# Customizing Xcom for data sharing between tasks

Vikram Koka and Ephraim Anierobi

Airflow
Summit 2021

# Introductions

## Vikram Koka

Apache Airflow Committer

Senior Vice President Engineering at Astronomer
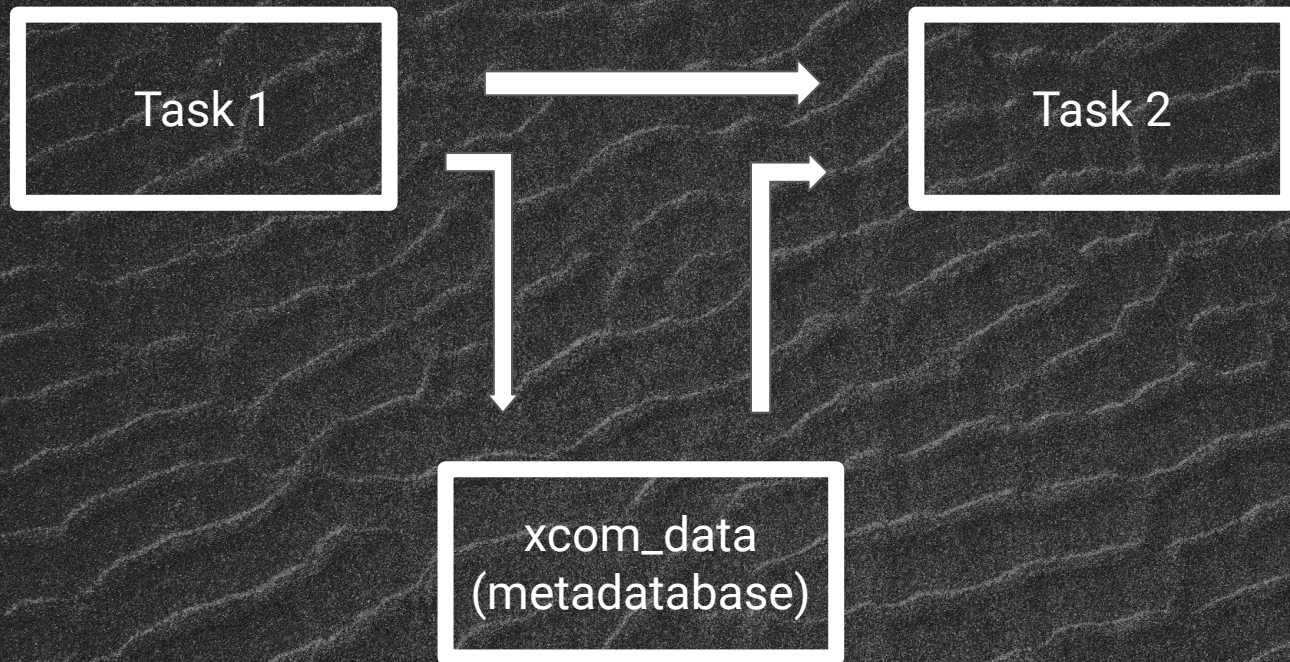
Silicon Valley

## Ephraim Anierobi

Apache Airflow Committer

Software Engineer, Open Source at Astronomer

Nigeria

# Xcom in a nutshell



Task 1

Task 2

xcom_data
(metadatabase)

# Xcom Overview

**Cross communication between tasks**

- Pass parameters from one task to another

- Supports multiple parameters

- Identified by key

- Intended for use within a single DAG

Usage:

- "push" and "pull"

Uses the Airflow metadatabase (Postgres / MySQL)

```
xcom_push(
    key = 'return_value',
    value = 'my value'
)


value = xcom_pull(
    task_ids='pushing_task',
    key='return_value'
)
```

# Xcom with TaskFlow API

## Greater Abstraction

- Return values implicitly use xcom

- Focused on the most common pattern

- Supports python native types including dict

Pythonic functional use
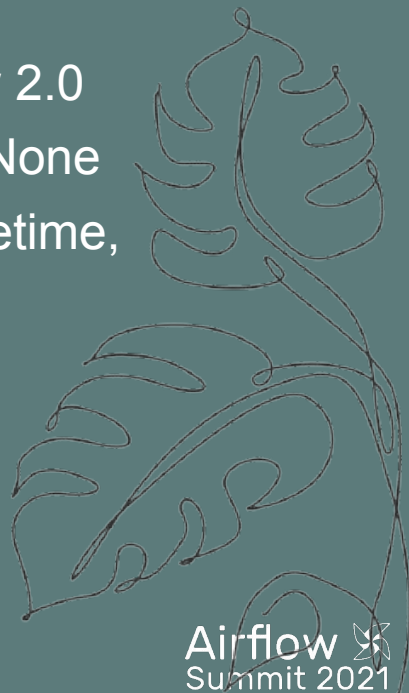
```
def extract:

    …

    return order_data

    …



order_data  = extract()

order_summary =
transform(order_data)
```

