# Airflow at Burns & McDonnell

## Orchestration from 0 to 100

**Bonnie Why**

# Bonnie Why

is just your average person with way too many interests and not enough time.

**Bonnie Why**

is just your average person with way too many interests and not enough time.

And that's my dog

# AGENDA

1. Data at Burns & McDonnell

2. Ingestion in one day (or less)

3. Airflow is awesome

**Slide Deck and Resources**

# ABOUT
# BURNS & MCDONNELL

**STRENGTH** — **14,000+** PROFESSIONALS

**DEPTH** — **75+** OFFICES WORLDWIDE

**EXCELLENCE** — **#7** TOP 500 DESIGN FIRMS
*Engineering News-Record*

**COMMITMENT** — **100%** EMPLOYEE-OWNED

# IVANPAH SOLAR THERMAL POWER FACILITY

## NRG ENERGY | NIPTON, CA

# LUNAR PRODUCTION & OPERATION CENTER

## INTUITIVE MACHINES| HOUSTON, TX

# STARTING LINE
## DATA AT BURNS & MCDONNELL

HARD TO
MAINTAIN

- **Multitude of varied source systems**

- **Multitude of differing tools and processes**

- **Teams working in silos**

HARD TO MAINTAIN

HARD TO TRUST

- **Duplicated effort, duplicated data**

- **Lack of discoverability**

- **Unclear as to what the data means**

HARD TO TRUST

HARD TO CHANGE

- **Brittle, interconnected systems**

- **Unsure of who is using the data, and how often**

- **Little testing, end users find the bugs**

HARD TO CHANGE

# FINISH LINE
## DATA AT BURNS & MCDONNELL

SCALEABLE

- **Scale with the business, as well as the data**

- **Centralized for uniform access**

- **Maintainable to keep up with high demand**

SCALEABLE

RELIABLE

- **Company-wide understanding, shared language**

- **Using the "right" data**

- **Obvious, secure, and accurate data lineage**

RELIABLE

EVOLVABLE

- **Composable data management strategy**
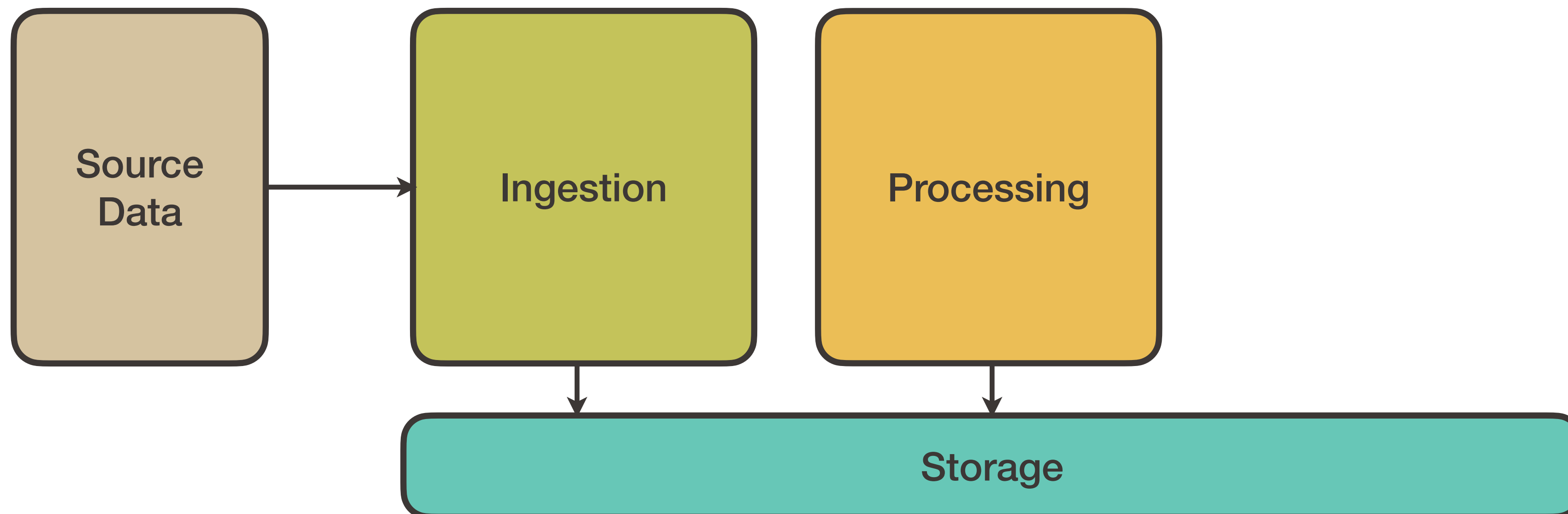
- **Grow with emerging technologies**

