

# Investigating the many loops of the Airflow Scheduler

Philippe Gagnon





# Philippe Gagnon

Your speaker today



Solutions Architecture at  
Astronomer, inc.



Based in Montreal, Canada



Works on data platform architecture  
and implementation in heavily  
regulated industries since 2017,  
mostly around open-source



# What is covered

1. Airflow task scheduling/execution components and their role
2. The scheduler initialization process
3. Task scheduling framework at a high level
4. The scheduler timers and what they do
5. How DagRuns are created
6. The scheduler "critical section" and TaskInstance handling
7. How do executors pick up TaskInstances? CeleryExec and KubeExec
8. How does the task actually "run"? CeleryExec and KubeExec

# Schedulers... DAG Processors..? Executors... Workers ?!

**Scheduler:** Responsible for adding jobs to the queue when their dependencies are met and triggering the execution of tasks.

**DAG Processor:** Parses, processes and serializes the DAG files. It can either run as part of the scheduler, or standalone.

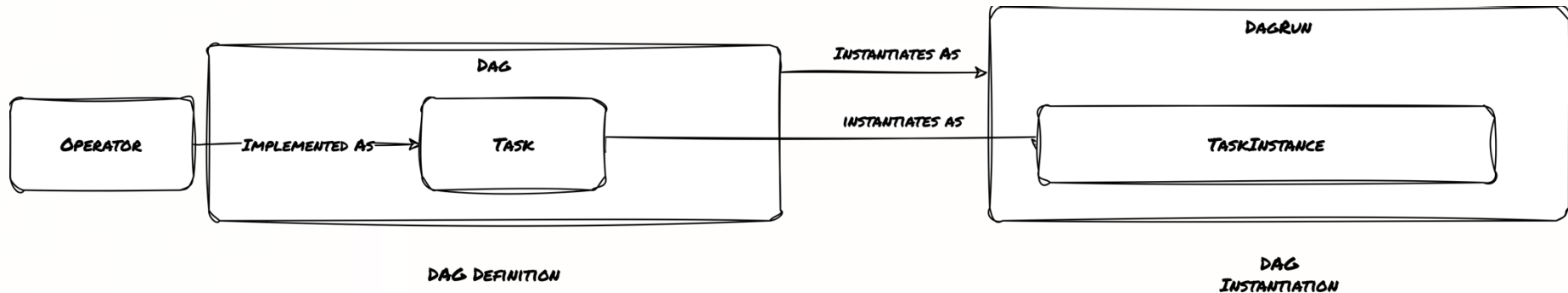
**Executor:** Component that actually runs or submits a task for execution. It runs as part of the scheduler.

**Worker:** Component that actually executes the tasks' payload. It runs an Operator's execute method.

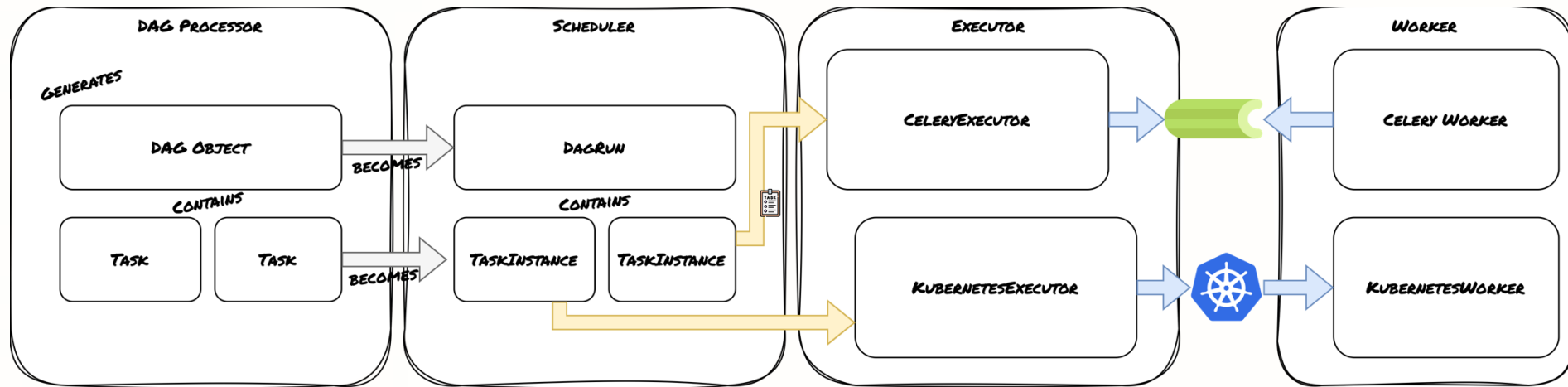
**Triggerer:** Runs and checks Triggers, which are asynchronous coroutines that monitor conditions after a task is deferred in order to resume it.

