```
$ cat weilee.py

__name__ = 李唯 / Wei Lee
__what_i_am_doing__ = [
    Software Engineer @ Astronomer,
    Committer @ Apache Airflow,
    First Time Speaker @ Airflow Summit
]
__github__ = Lee-W
__linkedin__ = clleew
__site__ = https://wei-lee.me
```

```
$ python weilee.py



              File "weilee.py", line 1
                __name__ = 李唯 / Wei Lee
                                          ^^^

              SyntaxError: invalid syntax
```

# QR Code links to this slide deck

# Let's start with how a typical task works now

Define a DAG

```python
from __future__ import annotations

import pendulum

from airflow import DAG
from airflow.operators.bash import BashOperator

with DAG(
    dag_id="example_dag",
    start_date=pendulum.datetime(2021, 1, 1, tz="UTC"),
    schedule=None,
    catchup=False,
):
    bash_task = BashOperator(task_id="bash_task", bash_command="echo example")
```

# Let's start with how a typical task works now

Define a DAG

```python
from __future__ import annotations

import pendulum

from airflow import DAG
from airflow.operators.bash import BashOperator

with DAG(
    dag_id="example_dag",
    start_date=pendulum.datetime(2021, 1, 1, tz="UTC"),
    schedule=None,
    catchup=False,
):
    bash_task = BashOperator(task_id="bash_task", bash_command="echo example")
```

# Let's start with how a typical task works now

## Define a DAG

```python
from __future__ import annotations

import pendulum

from airflow import DAG
from airflow.operators.bash import BashOperator

with DAG(
    dag_id="example_dag",
    start_date=pendulum.datetime(2021, 1, 1, tz="UTC"),
    schedule=None,
    catchup=False,
):
    bash_task = BashOperator(task_id="bash_task", bash_command="echo example")
```

# Let's start with how a typical task works now

Define a DAG

```python
1   from __future__ import annotations
1
2   import pendulum
3
4   from airflow import DAG
5   from airflow.operators.bash import BashOperator
6
7   with DAG(
8       dag_id="example_dag",
9       start_date=pendulum.datetime(2021, 1, 1, tz="UTC"),
10      schedule=None,
11      catchup=False,
12  ):
13      bash_task = BashOperator(task_id="bash_task", bash_command="echo example")
```

⬅ This page is for traditional operators.

# Let's start with how a typical task works now

Under the hood, it runs "execute".

```python
class BashOperator(BaseOperator):
    r"""

    template_fields: Sequence[str] = ("bash_command", "env", "cwd")
    template_fields_renderers = {"bash_command": "bash", "env": "json"}
    template_ext: Sequence[str] = (".sh", ".bash")
    ui_color = "#f0ede4"

    def __init__(

    @cached_property
    def subprocess_hook(self):

    @staticmethod
    def refresh_bash_command(ti: TaskInstance) -> None:

    def get_env(self, context):

    def execute(self, context: Context):

    def on_kill(self) -> None:
```

ASTRONOMER

# Let's start with how a typical task works now

🎉

# Let's start with how a typical task works now

# Since Airflow 2.2

deferrable operator was introduced

## Deferrable Tasks (AIP-40)

Deferrable tasks allows operators or sensors to defer themselves until a light-weight async check passes, at which point they can resume executing. Most importantly, this results in the worker slot, and most notably any resources used by it, to be returned to Airflow. This allows simple things like monitoring a job in an external system or watching for an event to be much cheaper.

To support this feature, a new component has been added to Airflow, the triggerer, which is the daemon process that runs the asyncio event loop.

Airflow 2.2.0 ships with 2 deferrable sensors, `DateTimeSensorAsync` and `TimeDeltaSensorAsync`, both of which are drop-in replacements for the existing corresponding sensor.

More information can be found at:

Deferrable Operators & Triggers

# But why?

## Non-Deferrable Operator

| Submit Job to Spark Cluster | Poll Spark Cluster for Job Status | Receive Terminal Status for Job on Spark Cluster |
|---|---|---|

**Worker Slot Allocated**

from: https://www.astronomer.io/docs/learn/deferrable-operators

ASTRONOMER

# But why?

## Non-Deferrable Operator

| Submit Job to Spark Cluster | Poll Spark Cluster for Job Status | Receive Terminal Status for Job on Spark Cluster |

**Worker Slot Allocated**

from: https://www.astronomer.io/docs/learn/deferrable-operators

# But why?

Non-Deferrable Operator

| Submit Job to Spark Cluster | Poll Spark Cluster for Job Status | Receive Terminal Status for Job on Spark Cluster |

**Worker Slot Allocated**

from: https://www.astronomer.io/docs/learn/deferrable-operators

ASTRONOMER

# But why?

## Non-Deferrable Operator

| Submit Job to Spark Cluster | Poll Spark Cluster for Job Status | Receive Terminal Status for Job on Spark Cluster |
|---|---|---|

**Worker Slot Allocated**

from: https://www.astronomer.io/docs/learn/deferrable-operators

ASTRONOMER

# But why? → Release worker slots

Deferrable Operator

| Submit Job to Spark Cluster | Poll Spark Cluster for Job Status | | Receive Terminal Status for Job on Spark Cluster |
|---|---|---|---|
| Worker Slot Allocated | Triggerer Process | | Worker Slot Allocated |

🔳 from: https://www.astronomer.io/docs/learn/deferrable-operators

# But why? → Release worker slots

Deferrable Operator

| Submit Job to Spark Cluster | Poll Spark Cluster for Job Status | Receive Terminal Status for Job on Spark Cluster |
|---|---|---|
| **Worker Slot Allocated** | **Triggerer Process** | **Worker Slot Allocated** |

This page is for deferrable operators.

⏳

from: https://www.astronomer.io/docs/learn/deferrable-operators
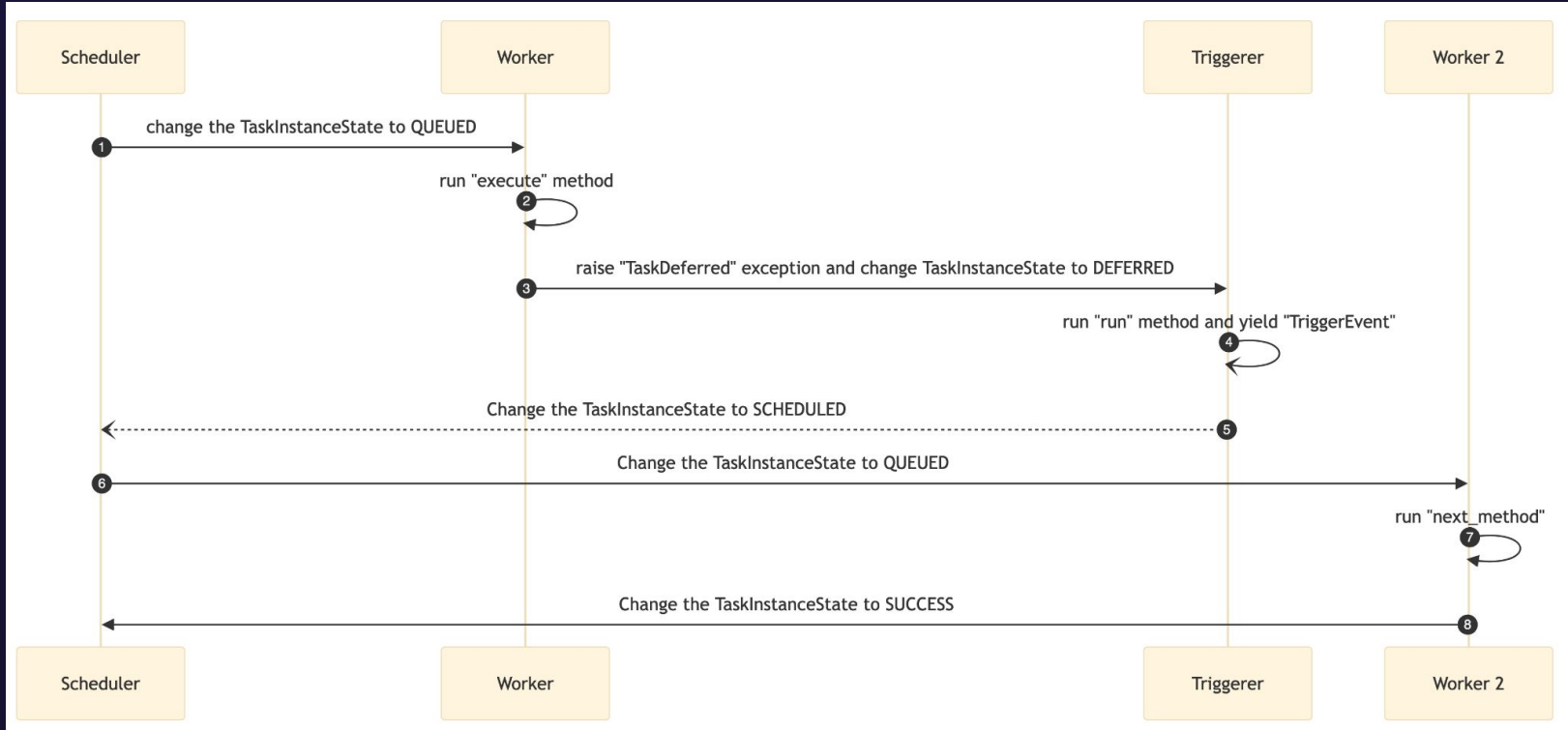
ASTRONOMER

