

Data & Analytics Case Study

XKCD Analytics Pipeline



Goals and high level overview

Goal: View and measure reader insights for xkcd comics.

High level implementation:

- Pull data from API endpoint
- Transform and create reader metrics
- Store in a data warehouse or DB for analysis

Requirements:

- Portable
- Easy to run
- Reliable results

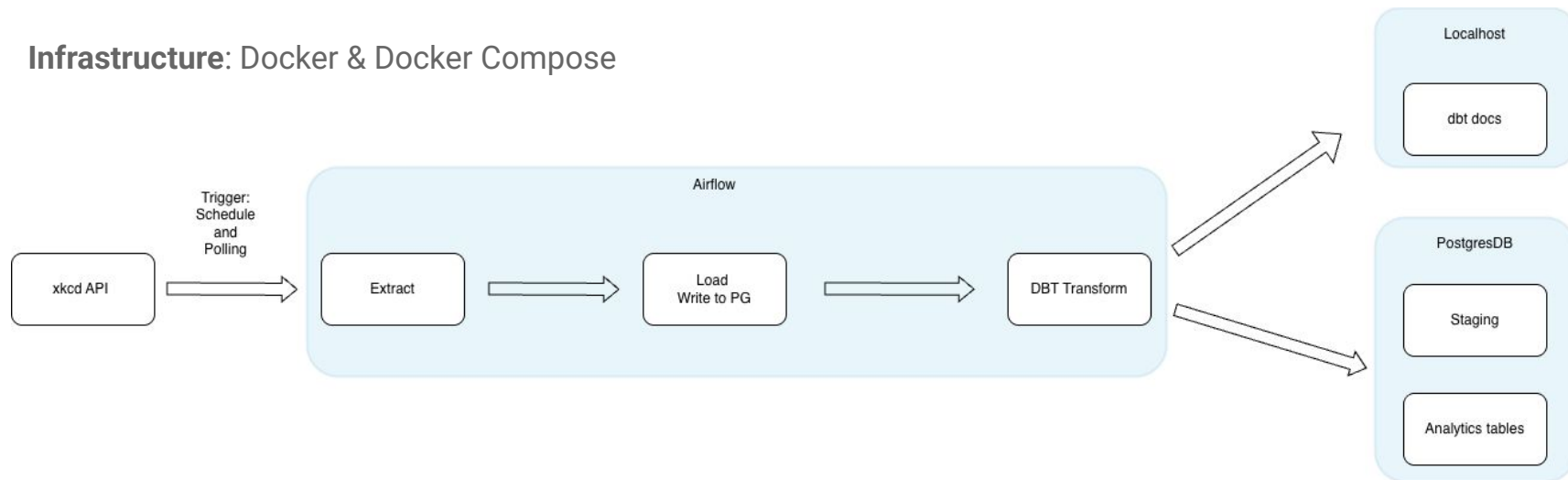
Flow Diagram and Tooling Choices

Orchestration: Apache Airflow

Database: PostgreSQL 16

Transformation: dbt Core

Infrastructure: Docker & Docker Compose



Extract and Load

Assumptions

- There is only one upload per comic
- Reader metrics are only uploaded once and not updated
- API structure is stable
- Every comic has a unique num
- Comics are published in numeric order
- API Results will always contain a D/M/Y

Extract and Load

Flow

- Calls API
- Checks if result is in DB, if not, load to comics table in PG.
- If ingested comic is not today's date, poll until the latest comic is in.
- If not, use a sensor to poll and detect when a new comic with today's date comes in.
- Once new comic is released continue with ingest and transforms.

Why this flow?

- Checking and updating the DB for comics that aren't today helps to catch if there was an unexpected or missed comic.
 - Not a full backfill but a nice addition for current scale
 - Useful for demos and locally running
- Sensor - Used for polling in airflow for newest comic
 - Includes a option to turn off completely. Useful for adhoc runs, demos, testing and the option to move to just scheduled runs if required.