# Tl;dr: DAGs, DagRuns, Tasks, TaskInstances

- A **DAG** is a DAG.
- A **task** is the implementation of an operator. It belongs to a **DAG**.
- A DagRun is the instantiation of a DAG, at runtime.
- A **TaskInstance** is the instantiation of a task, at runtime. It belongs to a **DagRun**.

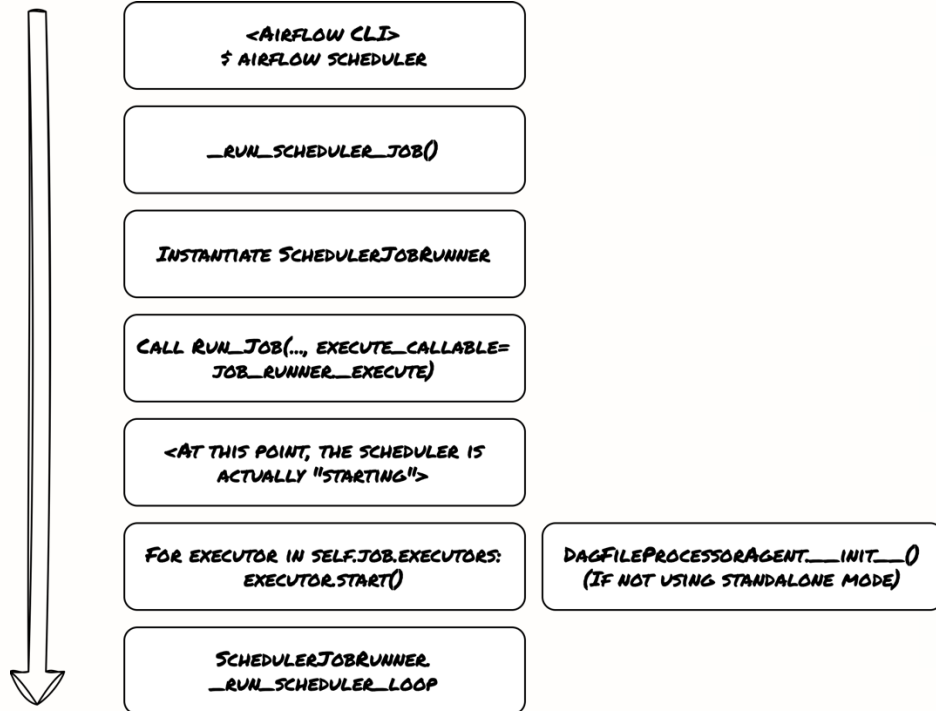


# Task scheduling and execution framework



# The scheduler initialization process

- The first step needed to schedule tasks is to start the scheduler.
- At a high level, the CLI command will invoke `_run_scheduler_job`, which will instantiate a `SchedulerJobRunner`.
- The job runner will (1) run the main scheduler loop and (2) start the executors.



# SchedulerJobRunner.\_run\_scheduler\_loop

INITIALIZE TIMERS  
(EVENTSCHEDULER)

FOR LOOP\_COUNT IN ITERTOOLS.COUNT(START=1):

SELF.DO\_SCHEDULING(...)

FOR EXECUTOR IN SELF.JOB.EXECUTORS:

EXECUTOR.HEARTBEAT(...)

PERFORM\_HEARTBEAT(JOB=SELF.JOB)

TIMERS.RUN(...)

ONLY IF NECESSARY!