# Release worker slots. And…?

Reduce resource usage

# How does deferrable operators work?

# How does deferrable operator work?

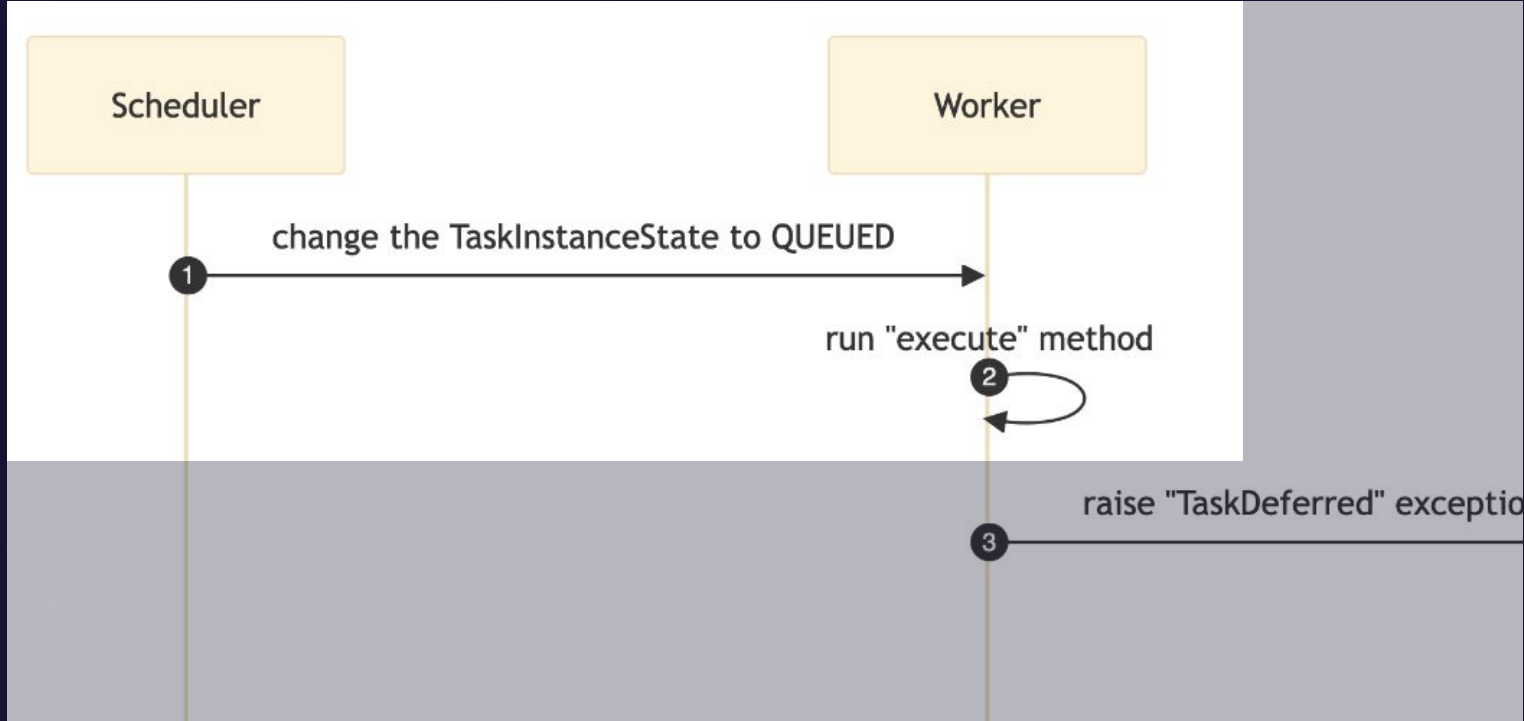Well… we still run "execute" first.

```
s3.py                                                    ×
85    class S3KeysUnchangedSensor(BaseSensorOperator):
6         def execute(self, context: Context) -> None:
5             """Airflow runs this method on the worker and defers using the trigger
4             if not self.deferrable:
3                 super().execute(context)
2             else:
1                 if not self.poke(context):
388                   self.defer(
1                         timeout=timedelta(seconds=self.timeout),
2                         trigger=S3KeysUnchangedTrigger(
3                             bucket_name=self.bucket_name,
4                             prefix=self.prefix,
5                             inactivity_period=self.inactivity_period,
6                             min_objects=self.min_objects,
7                             previous_objects=self.previous_objects,
8                             inactivity_seconds=self.inactivity_seconds,
```
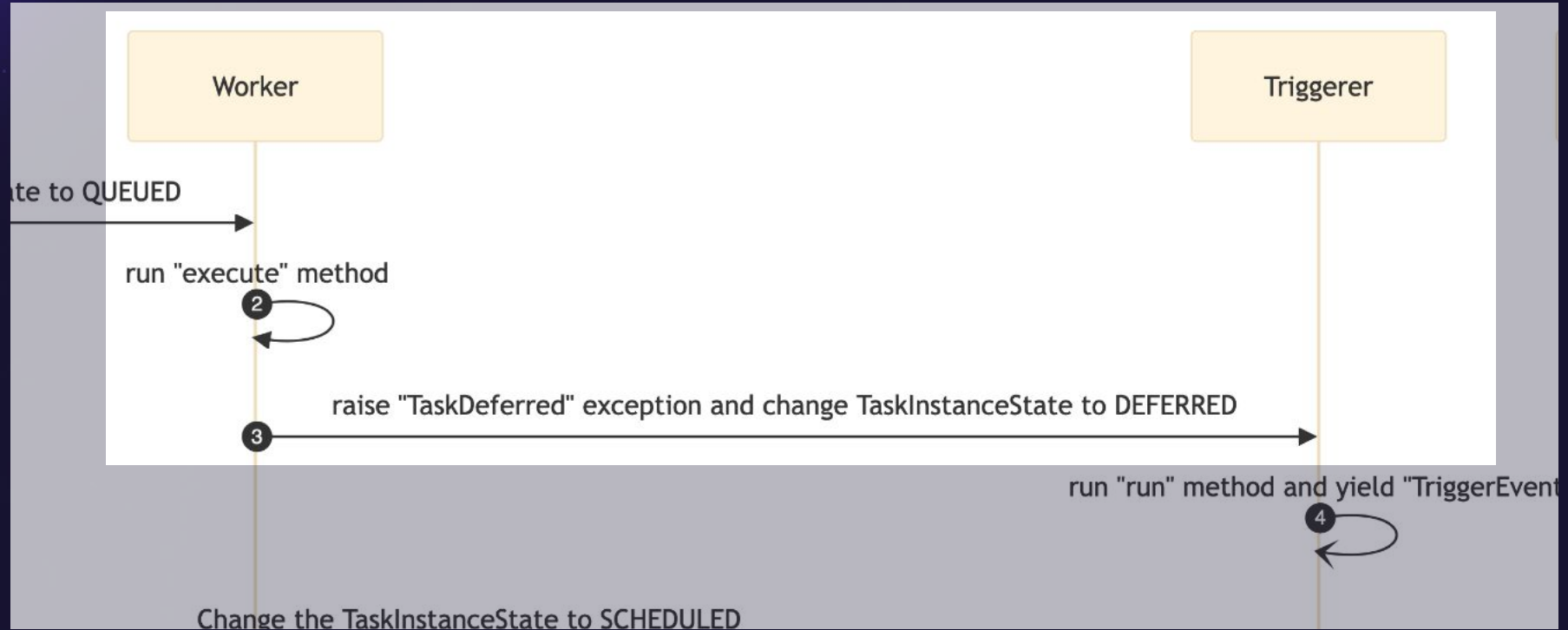
ASTRONOMER

# How does deferrable operators work?

# How does deferrable operator work?

The main difference is executing "self.defer" and
raise a TaskDeferred exception through it

```python
class S3KeysUnchangedSensor(BaseSensorOperator):
    def execute(self, context: Context) -> None:
        """Airflow runs this method on the worker and defers using the trigger
        if not self.deferrable:
            super().execute(context)
        else:
            if not self.poke(context):
                self.defer(
                    timeout=timedelta(seconds=self.timeout),
                    trigger=S3KeysUnchangedTrigger(
                        bucket_name=self.bucket_name,
                        prefix=self.prefix,
                        inactivity_period=self.inactivity_period,
                        min_objects=self.min_objects,
                        previous_objects=self.previous_objects,
                        inactivity_seconds=self.inactivity_seconds,
```

# How does deferrable operator work?

Then, it's the triggerer's turn to run.

```
85    class S3KeysUnchangedSensor(BaseSensorOperator):
 6        def execute(self, context: Context) -> None:
 5            """Airflow runs this method on the worker and defers using the trigger
 4            if not self.deferrable:
 3                super().execute(context)
 2            else:
 1                if not self.poke(context):
388                    self.defer(
 1                        timeout=timedelta(seconds=self.timeout),
 2                        trigger=S3KeysUnchangedTrigger(
 3                            bucket_name=self.bucket_name,
 4                            prefix=self.prefix,
 5                            inactivity_period=self.inactivity_period,
 6                            min_objects=self.min_objects,
 7                            previous_objects=self.previous_objects,
 8                            inactivity_seconds=self.inactivity_seconds,
```

s3.py

# How does deferrable operators work?

# How does deferrable operator work?

Execute the async "run" method in the triggerer

```python
# s3.py                    ×
14   class S3KeysUnchangedTrigger(BaseTrigger):
4        async def run(self) -> AsyncIterator[TriggerEvent]:
2            try:
1                async with self.hook.async_conn as client:
194                  while True:
1                        result = await self.hook.is_keys_unchanged_async(
2                            client=client,
3                            bucket_name=self.bucket_name,
4                            prefix=self.prefix,
5                            inactivity_period=self.inactivity_period,
6                            min_objects=self.min_objects,
7                            previous_objects=self.previous_objects,
8                            inactivity_seconds=self.inactivity_seconds,
9                            allow_delete=self.allow_delete,
10                           last_activity_time=self.last_activity_time,
11                       )
12                       if result.get("status") in ("success", "error"):
13                           yield TriggerEvent(result)
```

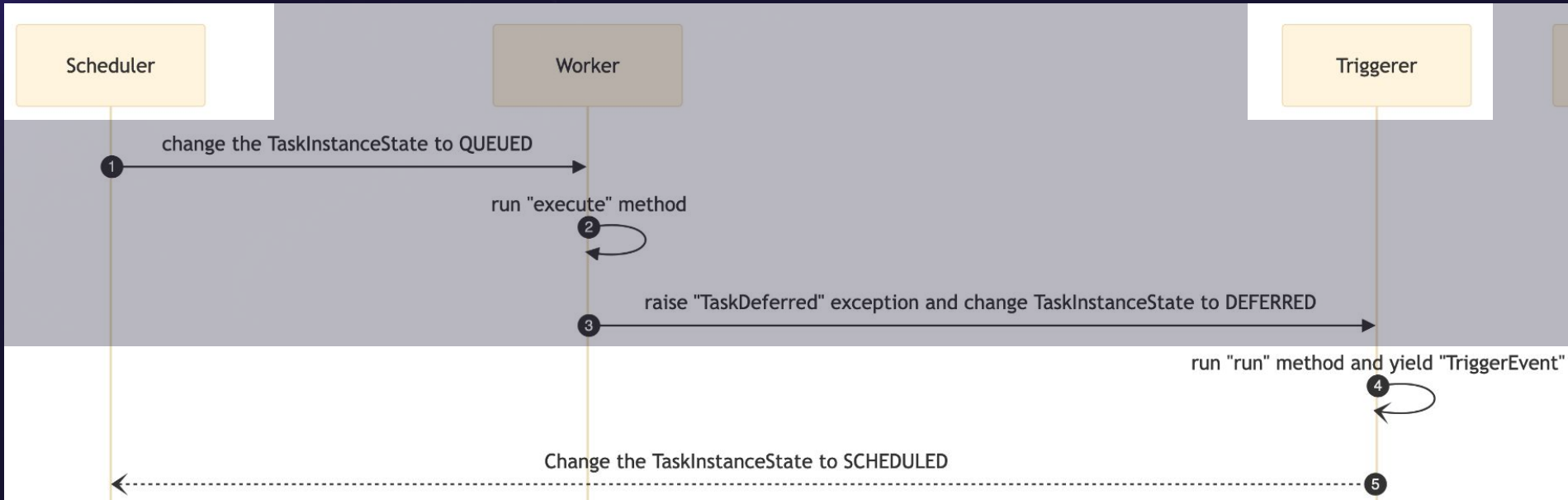# How does deferrable operator work?