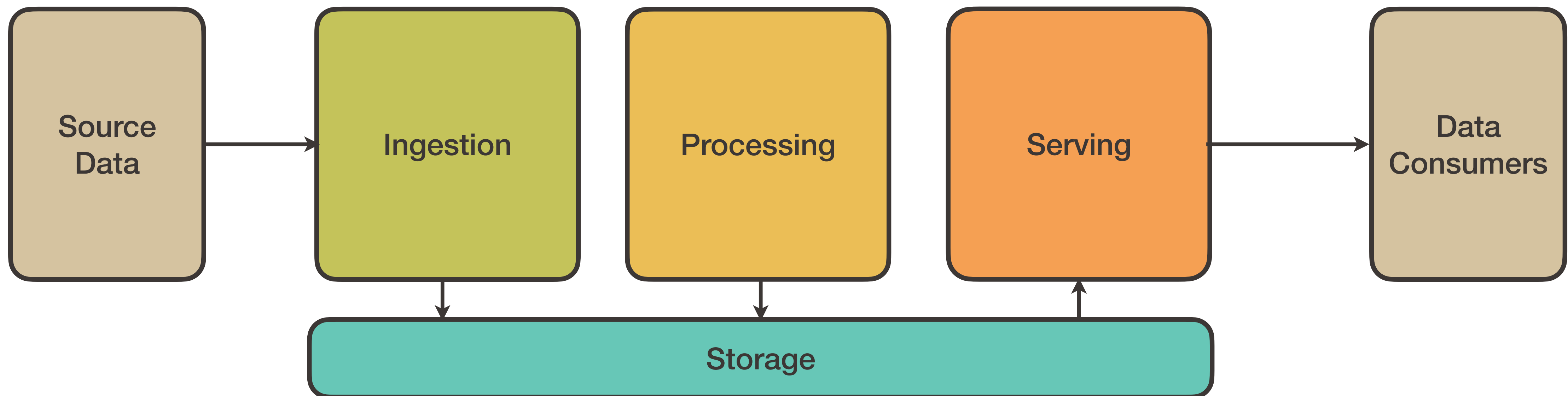- **Support any new use cases our business needs**