RUN ANY PENDING  
TIMED EVENTS

- At this point, we're ready to start the real scheduling loop.
- The core TaskInstance and DagRun scheduling logic is in the `_do_scheduling` method, and the executor logic is in the `heartbeat` method.
- We also run maintenance operations periodically.



# The loop timers

Method	Configuration Setting	Default
<code>adopt_or_reset_orphaned_tasks</code>	<code>orphaned_tasks_check_interval</code>	300
<code>check_trigger_timeouts</code>	<code>trigger_timeout_check_interval</code>	15
<code>_emit_pool_metrics</code>	<code>pool_metrics_interval</code>	5
<code>_find_zombies</code>	<code>zombie_detection_interval</code>	10
<code>_update_dag_run_state_for_paused_dags</code>	None!	60
<code>_fail_tasks_stuck_in_queued</code>	<code>task_queued_timeout_check_interval</code>	None!
<code>_orphan_unreferenced_datasets</code>	<code>parsing_cleanup_interval</code>	None!
<code>_cleanup_stale_dags</code>	<code>parsing_cleanup_interval</code>	None!

# SchedulerJobRunner.\_do\_scheduling(...)

`SELF.CREATE_DAGRUNS_FOR_DAGS(...)`

LIMIT: `SCHEDULER.MAX_DAGRUNS_TO_CREATE_PER_LOOP` (DEFAULT: 10)  
ORDERED BY `NEXT_DAGRUN_CREATE_AFTER`

`SELF.START_QUEUED_DAGRUNS(...)`

LIMIT: `MAX_DAGRUNS_PER_LOOP_TO_SCHEDULE` (DEFAULT: 20)

`SELF.GET_NEXT_DAGRUNS_TO_EXAMINE(RUNNING)`

LIMIT: `MAX_DAGRUNS_PER_LOOP_TO_SCHEDULE` (DEFAULT: 20)

`SELF.SCHEDULE_ALL_DAG_RUNS(...)`

SCHEDULE `TASKINSTANCES` FOR `RUNNING` `DAGRUNS` AND UPDATE `NEXT_DAGRUN_CREATE_AFTER`

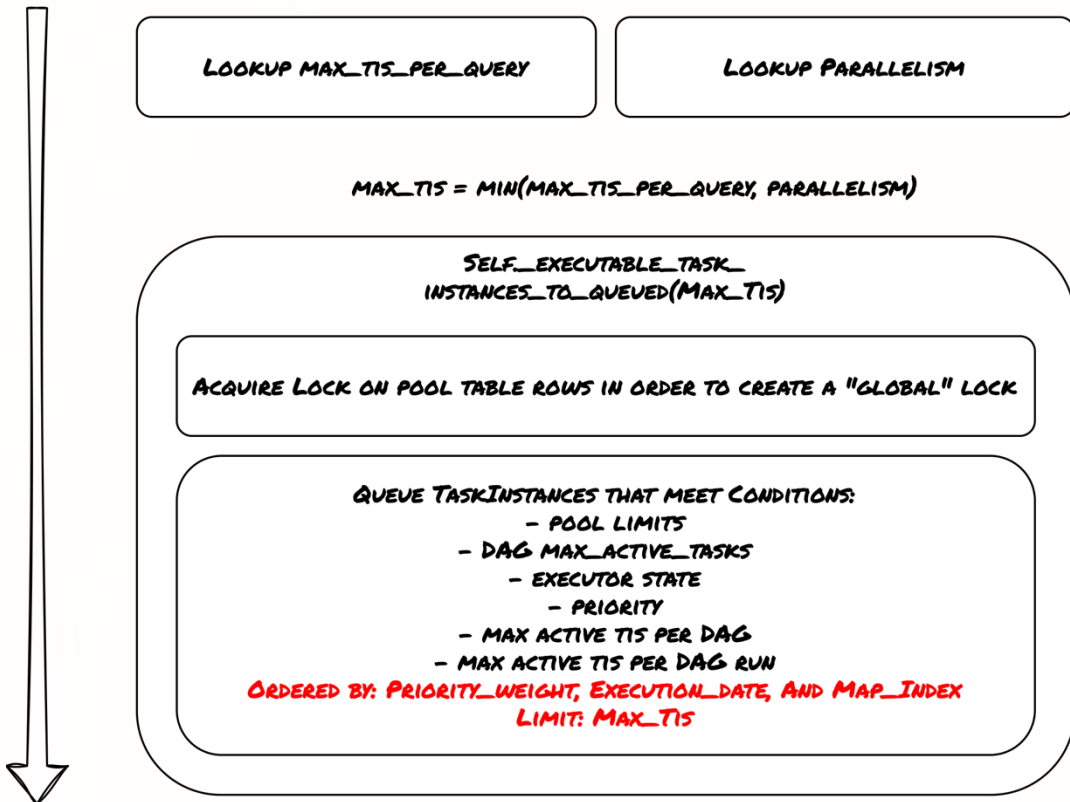
`SELF.CRITICAL_SECTION_ENQUEUE_TASK_INSTANCES(...)`