yield a TriggerEvent when it finished

```
s3.py                                ×
14   class S3KeysUnchangedTrigger(BaseTrigger):
 4       async def run(self) -> AsyncIterator[TriggerEvent]:
 2           try:
 1               async with self.hook.async_conn as client:
194                 while True:
 1                       result = await self.hook.is_keys_unchanged_async(
 2                           client=client,
 3                           bucket_name=self.bucket_name,
 4                           prefix=self.prefix,
 5                           inactivity_period=self.inactivity_period,
 6                           min_objects=self.min_objects,
 7                           previous_objects=self.previous_objects,
 8                           inactivity_seconds=self.inactivity_seconds,
 9                           allow_delete=self.allow_delete,
10                           last_activity_time=self.last_activity_time,
11                       )
12                       if result.get("status") in ("success", "error"):
13                           yield TriggerEvent(result)
```

# How does deferrable operators work?
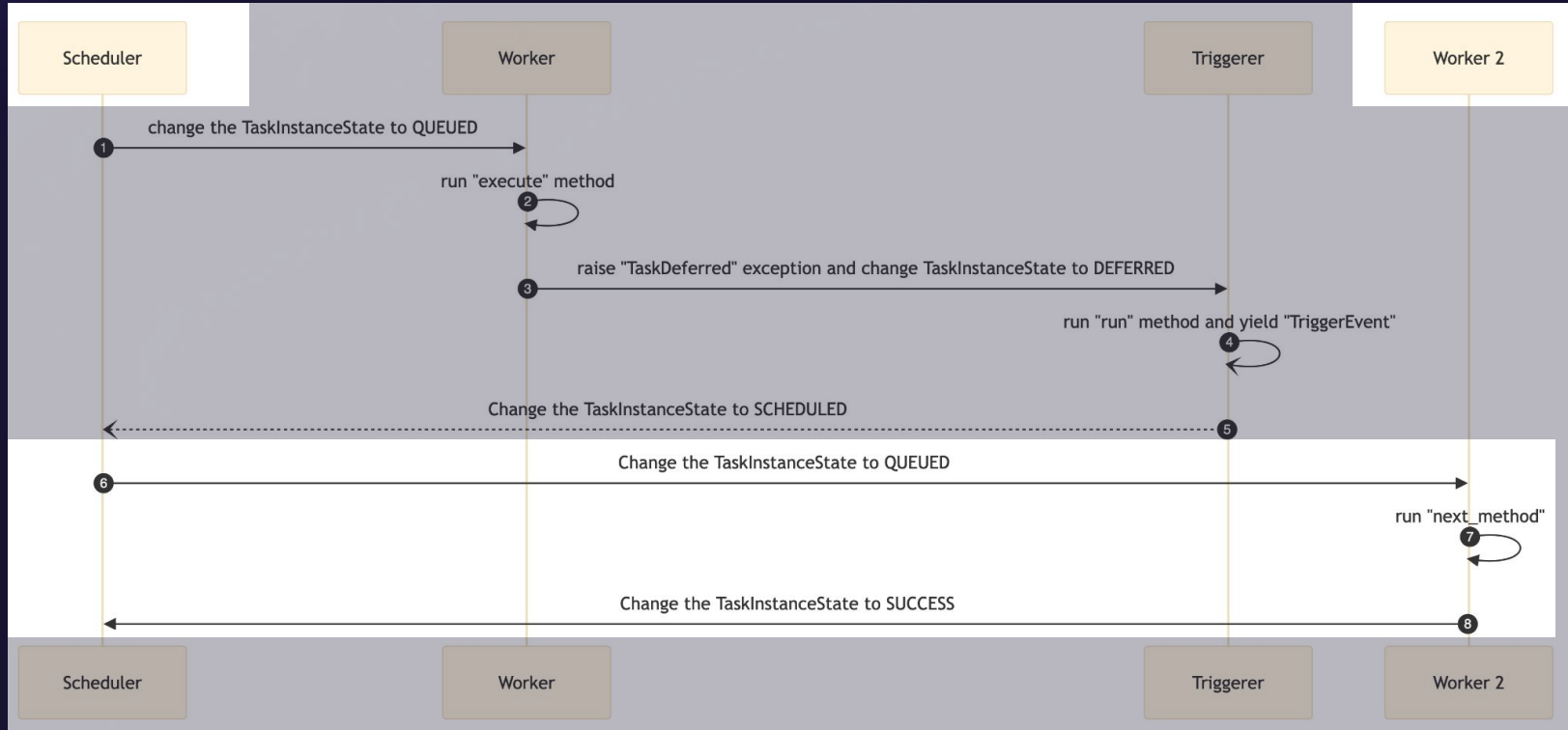
# How does deferrable operator work?

Execute the "execute_complete" method in another worker

```python
s3.py                                                                    ×
151    class S3KeysUnchangedSensor(BaseSensorOperator):
14         return self.is_keys_unchanged(set(self.hook.list_keys(self.bucket_name, prefix=self.prefix)))
13
12         def execute(self, context: Context) -> None:
11
10         def execute_complete(self, context: Context, event: dict[str, Any] | None = None) -> None:
9              """
8              Execute when the trigger fires - returns immediately.
7
6              Relies on trigger to throw an exception, otherwise it assumes execution was successful.
5              """
4              event = validate_execute_complete_event(event)
3
2              if event and event["status"] == "error":
1                  raise AirflowException(event["message"])
407            return None
```

ASTRONOMER

# How does deferrable operators work?



Scheduler | Worker | Triggerer | Worker 2

1. change the TaskInstanceState to QUEUED

2. run "execute" method

3. raise "TaskDeferred" exception and change TaskInstanceState to DEFERRED

4. run "run" method and yield "TriggerEvent"

5. Change the TaskInstanceState to SCHEDULED

6. Change the TaskInstanceState to QUEUED

7. run "next_method"

8. Change the TaskInstanceState to SUCCESS

ASTRONOMER

# Do we really need to run it in the worker first?

the only logic before deferring

```
149   class S3KeysUnchangedSensor(BaseSensorOperator):
  5       def execute(self, context: Context) -> None:
  3           if not self.deferrable:
  2               super().execute(context)
  1           else:
351             if not self.poke(context):
  1               self.defer(
```

ASTRONOMER

I DON'T THINK SO, SIR.
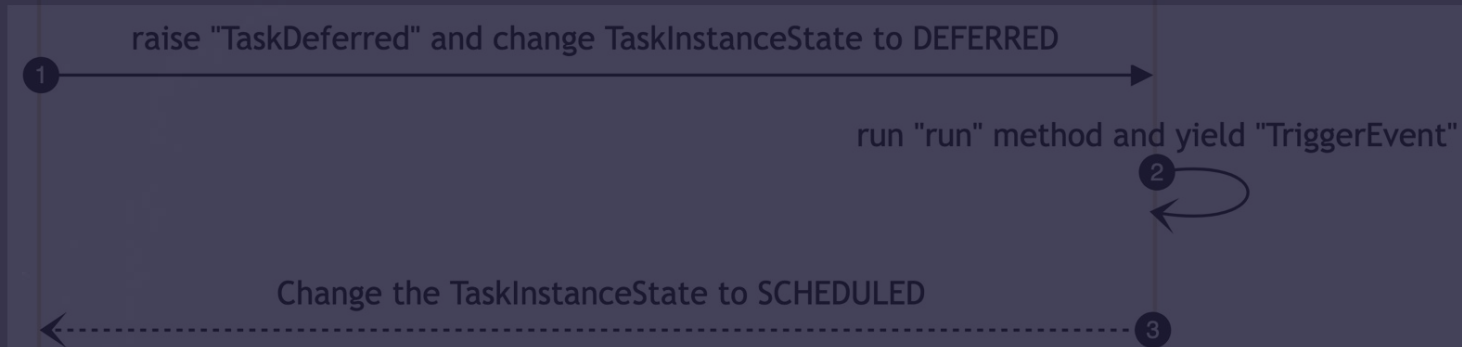
# What can we do?

# What can we do?

Start task execution in the triggerer

# What can we do?

Start task execution in the triggerer



This page is for new features.