# Xcom limitations

## Data types

As it stands, only the following datatypes are supported in Airflow 2.0

- Python native: dict, list, tuple, str, int, long, float, True, False, None

- Future: Airflow supported objects such as numpy objects, datetime, date, etc

- For security, pickling is no longer recommended

Airflow
Summit 2021

# Github issues

## Sample questions / problems

- Unable to store xcom because of MySQL Blob type limitation 65,535

- Data too long when pushing to XCOM

- Raise do_xcom_push size limit

- Lambda to transform response before xcom push

- Provide shared storage between task via pluggable storage providers akin to S3 remote logging

# Custom XCom backends

## Persistence class

- Python class specified in config

- Read at Airflow start up,

  Class needs to be in Airflow path

## Methods needed:

- `serialize_value`
- `deserialize_value`

Used for storing and restoring data
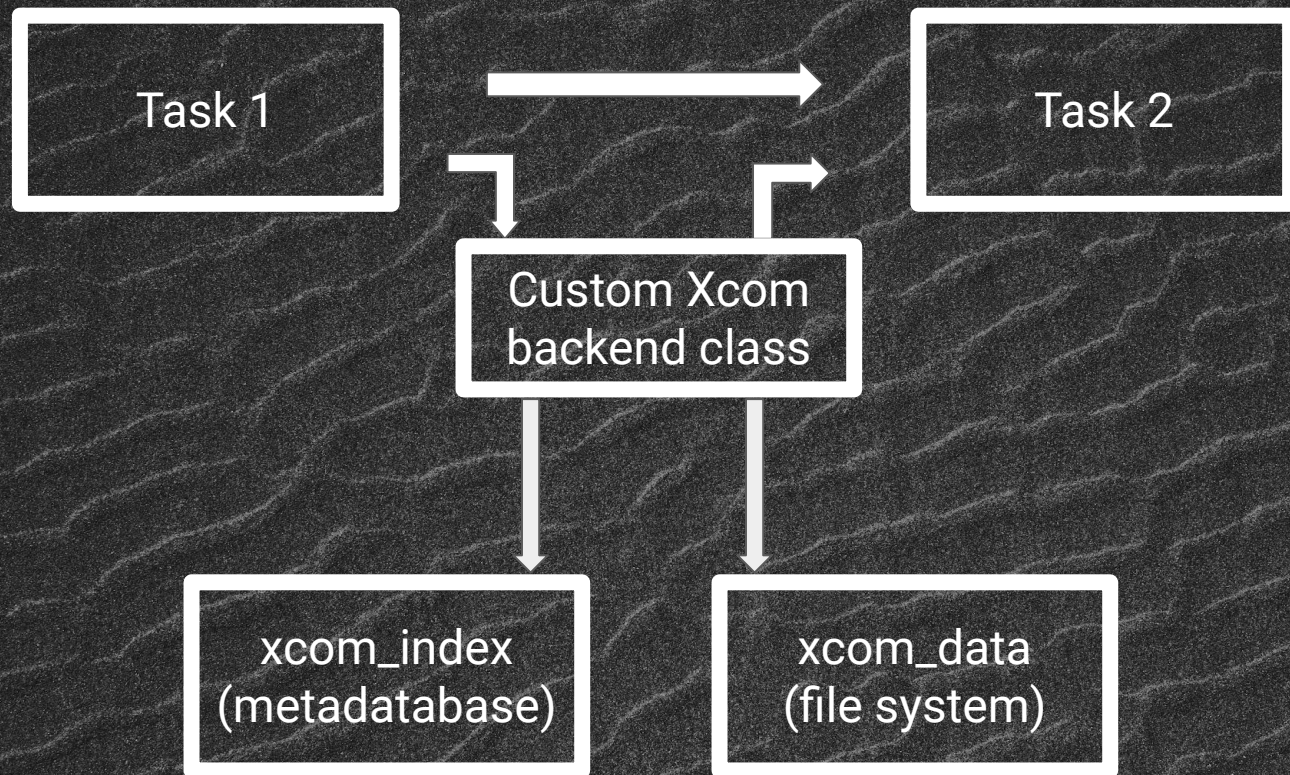
- `orm_deseralize_value`