EVOLVABLE

Source Data → Ingestion

Source Data → Ingestion → Storage

Source Data → Ingestion → Processing → Serving → Data Consumers

Storage

Adapted from "Designing Cloud Data Platforms" by
Daniel Zburivsky, Lynda Partner

Orchestration

Operational Metadata

| Source Data | Ingestion | Processing | Serving | Data Consumers |

Storage

Adapted from "Designing Cloud Data Platforms" by Daniel Zburivsky, Lynda Partner

# ON YOUR MARK
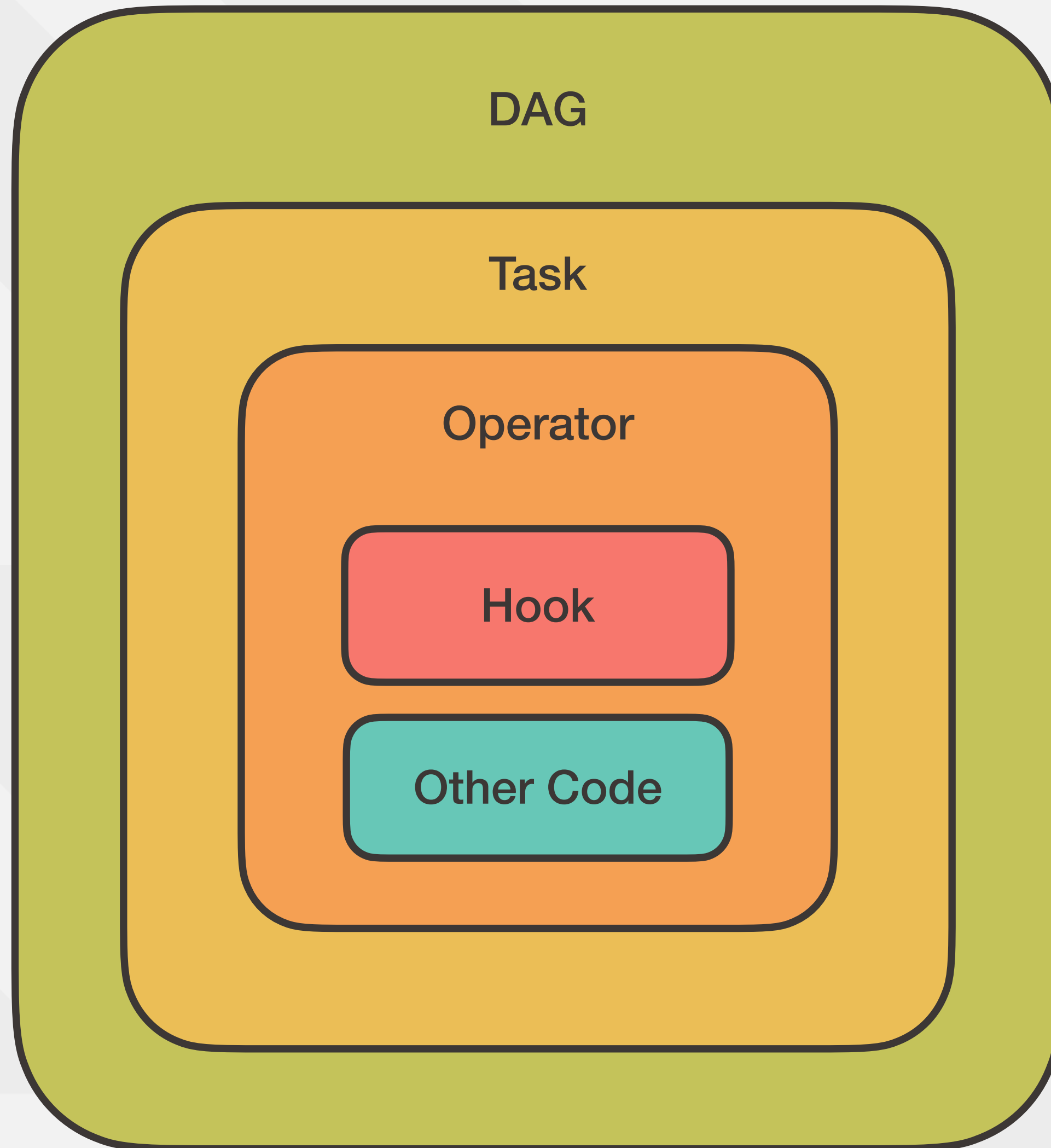
You have one day.

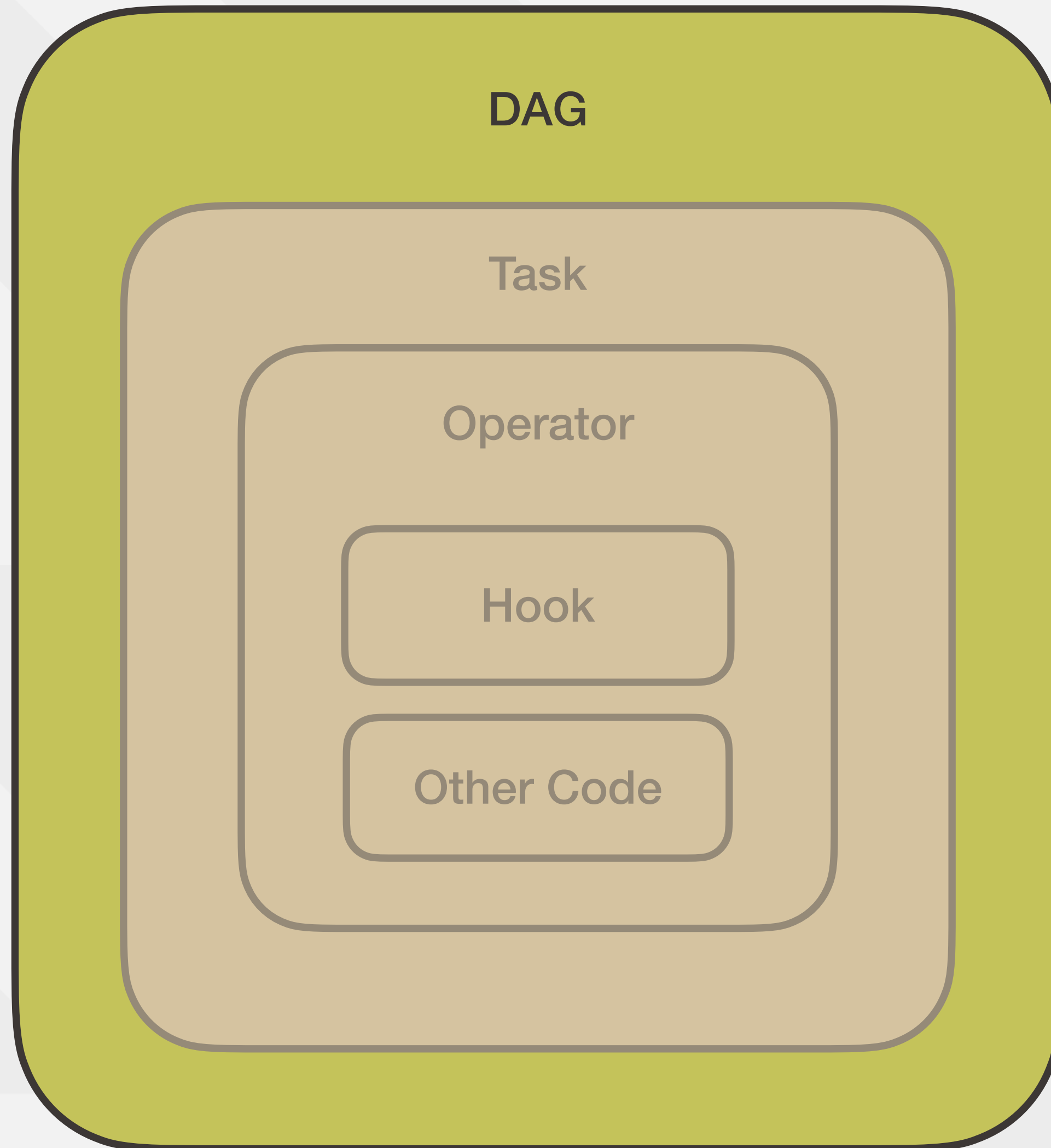STARTING AT ZERO