# Start task execution in the **worker**

```python
from datetime import timedelta
from typing import Any

from airflow.sensors.base import BaseSensorOperator
from airflow.triggers.temporal import TimeDeltaTrigger
from airflow.utils.context import Context


class WaitOneHourSensor(BaseSensorOperator):
    def execute(self, context: Context) -> None:
        self.defer(
            trigger=TimeDeltaTrigger(timedelta(hours=1)), method_name="execute_complete"
        )

    def execute_complete(
        self, context: Context, event: dict[str, Any] | None = None
    ) -> None:
        # We have no more work to do here. Mark as complete.
        return
```

# Start task execution in the triggerer

```python
class WaitHoursSensor(BaseSensorOperator):
    start_trigger_args = StartTriggerArgs(
        trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
        trigger_kwargs={"delta": timedelta(hours=1)},
        next_method="execute_complete",
        next_kwargs=None,
        timeout=None,
    )

    def __init__(
        self,
        *args: list[Any],
        trigger_kwargs: dict[str, Any] | None,
        start_from_trigger: bool,
        **kwargs: dict[str, Any],
    ) -> None:
        super().__init__(*args, **kwargs)
        self.start_trigger_args.trigger_kwargs = trigger_kwargs
        self.start_from_trigger = start_from_trigger

    def execute_complete(
        self, context: Context, event: dict[str, Any] | None = None
    ) -> None:
        # We have no more work to do here. Mark as complete.
        return
```

# Start task execution in the triggerer

StartTriggerArgs and start_from_trigger

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)


def __init__(
    self,
    *args: list[Any],
    trigger_kwargs: dict[str, Any] | None,
    start_from_trigger: bool,
    **kwargs: dict[str, Any],
) -> None:
    super().__init__(*args, **kwargs)
    self.start_trigger_args.trigger_kwargs = trigger_kwargs
    self.start_from_trigger = start_from_trigger
```

ASTRONOMER

# Start task execution in the triggerer

StartTriggerArgs and start_from_trigger

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)
```

# Start task execution in the triggerer

StartTriggerArgs vs self.defer

⌛

```python
18  class WaitHoursSensor(BaseSensorOperator):
17      start_trigger_args = StartTriggerArgs(
16          trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
15          trigger_kwargs={"delta": timedelta(hours=1)},
14          next_method="execute_complete",
13          next_kwargs=None,
12          timeout=None,
11      )
10
 9      def __init__(
 8          self,
 7          *args: list[Any],
 6          trigger_kwargs: dict[str, Any] | None,
 5          start_from_trigger: bool,
 4          **kwargs: dict[str, Any],
 3      ) -> None:
 2          super().__init__(*args, **kwargs)
 1          self.start_trigger_args.trigger_kwargs = trigger_kwargs
26          self.start_from_trigger = start_from_trigger
 1
 2      def execute_complete(
 3          self, context: Context, event: dict[str, Any] | None = None
 4      ) -> None:
 5          # We have no more work to do here. Mark as complete.
 6          return
```

```python
 7  from datetime import timedelta
 6  from typing import Any
 5
 4  from airflow.sensors.base import BaseSensorOperator
 3  from airflow.triggers.temporal import TimeDeltaTrigger
 2  from airflow.utils.context import Context
 1
 8
 1  class WaitOneHourSensor(BaseSensorOperator):
 2      def execute(self, context: Context) -> None:
 3          self.defer(
 4              trigger=TimeDeltaTrigger(timedelta(hours=1)),
 5              method_name="execute_complete"
 6          )
 7
 8      def execute_complete(
 9          self, context: Context, event: dict[str, Any] | None = None
10      ) -> None:
11          # We have no more work to do here. Mark as complete.
12          return
```

# Start task execution in the triggerer

StartTriggerArgs vs self.defer

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)
```

```python
def execute(self, context: Context) -> None:
    self.defer(
        trigger=TimeDeltaTrigger(timedelta(hours=1)),
        method_name="execute_complete"
    )
```

ASTRONOMER

# Start task execution in the triggerer

trigger_cls

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)
```

```python
def execute(self, context: Context) -> None:
    self.defer(
        trigger=TimeDeltaTrigger(timedelta(hours=1)),
        method_name="execute_complete"
    )
```

⌛

# Start task execution in the triggerer

trigger_kwargs

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)

def execute(self, context: Context) -> None:
    self.defer(
        trigger=TimeDeltaTrigger(timedelta(hours=1)),
        method_name="execute_complete"
    )
```

⏳

ASTRONOMER

# Start task execution in the triggerer

next_method

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)
```

```python
def execute(self, context: Context) -> None:
    self.defer(
        trigger=TimeDeltaTrigger(timedelta(hours=1)),
        method_name="execute_complete"
    )
```
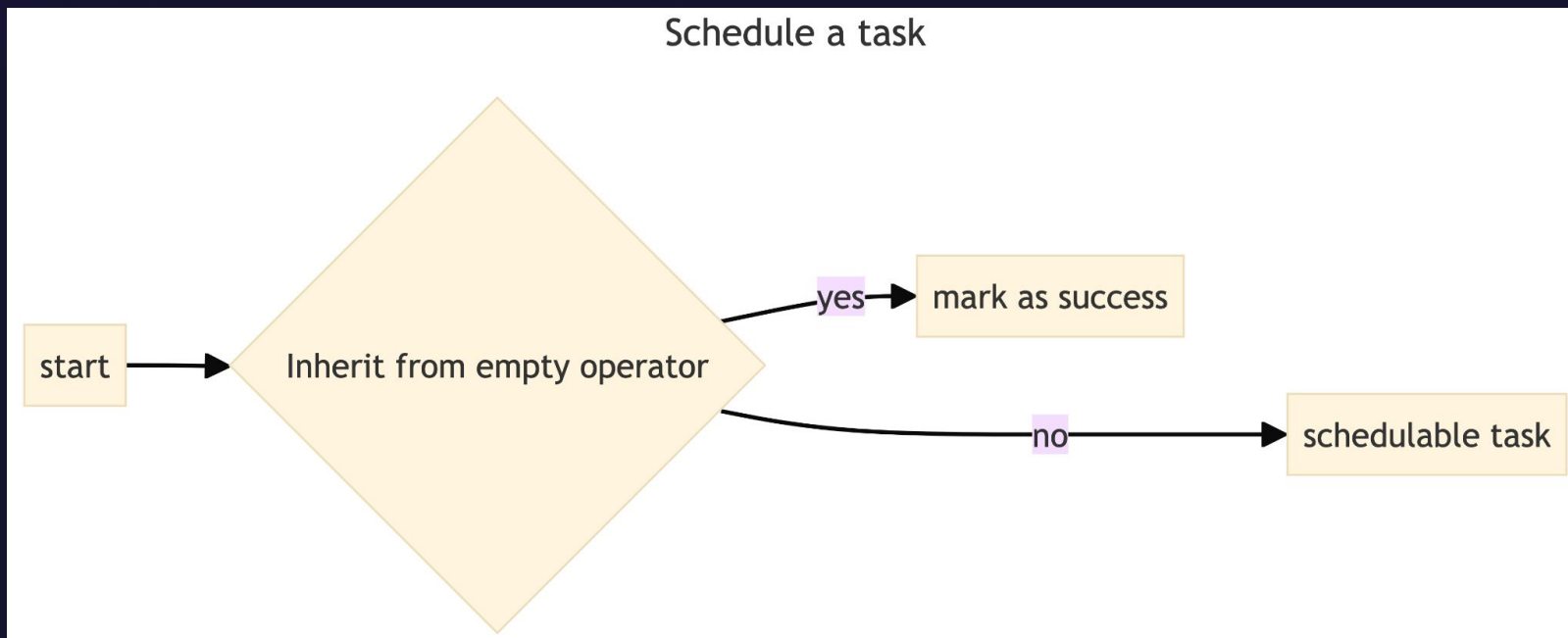
# Start task execution in the triggerer

args with default values

```python
start_trigger_args = StartTriggerArgs(
    trigger_cls="airflow.triggers.temporal.TimeDeltaTrigger",
    trigger_kwargs={"delta": timedelta(hours=1)},
    next_method="execute_complete",
    next_kwargs=None,
    timeout=None,
)
```

```python
def execute(self, context: Context) -> None:
    self.defer(
        trigger=TimeDeltaTrigger(timedelta(hours=1)),
        method_name="execute_complete"
    )
```

# Start task execution in the triggerer

## Under the hood

```python
  ✦ dagrun.py          ×
38      def schedule_tis(
37          self,
36          schedulable_tis: Iterable[TI],
18          for ti in schedulable_tis:
17              if TYPE_CHECKING:
16                  assert ti.task
15              if (
14                  ti.task.inherits_from_empty_operator
13                  and not ti.task.on_execute_callback
12                  and not ti.task.on_success_callback
11                  and not ti.task.outlets
10              ):
 9                  dummy_ti_ids.append((ti.task_id, ti.map_index))
 8              # check "start_trigger_args" to see whether the operator supports start execution from triggerer
 7              # if so, we'll then check "start_from_trigger" to see whether this feature is turned on and defer
 6              # this task.
 5              # if not, we'll add this "ti" into "schedulable_ti_ids" and later execute it to run in the worker
 4              elif ti.task.start_trigger_args is not None:
 3                  context = ti.get_template_context()
 2                  start_from_trigger = ti.task.expand_start_from_trigger(context=context, session=session)
 1
1588                  if start_from_trigger:
 1                      ti.start_date = timezone.utcnow()
 2                      if ti.state != TaskInstanceState.UP_FOR_RESCHEDULE:
 3                          ti.try_number += 1
 4                      ti.defer_task(exception=None, session=session)
 5                  else:
 6                      schedulable_ti_ids.append((ti.task_id, ti.map_index))
 7              else:
 8                  schedulable_ti_ids.append((ti.task_id, ti.map_index))
```
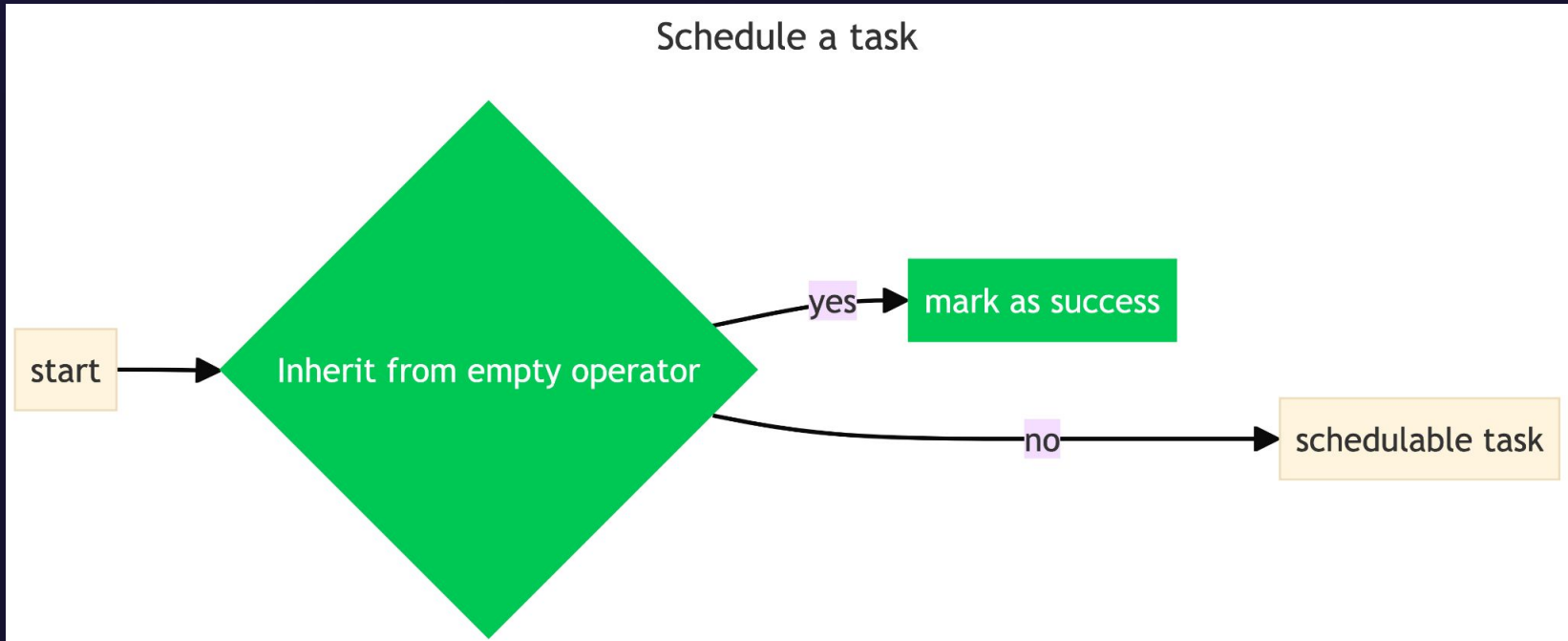
# Start task execution in the triggerer