SELECTS `TASKINSTANCES` IN `SCHEDULED` STATE AND `QUEUES` THEM (WHICH EFFECTIVELY MAKES THEM VISIBLE TO THE EXECUTOR)

BULK  
FETCH FOR

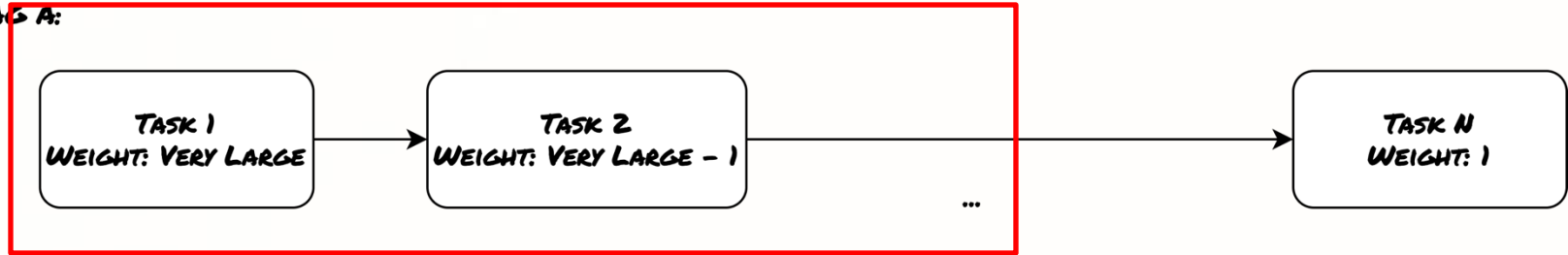


# The critical section

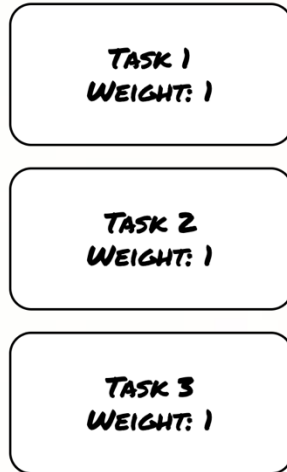


# Impact of priority weight

DAG A:

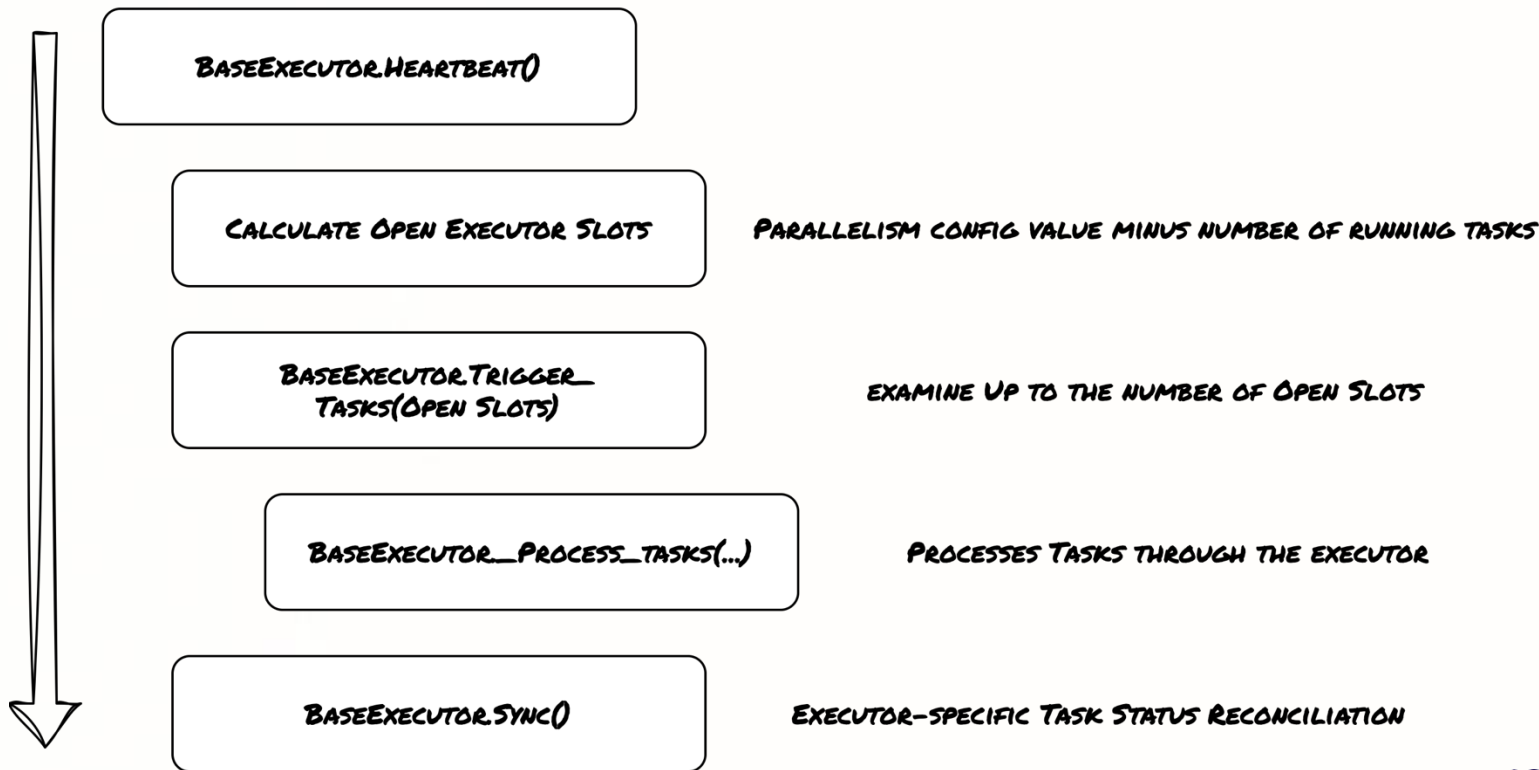


DAG B:



# The Executor Heartbeat Process

FOR EXECUTOR IN SELF.JOB.EXECUTORS:



## With CeleryExecutor (initialization)

- The CeleryExecutor initialization process is relatively simple.
- Everything is defined by the `__init__` method.
- Two important objects are initialized:
  - BulkStateFetcher
  - Tasks map
- Of note, the `start()` method in CeleryExecutor does not do anything.
- It's important to note that with CeleryExecutor, we also need to start workers!

# With CeleryExecutor

CELERYEXECUTOR\_PROCESS\_TASKS(...)

CELERYEXECUTOR OVERRIDES THE  
\_PROCESS\_TASKS METHOD

CELERYEXECUTOR\_SEND\_TASKS\_TO\_CELERY(...)

USES A PROCESS POOL EXECUTOR TO SEND TASKS TO CELERY IN PARALLEL (UP TO SYNC\_PARALLELISM)

WITH PROCESSPOOLEXECUTOR(...):

CELERYEXECUTOR\_SEND\_TASK\_TO\_EXECUTOR(...)

CELERY LIBRARY METHOD TO RUN CELERY TASKS

TASK\_TO\_RUN.APPLY\_ASYNC(...)