# GETTING STARTED

```python
 7
 8  def log_data_from_oracle():
 9      oracle_hook = OracleHook(oracle_conn_id="oracle")
10      sql_query = "SELECT * from A_REDACTED_TABLE"
11      records = oracle_hook.get_records(sql=sql_query)
12      for record in records:
13          print(record)
14
15
16  with DAG(
17      "oracle_connector_test_dag",
18      description="DAG to read data from Oracle and log it",
19      schedule=timedelta(days=1),
20      start_date=datetime(2023, 8, 23),
21      catchup=False,
22  ) as dag:
23      read_from_oracle = PythonOperator(
24          task_id="read_from_oracle",
25          python_callable=log_data_from_oracle,
26      )
27
28      read_from_oracle
29
```
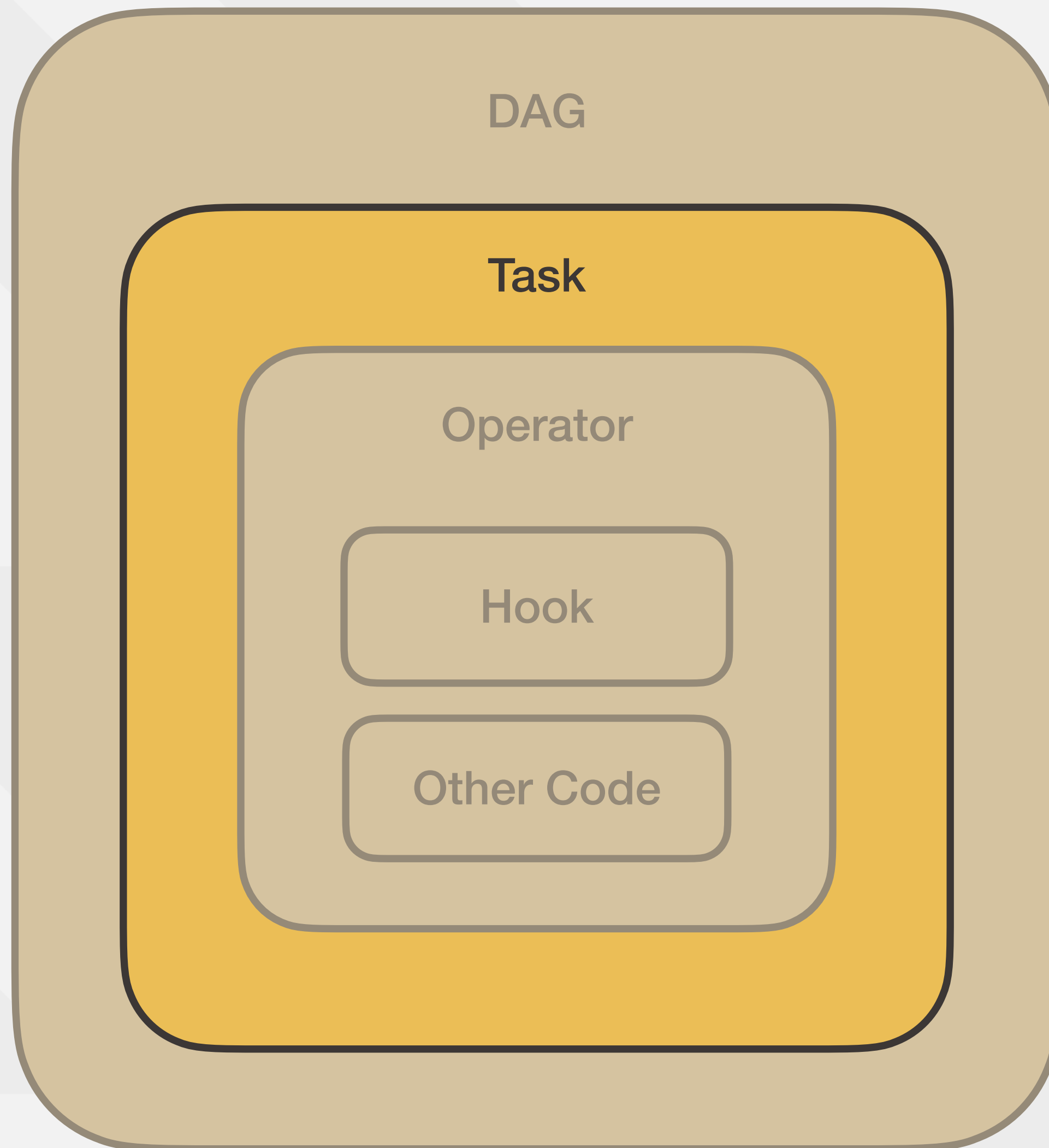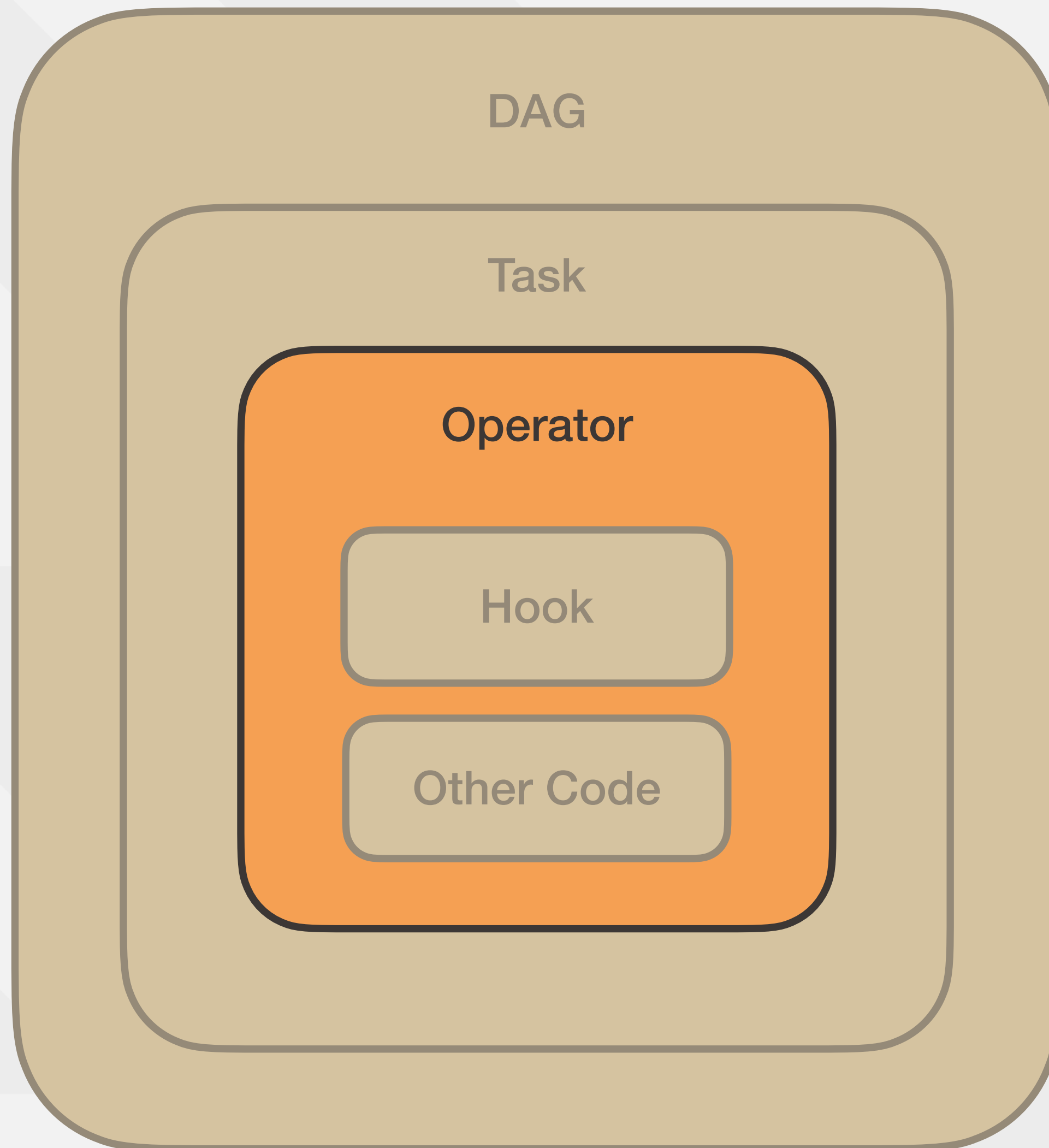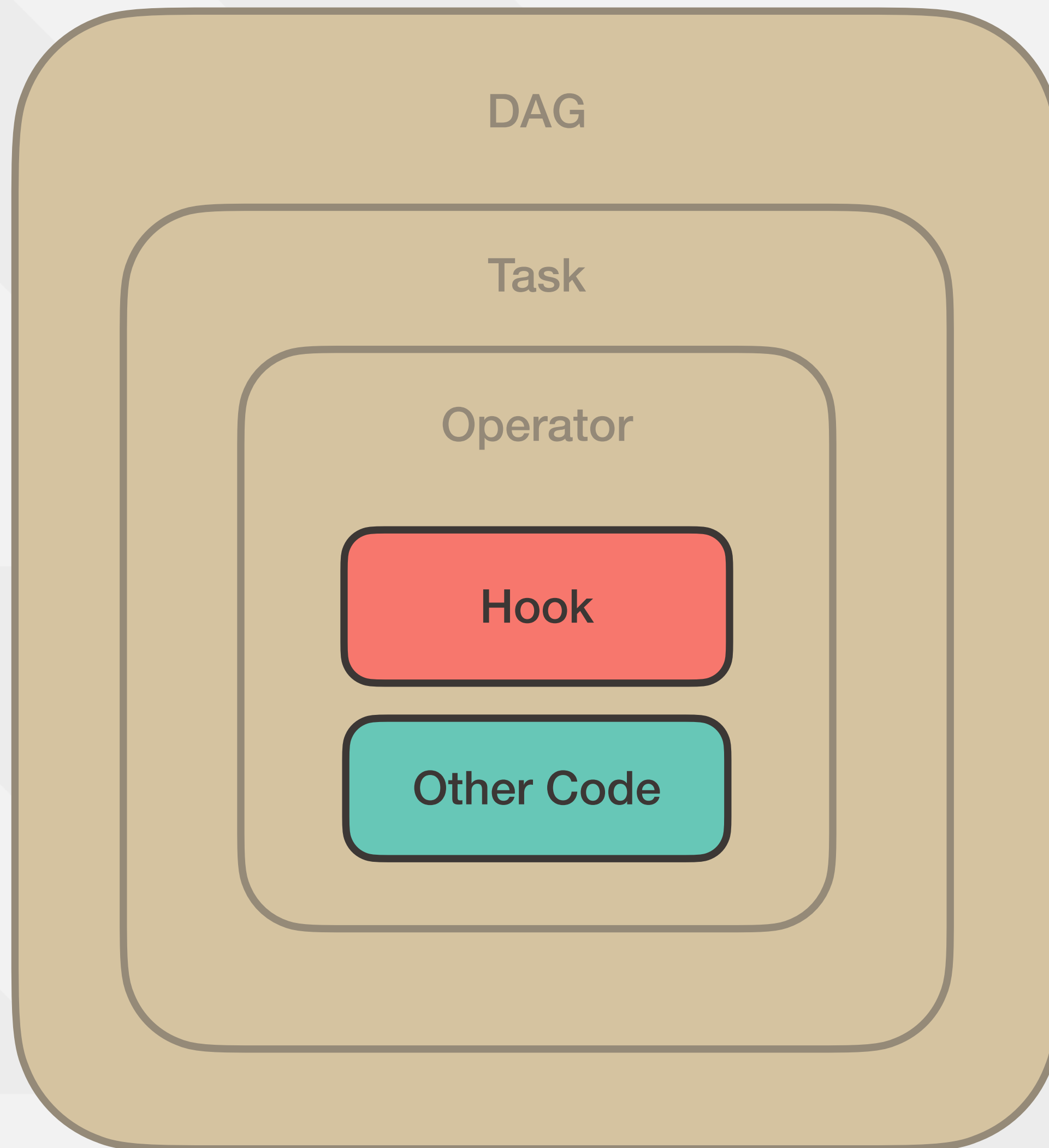
GETTING STARTED