Under the hood (it used to be…)

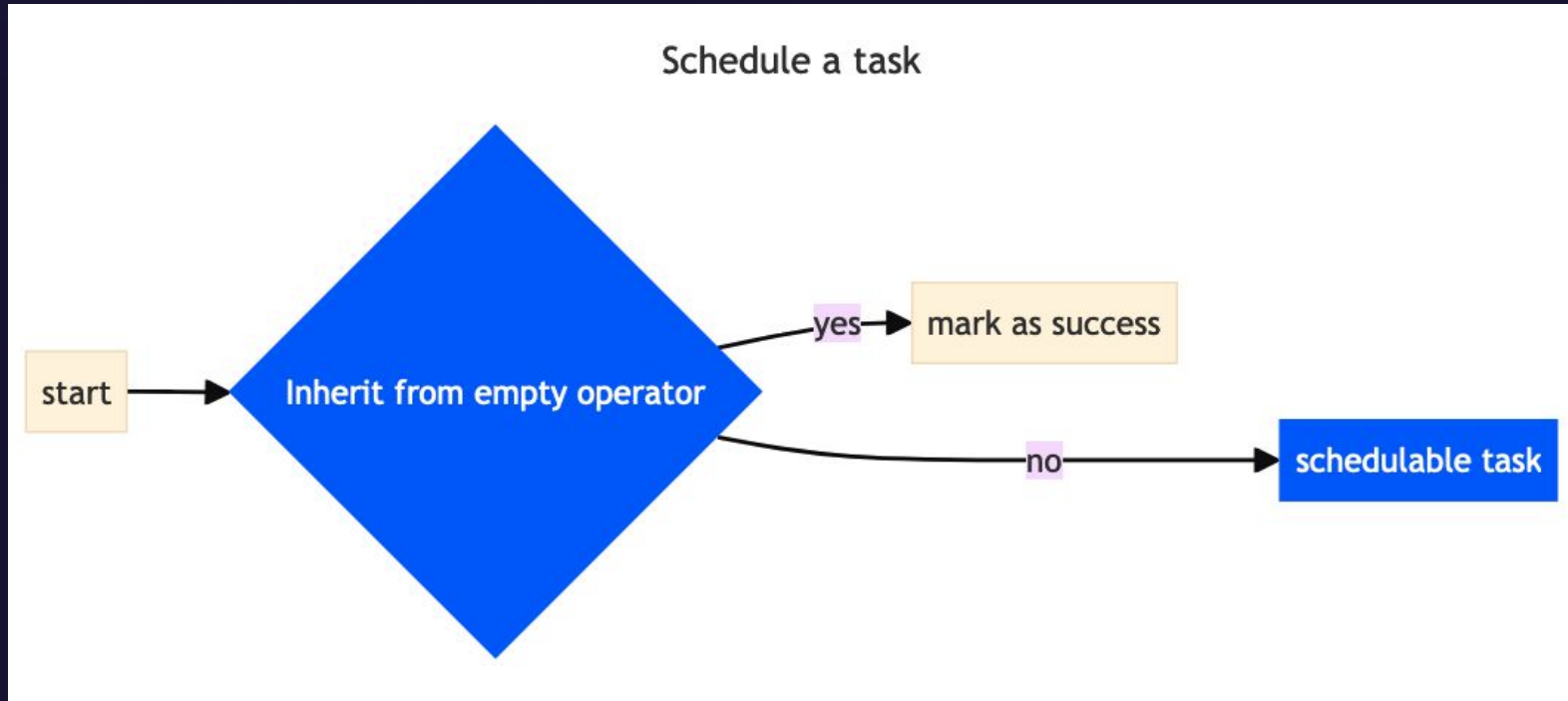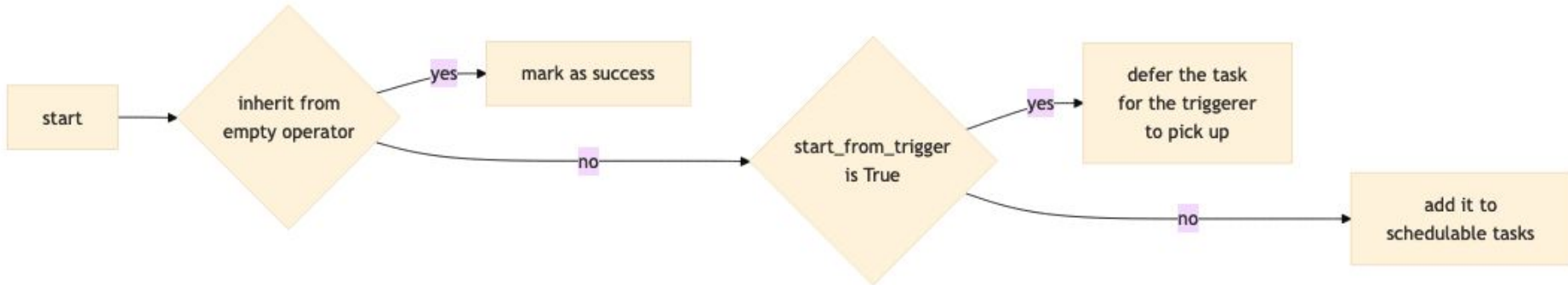# Start task execution in the triggerer

Under the hood (it used to be…)


Schedule a task

# Start task execution in the triggerer

Under the hood (it used to be…)

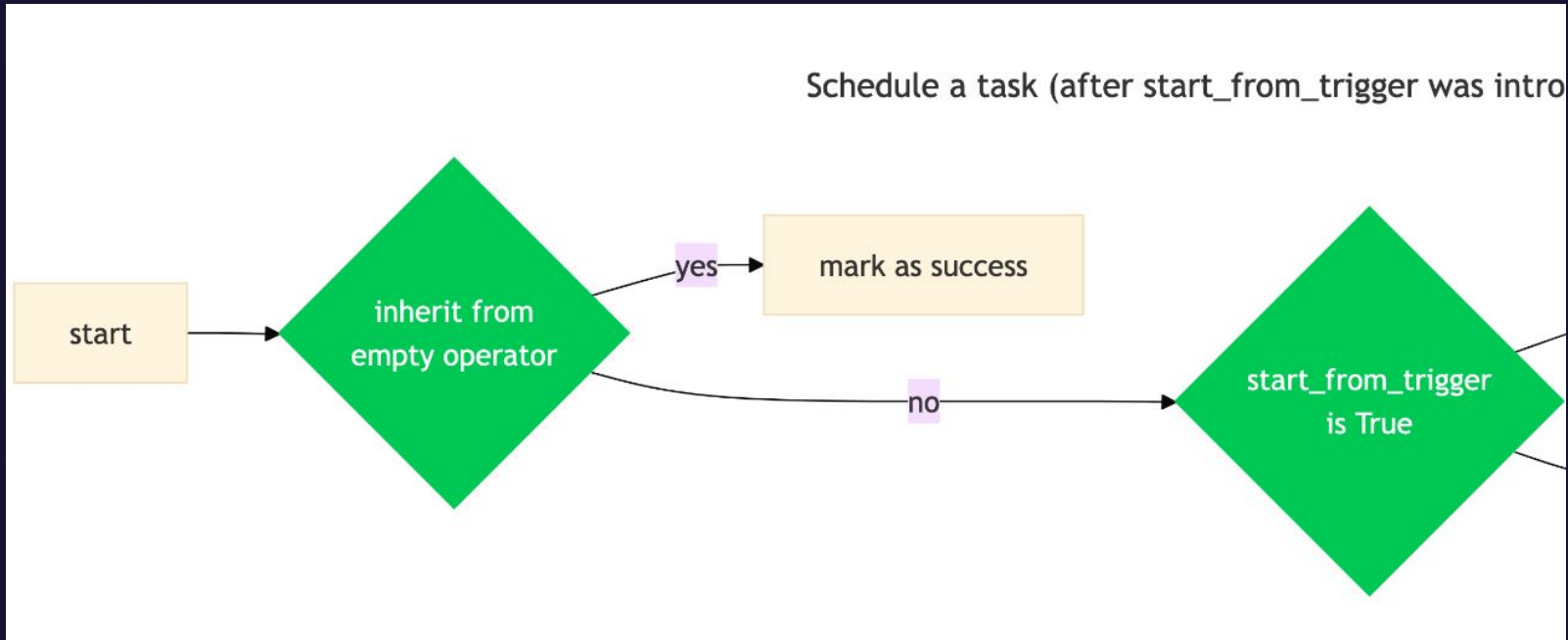# Start task execution in the triggerer

## Under the hood (it's now…)



Schedule a task (after start_from_trigger was introduced)