Extract and Load

```
10 class XkcdClient:
11     """
12     XKCD API Client with exponential backoff.
13     """
14
15     def __init__(
16         self, max_retries=5, base_delay=1, max_delay=60, timeout=10, base_url=None
17     ):
18         self.max_retries = max_retries
19         self.base_delay = base_delay
20         self.max_delay = max_delay
21         self.timeout = timeout
22         self.base_url = base_url or os.getenv("XKCD_BASE_URL", "https://xkcd.com")
23
24     def _request(self, url):
25         retries = 0
26
27         while retries < self.max_retries:
28             try:
29                 resp = requests.get(url, timeout=self.timeout)
30                 resp.raise_for_status()
31                 return resp.json()
32             except requests.exceptions.RequestException as e:
33                 retries += 1
34                 if retries == self.max_retries:
35                     logger.error(f"Max retries reached: {e}")
36                     raise
37
38                 # exponential backoff
39                 delay = min(self.max_delay, self.base_delay * (2 ** (retries - 1)))
40
41                 logger.warning(
42                     f"Request failed ({e}). "
43                     f"Retrying (retries)/{self.max_retries} in {delay:.2f}s"
44                 )
45                 time.sleep(delay)
46
47     def get_latest_comic(self):
48         url = f"{self.base_url}/info.0.json"
49         return self._request(url)
```

```
8
9 def insert_comic_if_not_exists(engine, comic_data):
10     """Insert comic into database if it doesn't already exist."""
11
12     Session = sessionmaker(bind=engine)
13     session = Session()
14
15     try:
16         existing = session.query(Comic).filter(Comic.num == comic_data["num"]).first()
17         if existing:
18             logger.info(f"Comic {comic_data['num']} already exists")
19             return False
20
21         comic = Comic(
22             num=comic_data["num"],
23             title=comic_data["title"],
24             transcript=comic_data.get("transcript", ""),
25             alt=comic_data.get("alt", ""),
26             img=comic_data.get("img", ""),
27             day=int(comic_data["day"]),
28             month=int(comic_data["month"]),
29             year=int(comic_data["year"]),
30             published_date=date(
31                 int(comic_data["year"]),
32                 int(comic_data["month"]),
33                 int(comic_data["day"]),
34             ),
35         )
36
37         session.add(comic)
38         session.commit()
39         logger.info(f"Successfully inserted comic {comic_data['num']}")
40         return True
41
42     except Exception as e:
43         session.rollback()
44         logger.error(f"Error inserting comic {comic_data['num']}: {e}")
45         raise
46     finally:
47         session.close()
```

Components

/ingest/

api.py - API client

ingest_xkcd.py - main ingestion logic, passes xcom data for comic id and date for airflow

