explanation!