RETURN VALUE IS ADDED TO SELF.TASKS COLLECTION FOR STATE TRACKING

AT THIS POINT, THE TASK SHOULD BE QUEUED IN THE CELERY BACKEND

FROM SCHEDULER'S \_DO\_SCHEDULING METHOD

CELERYEXECUTOR\_SYNC()

CELERYEXECUTOR\_UPDATE\_ALL\_TASK\_STATES()

BULK\_STATE\_FETCHER.GET\_MANY()

GET STATE FOR TASKS USING THE BEST AVAILABLE METHOD (DEPENDING ON RESULT BACKEND)

FOR TASK IN <RESULT FROM ABOVE>:

IF STATE:

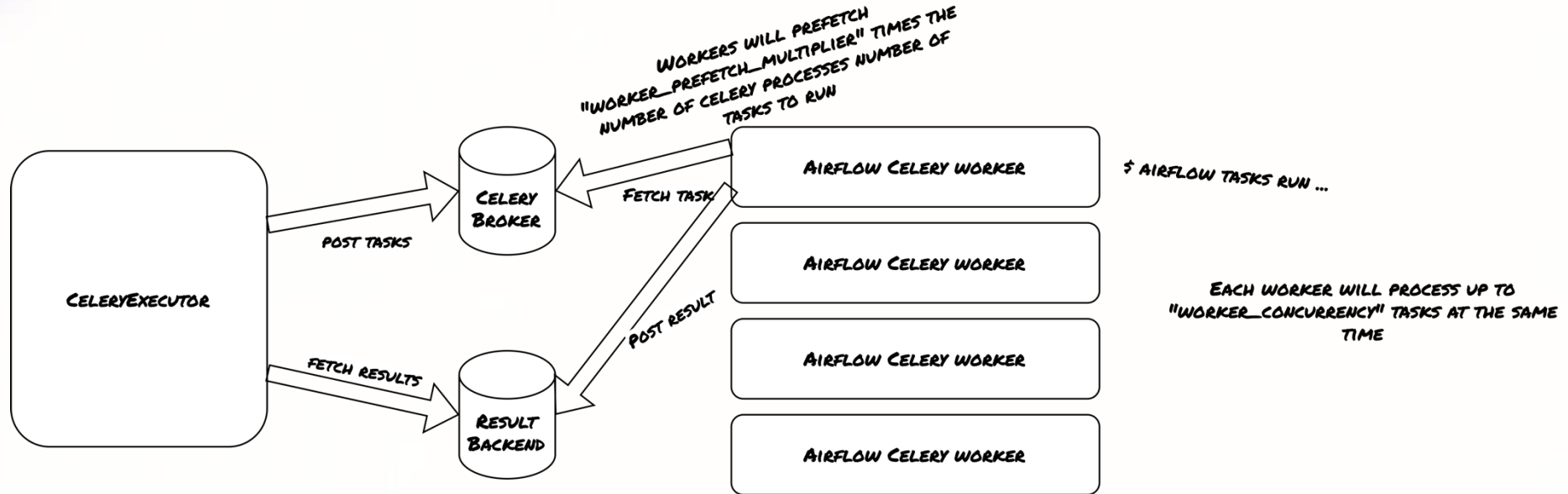
CELERYEXECUTOR\_UPDATE\_TASK\_STATE(...)