start → inherit from empty operator → yes → mark as success
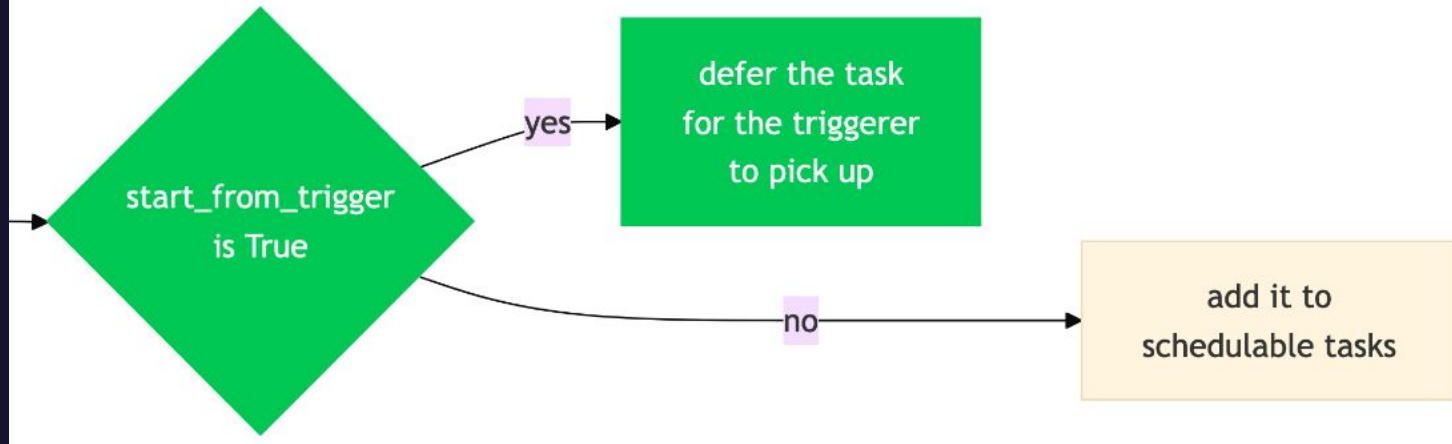
inherit from empty operator → no → start_from_trigger is True → yes → defer the task for the triggerer to pick up

start_from_trigger is True → no → add it to schedulable tasks

ASTRONOMER

# Start task execution in the triggerer

Under the hood (it's now…)

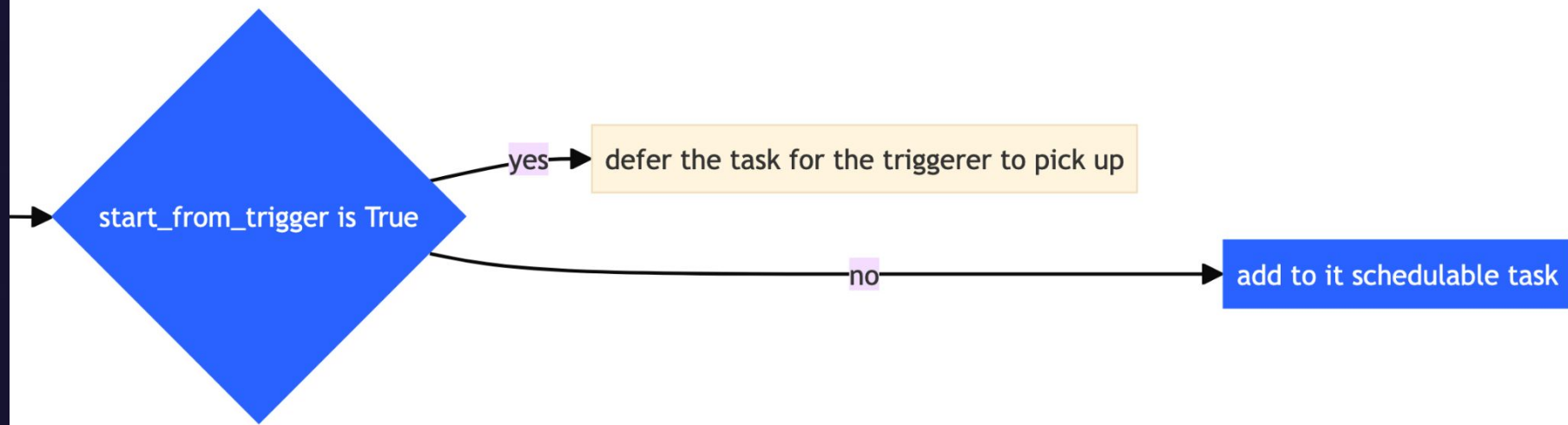# Start task execution in the triggerer

Under the hood (it's now…)

# Start task execution in the triggerer

Under the hood (it's now…)



Schedule a task

start_from_trigger is True

yes → defer the task for the triggerer to pick up

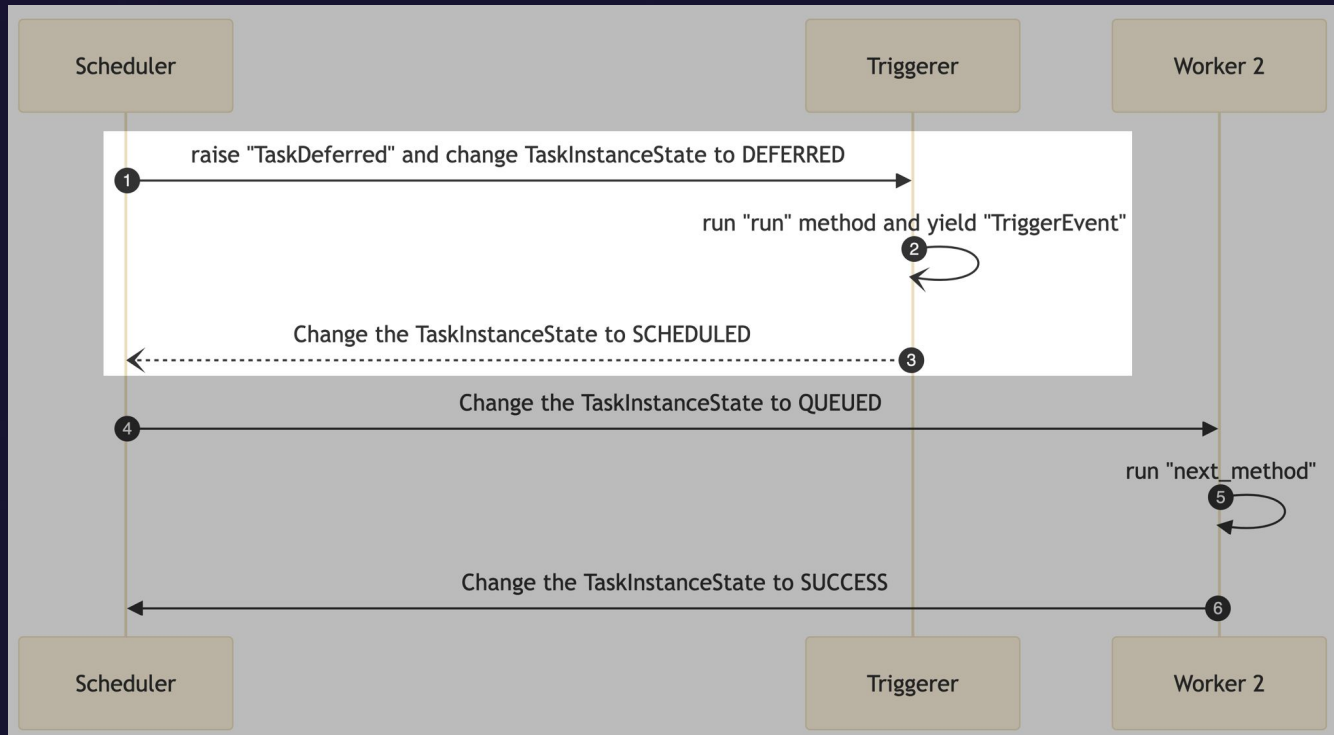no → add to it schedulable task

# Start task execution in the triggerer

**WHAT IF...?**

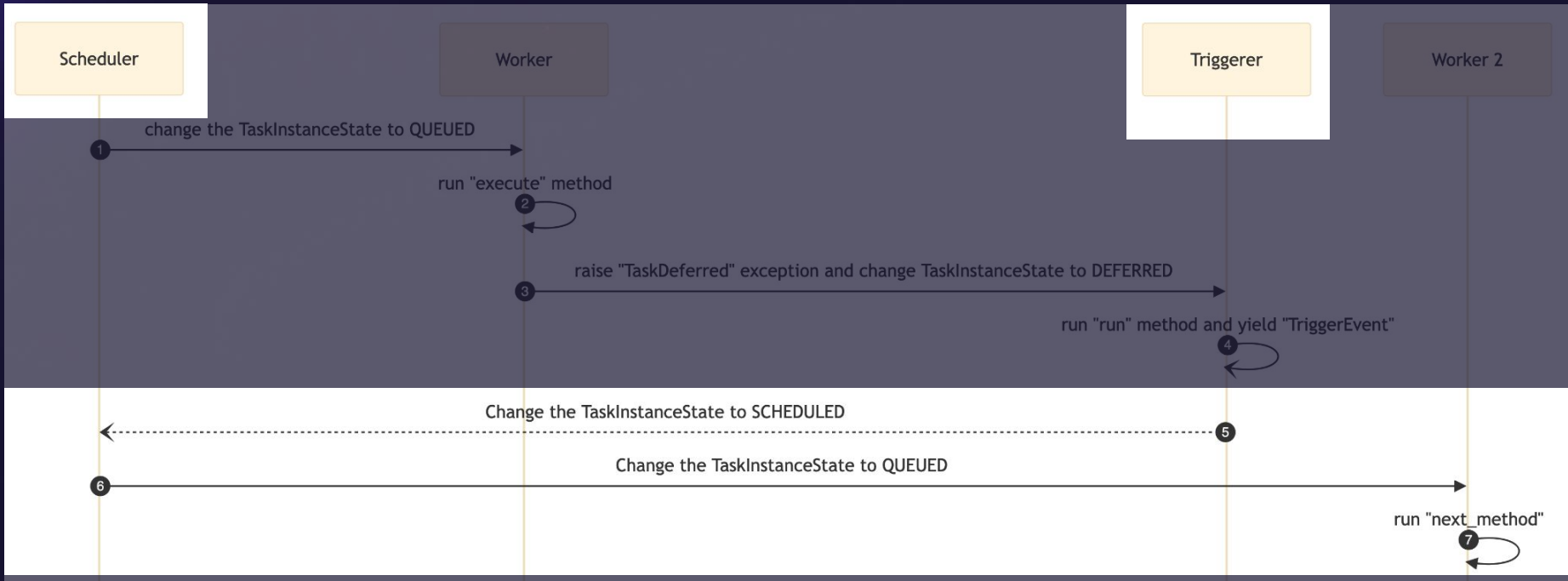We don't even have a next method to run after trigger finish its execution?