xkcd_sensor.py - sensor for polling for newest comic

/db/

connection.py - handles session starting in PG

models.py - model definition for raw comics table

operations.py - checks if comic already exists

Considerations and Improvements

Considerations

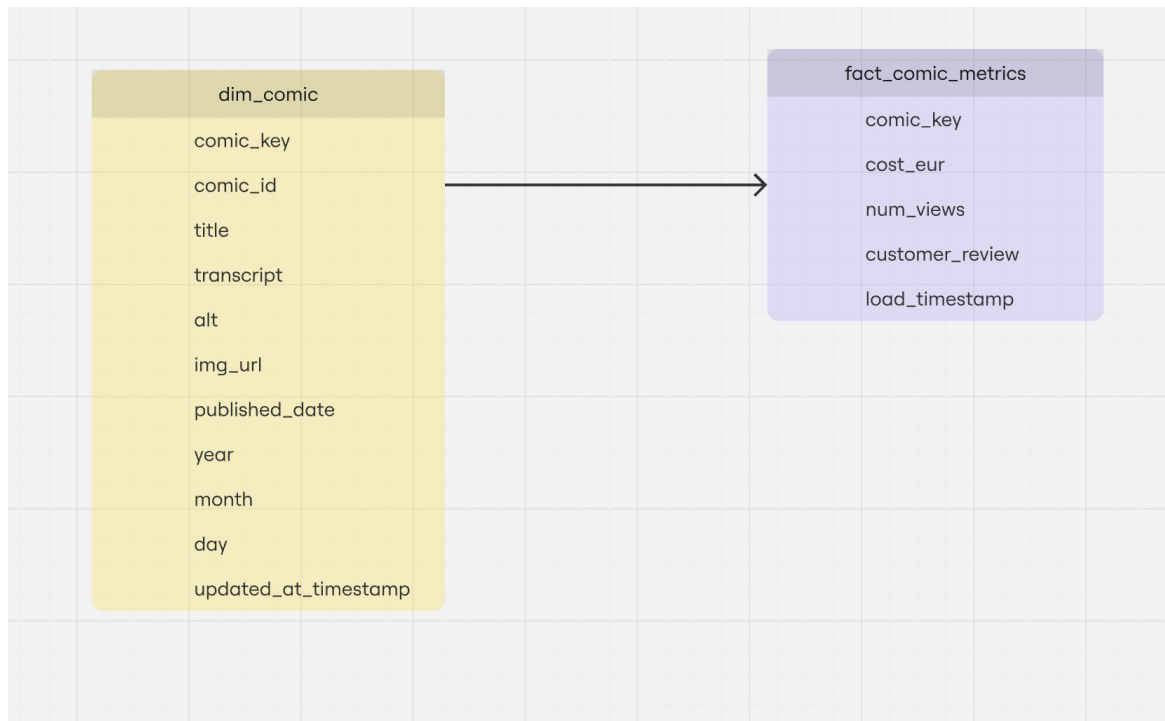
- API limitations
- Business scalability

Improvements

- Partitioning or clustering for large raw data sets to help with scaling
- Create larger backfill solution
- More unit and integration testing
- Schema validation
- Better logging and monitoring
- Connection pooling and bulk uploads

Transform – DWH (ER) Model

- Simple Star schema
- No physical constraints as we are using ELT.
- Validation is ran through DBT tests
- dim_comic uses Type 1 SDC
- comic_key is surrogate key



Transform

```
1 -- Test to ensure customer review scores are within valid range (1.0 to 10.0)
2 -- This test will FAIL if any rows are returned (indicating invalid review scores)
3 -- Customer reviews should be constrained between 1.0 (lowest) and 10.0 (highest)
4 SELECT *
5 FROM {{ ref('fact_comic_metrics') }}
6 WHERE customer_review < 1.0 OR customer_review > 10.0
7
```

```
1 -- Test to ensure customer review scores are within valid range (1.0 to 10.0)
2 -- This test will FAIL if any rows are returned (indicating invalid review scores)
3 -- Customer reviews should be constrained between 1.0 (lowest) and 10.0 (highest)
4 SELECT *
5 FROM {{ ref('fact_comic_metrics') }} f
6 LEFT JOIN {{ ref('dim_comic') }} d
7     ON f.comic_key = d.comic_key
8 WHERE d.comic_key IS NULL
```

```
1 {{ config(
2     materialized = 'table'
3 ) }}
4
5 SELECT
6     d.comic_key,
7     LENGTH(d.title) * 5 AS cost_eur,
8     FLOOR(random() * 10000) AS num_views,
9     ROUND((1 + (random() * 9))::numeric, 1) AS customer_review,
10    CURRENT_TIMESTAMP AS load_timestamp
11 FROM {{ ref('dim_comic') }} d
12
13 {% if is_incremental() %}
14     WHERE d.updated_at_timestamp > (SELECT MAX(load_timestamp) FROM {{ this }})
15 {% endif %}
```

dbt workflow:

- Staging layer (stg_comics view) - clean, renamed, type-cast data.
- Dimension table (dim_comic) - static metadata about comics.
- Fact table (fact_comic_metrics) - metrics: cost_eur, views, customer_review.