SUCCESS

FAILURE



# CeleryWorker



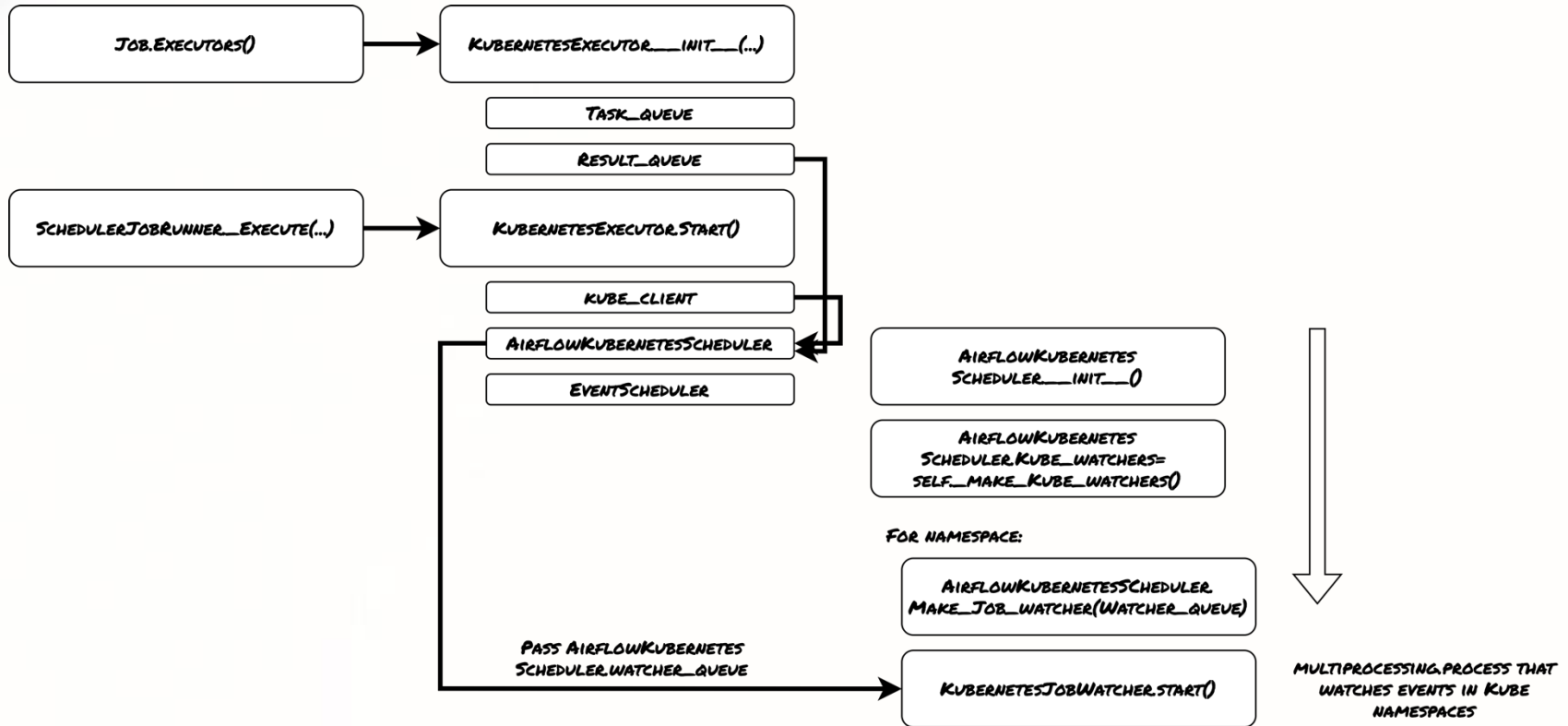


## With KubernetesExecutor (initialization)

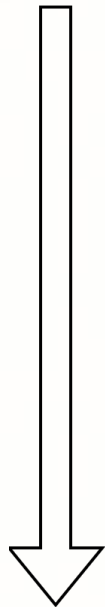
- The KubernetesExecutor initialization process on the other hand is a lot more complex. This is because it tracks a lot more state than the CeleryExecutor.
- The main subcomponents we instantiate are (1) a task queue, (2) a result queue, (3) an AirflowKubernetesScheduler, (4) a Kube client, and (5) an event scheduler.
- Since some of these components are relatively complex to instantiate, we make use of the `start()` method for the actual instantiation.