# Especially when "next_method" does nothing

```python
from datetime import timedelta
from typing import Any

from airflow.sensors.base import BaseSensorOperator
from airflow.triggers.temporal import TimeDeltaTrigger
from airflow.utils.context import Context


class WaitOneHourSensor(BaseSensorOperator):
    def execute(self, context: Context) -> None:
        self.defer(
            trigger=TimeDeltaTrigger(timedelta(hours=1)), method_name="execute_complete"
        )

    def execute_complete(
        self, context: Context, event: dict[str, Any] | None = None
    ) -> None:
        # We have no more work to do here. Mark as complete.
        return
```

ASTRONOMER

# What can we do?

# What're the things we want to change?

End task execution in the triggerer

# End task execution in the worker

```python
from datetime import timedelta
from typing import Any

from airflow.sensors.base import BaseSensorOperator
from airflow.triggers.temporal import TimeDeltaTrigger
from airflow.utils.context import Context


class WaitOneHourSensor(BaseSensorOperator):
    def execute(self, context: Context) -> None:
        self.defer(
            trigger=TimeDeltaTrigger(timedelta(hours=1)), method_name="execute_complete"
        )

    def execute_complete(
        self, context: Context, event: dict[str, Any] | None = None
    ) -> None:
        # We have no more work to do here. Mark as complete.
        return
```
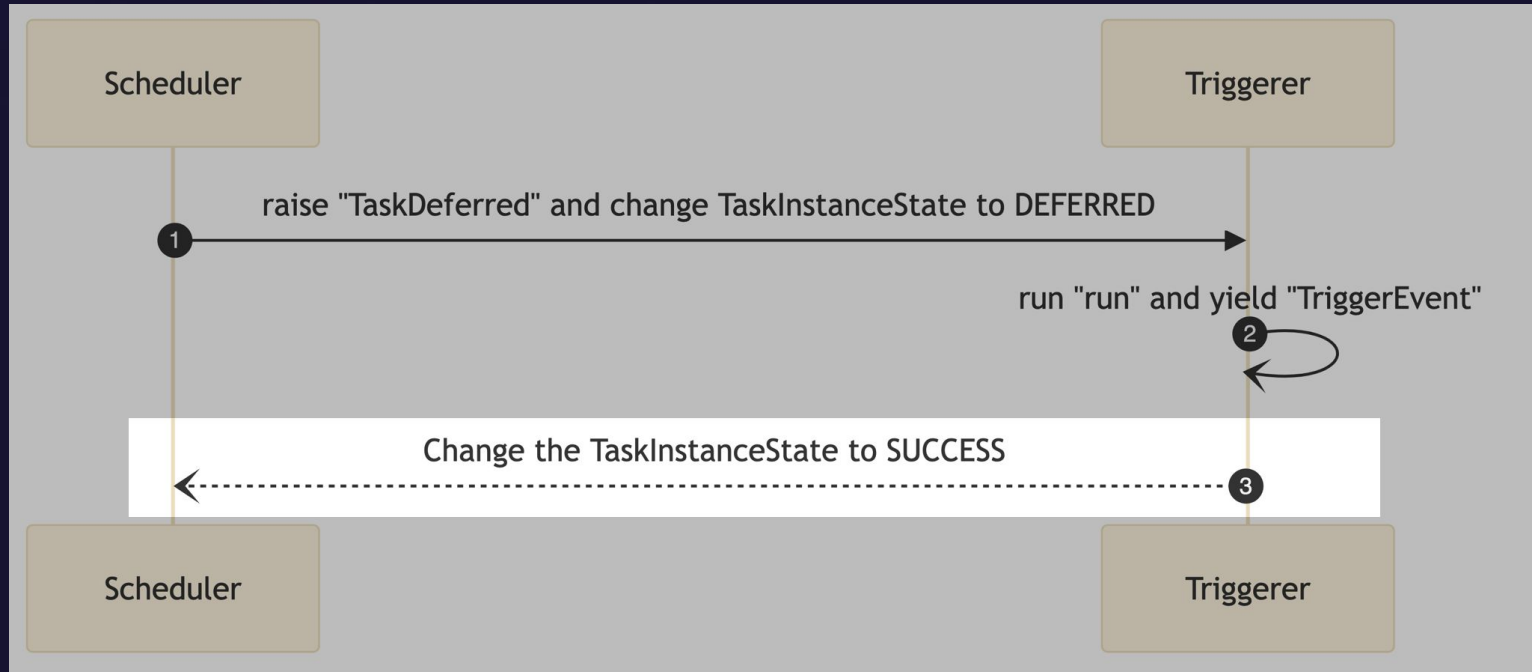
ASTRONOMER

# End task execution in the **worker**

```python
from datetime import timedelta
from typing import Any

from airflow.sensors.base import BaseSensorOperator
from airflow.triggers.temporal import TimeDeltaTrigger
from airflow.utils.context import Context


class WaitOneHourSensor(BaseSensorOperator):
    def execute(self, context: Context) -> None:
        self.defer(
            trigger=TimeDeltaTrigger(timedelta(hours=1)), method_name="execute_complete"
        )

    def execute_complete(
        self, context: Context, event: dict[str, Any] | None = None
    ) -> None:
        # We have no more work to do here. Mark as complete.
        return
```

ASTRONOMER

# End task execution in the triggerer

Well… we need to make some change in the trigger this time

```python
# temporal.py

13
12    class TimeDeltaTrigger(DateTimeTrigger):
11        """
10
9         def __init__(
8             self,
7             delta: datetime.timedelta,
6             *,
5             end_from_trigger: bool = False
4         ) -> None:
3             super().__init__(
2                 moment=timezone.utcnow() + delta,
1                 end_from_trigger=end_from_trigger
113           )
```

ASTRONOMER

# End task execution in the triggerer

Well… we need to make some change in the trigger

```python
# temporal.py

class DateTimeTrigger(BaseTrigger):
    async def run(self) -> AsyncIterator[TriggerEvent]:
        for step in 3600, 60, 10:
            seconds_remaining = (
                self.moment - pendulum.instance(timezone.utcnow())
            ).total_seconds()
            while seconds_remaining > 2 * step:
                self.log.info("%d seconds remaining; sleeping %s seconds", seconds_remaining, step)
                await asyncio.sleep(step)
                seconds_remaining = (
                    self.moment - pendulum.instance(timezone.utcnow())
                ).total_seconds()
        # Sleep a second at a time otherwise
        while self.moment > pendulum.instance(timezone.utcnow()):
            self.log.info("sleeping 1 second...")
            await asyncio.sleep(1)
        if self.end_from_trigger:
            self.log.info("Sensor time condition reached; marking task successful and exiting")
            yield TaskSuccessEvent()
        else:
            self.log.info("yielding event with payload %r", self.moment)
            yield TriggerEvent(self.moment)
```

# End task execution in the triggerer

Yield a TaskSuccessEvent

```python
if self.end_from_trigger:
    self.log.info("Sensor time condition reached; marking task successful and exiting")
    yield TaskSuccessEvent()
else:
    self.log.info("yielding event with payload %r", self.moment)
    yield TriggerEvent(self.moment)
```

# End task execution in the triggerer

Newly supported TriggerEvents

```
1  class TaskSuccessEvent(BaseTaskEndEvent):
2      """Yield this event in order to end the task successfully."""
3
4      task_instance_state = TaskInstanceState.SUCCESS
5
6
7  class TaskFailedEvent(BaseTaskEndEvent):
8      """Yield this event in order to end the task with failure."""
9
10     task_instance_state = TaskInstanceState.FAILED
11
12
13 class TaskSkippedEvent(BaseTaskEndEvent):
14     """Yield this event in order to end the task with status 'skipped'."""
15
16     task_instance_state = TaskInstanceState.SKIPPED
```

# End task execution in the triggerer

## Under the hook

```
163  class TriggererJobRunner(BaseJobRunner, LoggingMixin):
  7      @add_span
  6      def handle_events(self):
  5          """Dispatch outbound events to the Trigger model which pushes
  4          while self.trigger_runner.events:
  3              # Get the event and its trigger ID
  2              trigger_id, event = self.trigger_runner.events.popleft()
  1              # Tell the model to wake up its tasks
409              Trigger.submit_event(trigger_id=trigger_id, event=event)
  1              # Emit stat event
  2              Stats.incr("triggers.succeeded")
```

# End task execution in the triggerer

## Under the hook

```
🐍 triggerer_job_runner.py        ×          🐍 trigger.py        ×
164  class Trigger(Base):
 10      @classmethod
  9      @internal_api_call
  8      @provide_session
  7      def submit_event(cls, trigger_id, event, session: Session = NEW_SESSION) -> None:
  6          """Take an event from an instance of itself, and trigger all dependent tasks to resume."""
  5          for task_instance in session.scalars(
  4              select(TaskInstance).where(
  3                  TaskInstance.trigger_id == trigger_id, TaskInstance.state == TaskInstanceState.DEFERRED
  2              )
  1          ):
208              event.handle_submit(task_instance=task_instance)
```

# End task execution in the triggerer

Under the hook, it updates the state...

```python
# triggerer_job_runner.py        # trigger.py        # base.py
33  class BaseTaskEndEvent(TriggerEvent):
12      @provide_session
11      def handle_submit(self, *, task_instance: TaskInstance, session: Session = NEW_SESSION) -> None:
10          """
9           Submit event for the given task instance.
8
7           Marks the task with the state `task_instance_state` and optionally pushes xcom if applicable.
6
5           :param task_instance: The task instance to be submitted.
4           :param session: The session to be used for the database callback sink.
3           """
2           # Mark the task with terminal state and prevent it from resuming on worker
1           task_instance.trigger_id = None
206         task_instance.state = self.task_instance_state
1           self._submit_callback_if_necessary(task_instance=task_instance, session=session)
2           self._push_xcoms_if_necessary(task_instance=task_instance)
```

