



GitHub Archive ETL for Sparrow Analytics

Team MINIRO



GitHub Archive Data

a project to **record** the public GitHub timeline, **archive it**, and **make it easily accessible** for further analysis

Event Examples

- Push
commits are pushed to a repository branch or tag.
- Gollum
human-readable files created or updated
- Create
a git branch or is created
- Pull Request
a contributor requests permission to merge a protected branch



Data Lake

Azure Data Lake Storage Gen2 = ADLS Gen1 + Blob storage

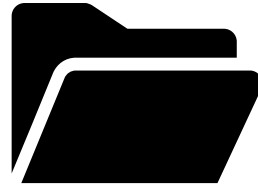
- hierarchical namespace, file-level security
- Blob storage



file system performance at object storage scale and prices

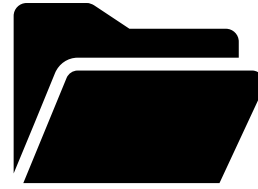


File Structure



LANDING ZONE

Compressed JSON format
Raw GitHub Archives



BRONZE LAYER

Parquet format
Partitioned by day



SILVER LAYER

Cleaned data
Snowflake structure



GOLD LAYER

Aggregated data
Ready for data warehouse

Cleaning

Goal: Process the data so that it is in a flattened tabular format

- Keep as much data as possible
 - Futureproof for any new business usages
- Drop any unnecessary columns
 - Columns that are all null, columns that are a url or contain data on 'gravatar' (globally recognized avatar)



Flattening the data

Issue: Nested Values

```
df: pyspark.sql.dataframe.DataFrame
└─ actor: struct
  │ created_at: string
  │ id: string
  │ org: struct
  └─ payload: struct
    │ action: string
    │ before: string
    │ comment: struct
    │ commits: array
    │ │ element: struct
    │ │ │ author: struct
    │ │ │ │ email: string
    │ │ │ │ name: string
    │ │ │ │ distinct: boolean
    │ │ │ │ message: string
    │ │ │ │ sha: string
    │ │ │ │ url: string
    │ │ description: string
    │ │ distinct_size: long
    │ forkee: struct
    │ │ archive_url: string
    │ │ assignees_url: string
```

Solution:

Unnest the JSON structs into separate columns



Issue: JSON data is Flexible, Tables are not

```
{  
  "id": 34,  
  "type": "car",  
  "make": "Honda",  
  "opts": {  
    "seats": 5,  
    "airbags": 0,  
    "ac": True  
  }  
}
```

```
{  
  "id": 87,  
  "type": "motorbike",  
  "make": "Ducati",  
  "opts": {  
    "sidecar": True,  
    "transmission": "chain"  
  }  
}
```

Vehicle DataFrame							
id	type	make	opts_seats	opts_airbags	opts_ac	opts_sidecar	opts_windshield
34	car	honda	5	0	true	NULL	NULL
87	motorbike	ducati	NULL	NULL	NULL	True	chain

Flattening the data

Issue: Event Types contain only some payload columns

payload_distinct_size ▲	payload_forkee ▲	payload_head ▲	payload_issue ▲
null	null	null	null
null	null	null	null
1	null	a76e751b921a74f0c58f0c5f680c4415d88f04fd	null
null	null	null	null
1	null	d9d3b02535f14610efdbd9a3e543d7318d34e774	null
null	null	null	▶ {"assignee": null, "assignees": [], "body": "This was crea "comments_url": "https://api.github.com/repos/SouthAf "2017-01-01T18:45:18Z", "events_url": "https://api.githu "html_url": "https://github.com/SouthAfricaDigitalScienc "default": false, "id": 308743652, "name": "build", "url": "f deploy/labels/build"}, {"color": "ededed", "default": false "https://api.github.com/repos/SouthAfricaDigitalScience "https://api.github.com/repos/SouthAfricaDigitalScience pull_number": 123, "pull_request": null, "repository_url"

Solution:

- > Separate the data into 14 different Event Type tables
- > Flatten one layer only
- > Drop null columns
- > Flatten completely

Two Approaches to Flatten

Hard Coding

Pros:

- ✓ Explicit about the kept columns
- ✓ Faster run time

Cons:

- Manually selecting each column for every flattened layer
- Needs to be changed if the schema of the data changes

Automating

Pros:

- ✓ Reusable code for data with JSON type columns
- ✓ Data does not have to be manually explored for each flattened layer
- ✓ Templates can be used to fix the issue of changing data schemas

Cons:

- Slower run time
- More difficult to debug

Issue: Data Schema Changes

The data added....

- contains extra columns
 - Any extra columns need to be dropped
- does not contain the columns from the current schema
 - Add null columns if data doesn't contain a column that is in the template

Solution: Templates Ensure Integrity

Template DataFrame			
a	b	c	d
...
...

New DataFrame			
d	c	a	e
...
...
...

Add template columns
that are not in the
DataFrame

New DataFrame				
d	c	a	e	b
...	null
...	null
...	null

Keep only template
columns, in order

New DataFrame			
a	b	c	d
...	null
...	null
...	null

ERD