Used to display XCom data in UI

Airflow
Summit 2021

# Custom XCom for Local Execution

- Write / read local file system
- Essential for development
- Local Executor

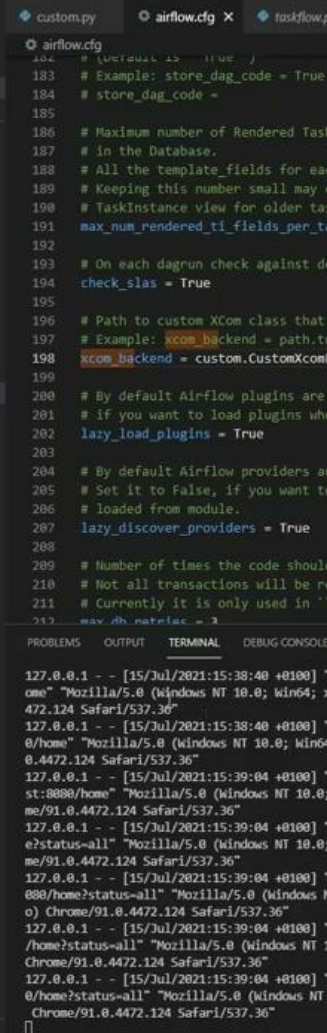Not for distributed deployments with celery and Kubernetes Executors

# Xcom stored in local file system



Task 1

Task 2

Custom Xcom backend class

xcom_index (metadatabase)

xcom_data (file system)

Airflow
Summit 2021

Code walk through and demo

Airflow
Summit 2021

Code editor view — airflow.cfg

```
182   # (default is   True )
183   # Example: store_dag_code = True
184   # store_dag_code =
185
186   # Maximum number of Rendered Task Instance Fields (Template Fields) per task to store
187   # in the Database.
188   # All the template_fields for each of Task Instance are stored in the Database.
189   # Keeping this number small may cause an error when you try to view ``Rendered`` tab in
190   # TaskInstance view for older tasks.
191   max_num_rendered_ti_fields_per_task = 30
192
193   # On each dagrun check against defined SLAs
194   check_slas = True
195
196   # Path to custom XCom class that will be used to store and resolve operators results
197   # Example: xcom_backend = path.to.CustomXCom
198   xcom_backend = custom.CustomXcomFS
199
200   # By default Airflow plugins are lazily-loaded (only loaded when required). Set it to ``False``,
201   # if you want to load plugins whenever 'airflow' is invoked via cli or loaded from module.
202   lazy_load_plugins = True
203
204   # By default Airflow providers are lazily-discovered (discovery and imports happen only when required).
205   # Set it to False, if you want to discover providers whenever 'airflow' is invoked via cli or
206   # loaded from module.
207   lazy_discover_providers = True
208
209   # Number of times the code should be retried in case of DB Operational Errors.
210   # Not all transactions will be retried as it can cause undesired state.
211   # Currently it is only used in ``DagFileProcessor.process_file`` to retry ``dagbag.sync_to_db``.
212   max_db_retries = 3
```

# Airflow distributed execution

# Airflow distributed execution - kubernetes

# Custom XCom for Distributed Execution

- Write / read cloud storage
- Accessible from any configured node
- Can be used with Celery and Kubernetes Executors

Higher latency, so could cause delays when used with short running tasks

More expensive than other options

# Custom Xcom for Distributed Execution

- Write / read from Redis
- Accessible from any configured node
- Can be used with Celery and Kubernetes Executors
- Already part of the Airflow stack

Size limit of 512MB, so ideal for smaller dataset between short running tasks.

Another caveat that Redis keeps everything in memory.

# Xcom stored in Redis or Cloud Storage

Code walk through and demo

# Clutter

- Clean-up of old data in cloud storage or elsewhere

- As data gets larger, data cleanup becomes more important

- System performance can degrade

# Clean-up DAG

- Maintenance DAG to clean-up old Xcom data

- Deletes data from metadatabase and external locations

- Not tied to DAG lifecycle- needs to be configured carefully

- Downside if trying to rerun old tasks

Code walk through and demo

Airflow
Summit 2021

File  Edit  Selection  View  Go  Run  Terminal  Help

EXPLORER

∨ OPEN EDITORS
  × ● maintenance.py  dags
    ● custom_redis.py  config
∨ XCOMBACKEND [WSL:...
  ∨ config
    > __pycache__
    ● custom_gcs.py
    ● custom_redis.py
    ● custom.py
  ∨ dags
    > __pycache__
    ● example_xcom_pandas.py
    ● maintenance.py
    ● taskflow.py
  > env
  > logs
  ∨ tmp
    ≡ airflow-webserver.pid
    ✿ airflow.cfg
    ≡ airflow.db
    ≡ dump.rdb
    ● webserver_config.py