# With KubernetesExecutor (initialization)



# With KubernetesExecutor (execute\_async)



`BASEEXECUTOR.PROCESS_TASKS(...)`

GETS A TASKINSTANCE FROM  
BASEEXECUTOR.QUEUED\_TASKS

`KUBERNETESEXECUTOR.EXECUTE_ASYNC(...)`

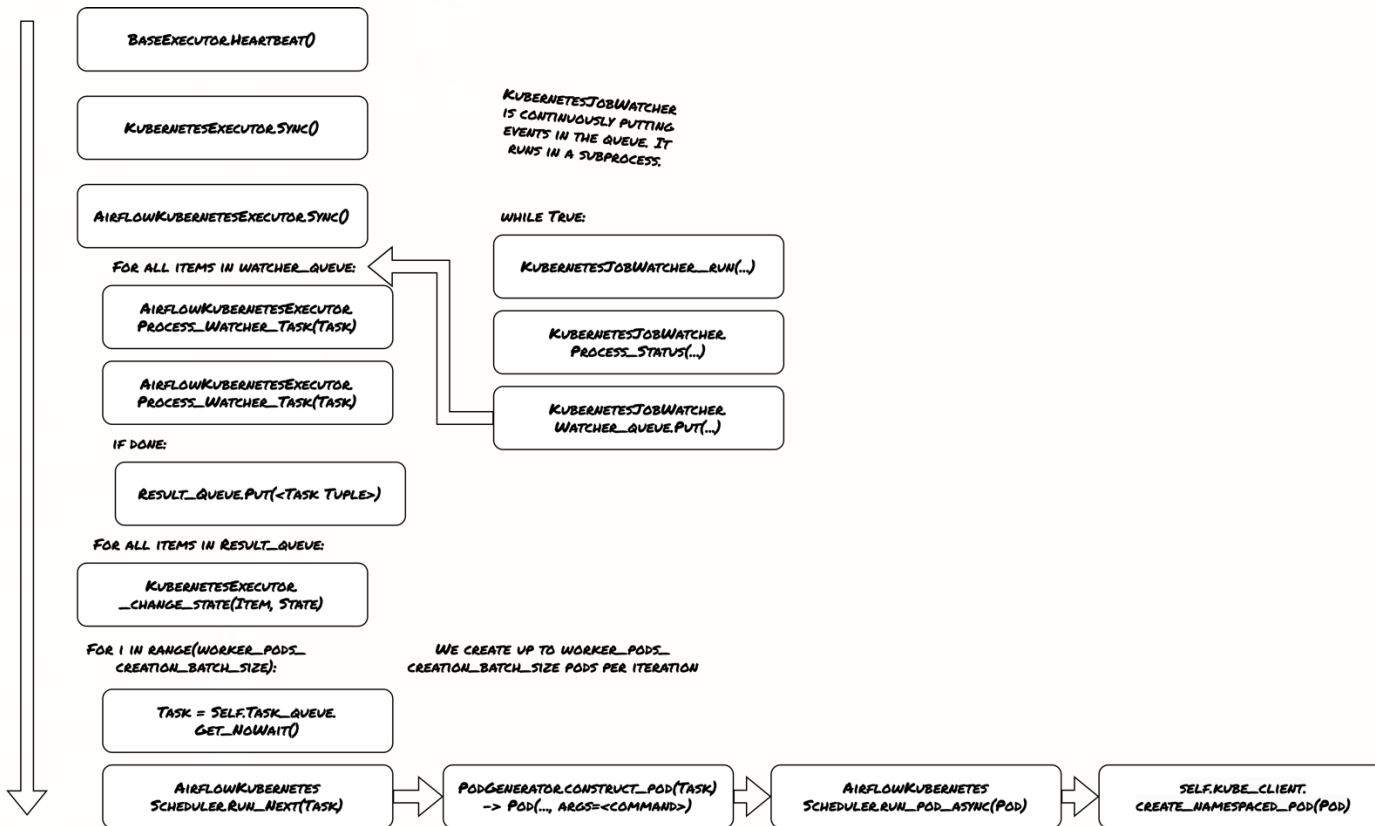
GENERATES POD SPECS AND COMMAND

`KUBERNETESEXECUTOR.TASK_QUEVE.PUT(TASK)`

PUTS (KEY, COMMAND, KUBE\_EXECUTOR\_CONFIG,  
POD\_TEMPLATE\_FILE) IN THE QUEUE.



# With KubernetesExecutor (sync method)



# Conclusions and Takeaways

- The scheduling process follows the same steps at each iteration:
  - We create DagRuns
  - Queue DagRuns
  - Queue TaskInstances
  - Create new TaskInstances
  - Run the executor
- A task always “travels” from the scheduler to the executor to the worker.
- Configuration parameters are numerous and need to be tuned carefully according to your workload patterns.

# Questions?

[www.linkedin.com/in/pfgagnon](https://www.linkedin.com/in/pfgagnon)