Business logic implemented:

- Cost: length(title) * 5 EUR.
- Views: random() * 10000.
- Customer reviews: random 1.0–10.0.

Data Quality Checks:

- dbt tests: not_null, unique, relationships, accepted_values, expression_is_true.

Considerations and Improvements

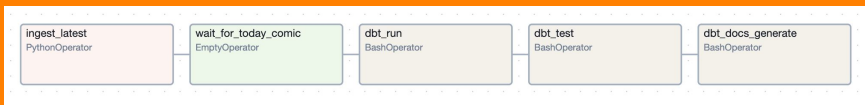
Considerations

- Business requirements

Improvements

- Randomised metrics
 - Randomised metrics may fail certain dbt tests if misconfigured.
 - Currently not idempotent (okay for this brief but not if business requirements change)
- Data quality models (e.g. comics with missing images, comics with invalid dates)
 - Helps monitor and measure “bad” data
- Performance boosting
 - Indexing for common queries
- More complex business metrics or additional dimensions (authors, tags).

Orchestration and Infrastructure



Workflow

- Airflow DAG that triggers both ingestion and dbt transformations.
- DBT tests run automatically at the end of the DAG.
- All ran within a dockerfile along with DBT and PG

Why is it beneficial

- Replicates production environment
- Reduces manual work and custom orchestration code (e.g. restarts, polling, alerting)
- Automates all processes
- Portability

Challenges

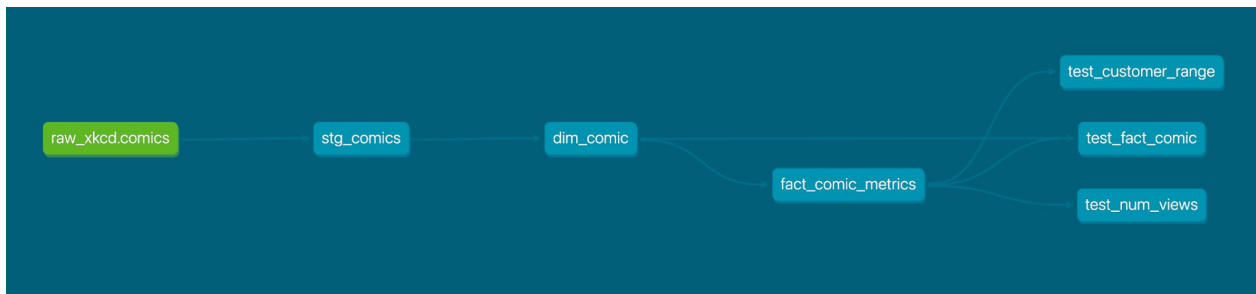
Ingestion

- Finding a clean way to balance ad-hoc runs and scheduled runs with polling

Infrastructure

- Setting up airflow/dbt/pg docker
 - Issues with docker-compose, had to break down Airflow deployment a lot.
 - Errors with user creation and unhealthy schedulers

End Output



Script Editor: *xkcd> Script

```
SELECT
  d.comic_id,
  d.title,
  d.published_date,
  f.cost_eur,
  f.num_views,
  f.customer_review
FROM dim_comic d
LEFT JOIN fact_comic_metrics f USING (comic_key);
```

Table: dim_comic(+) 1 X

Query: SELECT d.comic_id, d.title, d.published_date, d.cost_eur, d.num_views, d.customer_review

	comic_id	title	published_date	cost_eur	num_views	customer_review
1	3,171	Geologic Core Sample	2025-11-21	100	5,425	6.5

Productionisation

Further Production environment considerations:

- CI/CD
- Monitoring / alerting on DAG failures.
- Tooling choices (Use managed services)
- Container orchestration
- Data governance
- Access control (secrets manager, Airflow RBAC)
- Caching for API calls
- IAC

Thank you!