● maintenance.py ×    ● custom_redis.py

dags > ● maintenance.py > ...

```python
19                prefix= GCSXCOM_ )
20                dag=dag
21            )
22
23    @dag.task()
24    def clean_fs_xcom():
25        for filename in glob.glob('tmp/XCOM_*'):
26            print("Deleting file :", filename)
27            os.remove(filename)
28
29    @dag.task()
30    def clean_redis_xcom():
31        r = redis.Redis()
32        keys = r.keys(pattern='RedisXCOM_*')
33        if len(keys):
34            print('XCOM keys found: ', keys)
35            r.delete(*keys)
36
37    @dag.task()
38    def clean_db_xcom():
39        with create_session() as session:
40            session.query(XCom).delete()
41
42    [clear_gcs_xcom, clean_fs_xcom(),  clean_redis_xcom()] >> clean_db_xcom()
```

PROBLEMS   OUTPUT   TERMINAL   DEBUG CONSOLE

```
lean_fs_xcom&execution_date=2021-07-15T20%3A11%3A
17.330846%2B00%3A00" "Mozilla/5.0 (Windows NT 10.
0; Win64; x64) AppleWebKit/537.36 (KHTML, like Ge
cko) Chrome/91.0.4472.124 Safari/537.36"
127.0.0.1 - - [15/Jul/2021:21:16:31 +0100] "GET /
get_logs_with_metadata?dag_id=example_maintenance
&task_id=clean_fs_xcom&execution_date=2021-07-15T
20%3A11%3A17.330846%2B00%3A00&try_number=1&metada
ta=null HTTP/1.1" 200 2750 "http://localhost:8080
/log?dag_id=example_maintenance&task_id=clean_fs_
xcom&execution_date=2021-07-15T20%3A11%3A17.33084
6%2B00%3A00" "Mozilla/5.0 (Windows NT 10.0; Win64
; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chr
ome/91.0.4472.124 Safari/537.36"
127.0.0.1 - - [15/Jul/2021:21:16:32 +0100] "GET /
log?dag_id=example_maintenance&task_id=clean_fs_x
com&execution_date=2021-07-15T20%3A11%3A17.330846
%2B00%3A00 HTTP/1.1" 200 40283 "-" "Mozilla/5.0 (
Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (
KHTML, like Gecko) Chrome/91.0.4472.124 Safari/53
7.36"
```

```
pool', 'default_pool', '--subdir', '/home/ephraimb
uddy/Documents/xcombackend/dags/maintenance.py']
[2021-07-15 21:11:21,276] {dagbag.py:496} INFO - F
illing up the DagBag from /home/ephraimbuddy/Docum
ents/xcombackend/dags/maintenance.py
Running <TaskInstance: example_maintenance.clean_d
b_xcom 2021-07-15T20:11:17.330846+00:00 [queued]>
on host DESKTOP-RMQKCOB.localdomain
[2021-07-15 21:11:22,042] {dagrun.py:444} INFO - M
arking run <DagRun example_maintenance @ 2021-07-1
5 20:11:17.330846+00:00: manual__2021-07-15T20:11:
17.330846+00:00, externally triggered: True> succe
ssful
[2021-07-15 21:11:22,053] {scheduler_job.py:1222}
INFO - Executor reports execution of example_maint
enance.clean_db_xcom execution_date=2021-07-15 20:
11:17.330846+00:00 exited with status success for
try_number 1
[2021-07-15 21:15:46,857] {scheduler_job.py:1839}
INFO - Resetting orphaned tasks for active dag run
s
```

```
44778:M 15 Jul 2021 20:31:48.410 # Server initiali
zed
44778:M 15 Jul 2021 20:31:48.410 * Loading RDB pro
duced by version 6.2.3
44778:M 15 Jul 2021 20:31:48.410 * RDB age 34 seco
nds
44778:M 15 Jul 2021 20:31:48.410 * RDB memory usag
e when created 0.51 Mb
44778:M 15 Jul 2021 20:31:48.410 * DB loaded from
disk: 0.000 seconds
44778:M 15 Jul 2021 20:31:48.410 * Ready to accept
 connections
```

```
(env) → xcombackend redis-cli
127.0.0.1:6379> ping
PONG
127.0.0.1:6379> get RedisXCOM_db66dc8d-8e89-4a27-
bc0d-ea1d56ce3b8a
"{\"columns\":[\"a\",\"b\",\"c\"],\"index\":[0,1,
2],\"data\":[[1,2,3],[4,5,6],[7,8,9]]}"
127.0.0.1:6379> get RedisXCOM_db66dc8d-8e89-4a27-
bc0d-ea1d56ce3b8a
"{\"columns\":[\"a\",\"b\",\"c\"],\"index\":[0,1,
2],\"data\":[[1,2,3],[4,5,6],[7,8,9]]}"
127.0.0.1:6379> get RedisXCOM_db66dc8d-8e89-4a27-
bc0d-ea1d56ce3b8a
(nil)
127.0.0.1:6379>
```

python3
python3
redis-server
redis-cli

# Success

We have addressed the core questions raised
- Handling of non-native objects such as Dataframes
- Large data sets between tasks
- Leveraging cloud storage
- Maintenance and cleanup