ASTRONOMER

# End task execution in the triggerer

based on the TriggerEvent type

```python
class TaskSuccessEvent(BaseTaskEndEvent):
    """Yield this event in order to end the task successfully."""

    task_instance_state = TaskInstanceState.SUCCESS


class TaskFailedEvent(BaseTaskEndEvent):
    """Yield this event in order to end the task with failure."""

    task_instance_state = TaskInstanceState.FAILED


class TaskSkippedEvent(BaseTaskEndEvent):
    """Yield this event in order to end the task with status 'skipped'."""

    task_instance_state = TaskInstanceState.SKIPPED
```

# End task execution in the triggerer

which used to be always set as SCHEDULED



```
    9  ████□  airflow/models/trigger.py

        @@ -203,14 +203,7 @@ def submit_event(cls, trigger_id, event, session: Session = NEW_SESSION) -> None
203 203              TaskInstance.trigger_id == trigger_id, TaskInstance.state == TaskInstanceState.DEFERRED
204 204          )
205 205      ):
206     -       # Add the event's payload into the kwargs for the task
207     -       next_kwargs = task_instance.next_kwargs or {}
208     -       next_kwargs["event"] = event.payload
209     -       task_instance.next_kwargs = next_kwargs
210     -       # Remove ourselves as its trigger
211     -       task_instance.trigger_id = None
212     -       # Finally, mark it as scheduled so it gets re-queued
213     -       task_instance.state = TaskInstanceState.SCHEDULED
    206 +       event.handle_submit(task_instance=task_instance)
```

# Credit

author of the end from trigger feature

# Ankit Chaurasia

## Senior Software Engineer at Astronomer

Ankit Chaurasia is a Senior Software Engineer at Astronomer, where he focuses on the design and engineering of Apache Airflow. He is an advocate for open-source projects and has contributed to initiatives such as Apache Airflow, Ask-Astro, and OpenCV CVAT. Previously, Ankit led teams at Wadhwani AI, developing AI solutions for healthcare and agriculture, which resulted in winning a $2 million Google AI Challenge grant.

More at https://ankitchaurasia.info/

## Sessions by Ankit Chaurasia

- Mastering Advanced Dataset Scheduling in Apache Airflow (2024)

ASTRONOMER

# 17:40 at the same room

## Mastering Advanced Dataset Scheduling in Apache Airflow

Speaker(s):



Ankit Chaurasia

Sep-11 17:40-18:05 in Elizabethan A+B    📅 Add to Calendar

Are you looking to harness the full potential of data-driven pipelines with Apache Airflow? This session will dive into the newly introduced conditional expressions for advanced dataset scheduling in Airflow - a feature highly requested by the Airflow community. Attendees will learn how to effectively use logical operators to create complex dependencies that trigger DAGs based on the dataset updates in real-world scenarios. We'll also explore the innovative DatasetOrTimeSchedule, which combines time-based and dataset-triggered scheduling for unparalleled flexibility. Furthermore, attendees will discover the latest API endpoints that facilitate external updates and resets of dataset events, streamlining workflow management across different deployments.

This talk also aims to explain:

- The basics of using conditional expressions for dataset scheduling.

- How do we integrate time-based schedules with dataset triggers?

- Practical applications of the new API endpoints for enhanced dataset management.

- Real-world examples of how these features can optimize your data workflows.

# How does it affect DAG authors?

- Release more worker slot
- Improve operators and sensors for efficiency
- Reduce resource usage, which indicates cost saving
- More new use cases to come after more operator authors apply this new feature

# How does it affect operator authors?

- A new way to implement operators in an asynchronous manner
- Simplify operators / sensors by reducing unnecessary "execute" and "execute_complete" methods
(most applicable to sensors I think)

# How does it affect Airflow?

Potential to run all tasks in async

# How does it affect Airflow?

Potential to run all tasks in async ?

# They're included in Airflow 2.10.0 🎉

You can find it if you scroll down to the end of
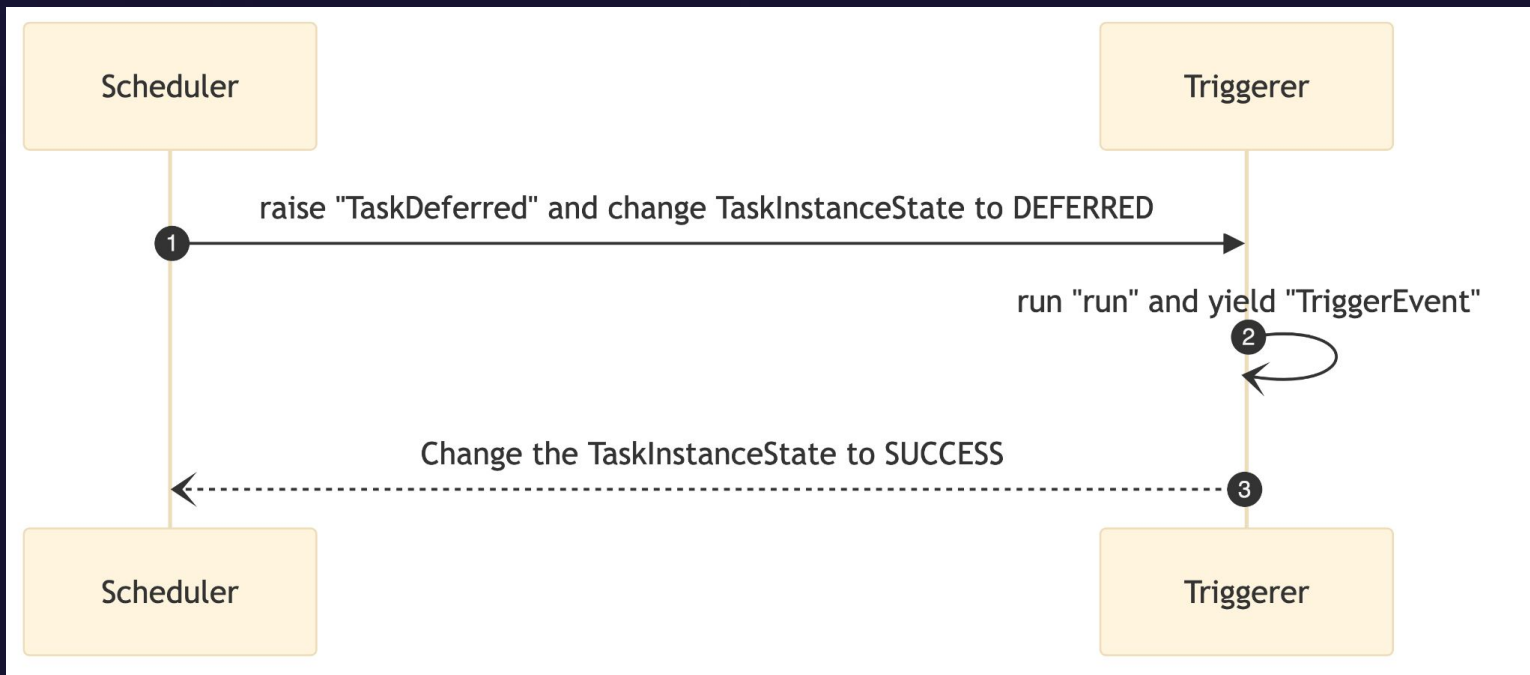Airflow 2.10.0 post

## Additional new features

Here are just a few interesting new features since there are too many to list in full:

- Deferrable operators can now execute directly from the triggerer without needing to go through the worker. This is especially efficient for certain operators, like sensors, and can help teams save both time and money.

ASTRONOMER

# Running Airflow Tasks without the workers

There's no more "What if…"

THE CIRCLE IS NOW COMPLETE.

Well yes, but actually no

# Limitation (Start execution from the trigger)

Limited dynamic task mapping support

```python
class WaitHoursSensor(BaseSensorOperator):
    # You'll need to change trigger_cls to the actual path to HourDeltaTrigger.
        timeout=None,
    )

    def __init__(
        self,
        *args: list[Any],
        trigger_kwargs: dict[str, Any] | None,
        start_from_trigger: bool,
        **kwargs: dict[str, Any],
    ) -> None:
        # This whole method will be skipped during dynamic task mapping.

        super().__init__(*args, **kwargs)
        self.start_trigger_args.trigger_kwargs = trigger_kwargs
        self.start_from_trigger = start_from_trigger
```

# Limitation (Start execution from the trigger)

trigger_kwargs, start_from_trigger required in __init__

```python
class WaitHoursSensor(BaseSensorOperator):
    # You'll need to change trigger_cls to the actual path to HourDeltaTrigger.
        timeout=None,
    )

    def __init__(
        self,
        *args: list[Any],
        trigger_kwargs: dict[str, Any] | None,
        start_from_trigger: bool,
        **kwargs: dict[str, Any],
    ) -> None:
        # This whole method will be skipped during dynamic task mapping.

        super().__init__(*args, **kwargs)
        self.start_trigger_args.trigger_kwargs = trigger_kwargs
        self.start_from_trigger = start_from_trigger
```

# Limitation (Start execution from the trigger)

the whole **__init__** method skipped before execution

```python
class WaitHoursSensor(BaseSensorOperator):
    # You'll need to change trigger_cls to the actual path to HourDeltaTrigger.
        timeout=None,
    )

    def __init__(
        self,
        *args: list[Any],
        trigger_kwargs: dict[str, Any] | None,
        start_from_trigger: bool,
        **kwargs: dict[str, Any],
    ) -> None:
        # This whole method will be skipped during dynamic task mapping.

        super().__init__(*args, **kwargs)
        self.start_trigger_args.trigger_kwargs = trigger_kwargs
        self.start_from_trigger = start_from_trigger
```

# Limitation (Start execution from the trigger)

slightly different syntax

```python
WaitHoursSensor.partial(
    task_id="wait_for_n_hours",
    start_from_trigger=True
).expand(
    trigger_kwargs=[
        {"hours": 1},
        {"hours": 2}
    ]
)
```

If you don't know if you need to combine this feature with dynamic task mapping

# Limitation (End execution from the trigger)

## Doesn't support listeners

> **❗ Note**
>
> Exiting from the trigger works only when listeners are not integrated for the deferrable operator. Currently, when deferrable operator has the `end_from_trigger` attribute set to `True` and listeners are integrated it raises an exception during parsing to indicate this limitation. While writing the custom trigger, ensure that the trigger is not set to end the task instance directly if the listeners are added from plugins. If the `end_from_trigger` attribute is changed to different attribute by author of trigger, the DAG parsing would not raise any exception and the listeners dependent on this task would not work. This limitation will be addressed in future releases.

# QR Code links to my posts related to this talk

Thank you!
Any questions?