# Limitations

Not tied to DAG life cycle management
- Data sharing across DAGs is difficult
- Maintenance DAGs for clean-up is a kludge

- Should be cleanly handled by Airflow when DAG is done

# Future: Top level data object in Airflow

- Result of DAGs from one team is data
- Can be used by DAGs from other teams

- Key for cross-DAG dependencies
- Availability can be used to trigger follow-on DAGs

Integrated with DAG life cycle management and with Event driven DAGs
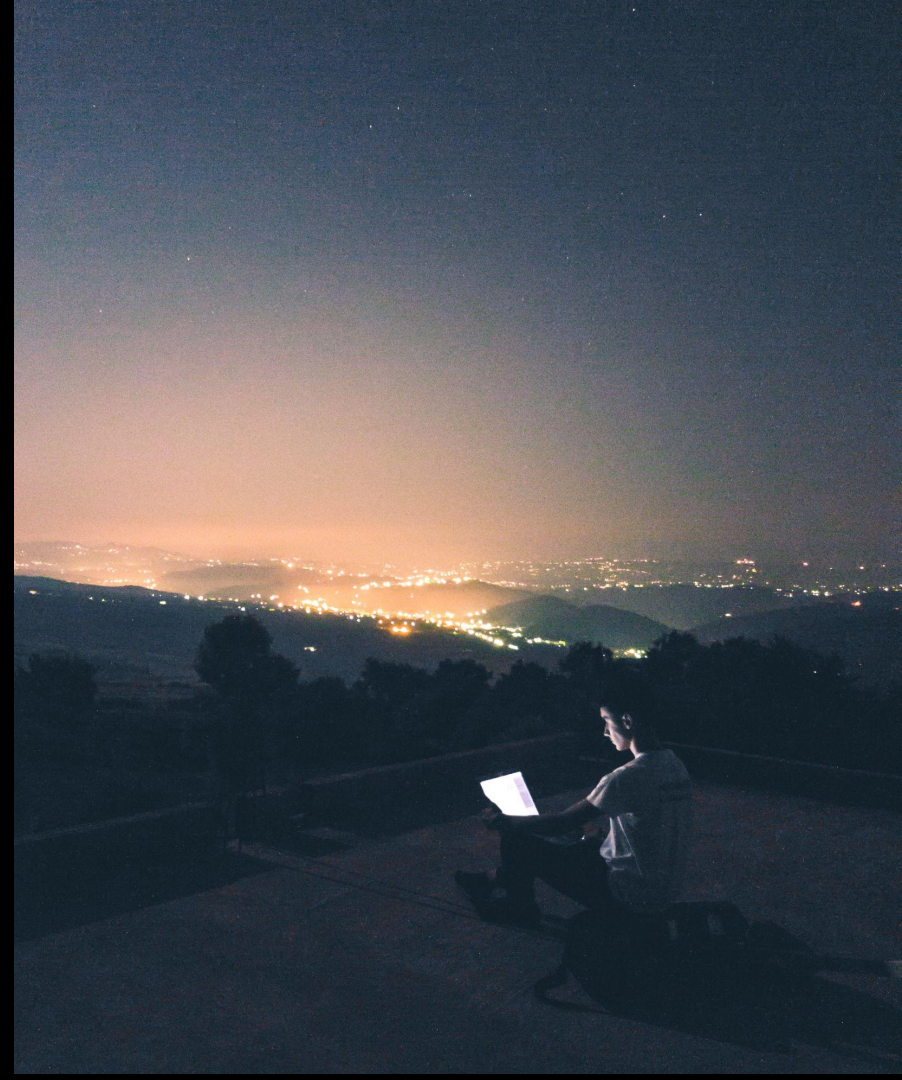
Airflow Improvement Proposal upcoming

# Jobs at Astronomer

We are hiring Airflowers all over the world!

https://careers.astronomer.io/

https://linkedin.com/vikramkoka

Contact us: We would love to hear from you!

ASTRONOMER

Questions?

Airflow
Summit 2021