DAG

Task

Operator

Hook

Other Code

GETTING STARTED

DAG

Task

Operator

Hook

Other Code

GETTING STARTED

DAG

Task

Operator

Hook

Other Code

GETTING STARTED

DAG

Task

Operator

Hook

Other Code

GETTING STARTED

DAG

Task

Operator

Hook

Other Code

GETTING STARTED

GETTING STARTED

GETTING STARTED

# CUSTOM OPERATOR

Source Table → Fetch → Query Results → Convert → Parquet Files → Upload → Landing Zone

**CUSTOM OPERATOR**

Fetch

Query
Results

Convert

Parquet
Files

Upload

```python
class CopyOperator(BaseOperator):
    def __init__(
        self,
        query,
        config_data,
        **kwargs,
    ):
        super().__init__(**kwargs),
        self._query = query
        self._config_data = config_data
        self._domain = config_data.get_batch_property("domain")
        self._target_container = self.config_data.get_batch_property("target_containe
        self._load_type = config_data.get_batch_property("load_type")

    def execute(self, **context):
        df = self.get_source_df()

        parquet_file = df.to_parquet(path=None, engine="pyarrow")

        self.upload_file(parquet_file)

    def get_source_df(self):
        hook = getattr(Hooks, self._config_data.get_batch_property("source"))
        source_file_path = self.get_source_file_path()
        sql_query = self.read_file(source_file_path)
        df = hook.get_pandas_df(sql=sql_query)
```

**CUSTOM OPERATOR**

CUSTOM OPERATOR

Fetch

Query
Results

Convert

Parquet
Files

Upload

```python
 8  class CopyOperator(BaseOperator):
 9      def __init__(
10          self,
11          query,
12          config_data,
13          **kwargs,
14      ):
15          super().__init__(**kwargs),
16          self._query = query
17          self._config_data = config_data
18          self._domain = config_data.get_batch_property("domain")
19          self._target_container = self.config_data.get_batch_property("target_container")
20          self._load_type = config_data.get_batch_property("load_type")
21
22      def execute(self, **context):
23          df = self.get_source_df()
24
25          parquet_file = df.to_parquet(path=None, engine="pyarrow")
26
27          self.upload_file(parquet_file)
28
29      def get_source_df(self):
30          hook = getattr(Hooks, self._config_data.get_batch_property("source"))
31          source_file_path = self.get_source_file_path()
32          sql_query = self.read_file(source_file_path)
33          df = hook.get_pandas_df(sql=sql_query)
```

**CUSTOM OPERATOR**

```python
37
@task
def fetch_the_things(hook, things, **context):

    things_to_return = get_things(things)

    for thing in things:
        things_to_return.append(things)

    return things_to_return

47
```

CUSTOM OPERATOR

# GENERATOR PATTERN

DAG

Task

Operator

Hook

Other
Code

**GENERATOR PATTERN**

GENERATOR PATTERN

```yaml
batch:
  source_system: oracle
  source_schema: make_believe
  migration_system: sunset_blvd
  migration_schema: california
  load_type: incremental

connection:
  source_conn_id: oracle
  target_conn_id: azure

dag:
  name: sql_source_dag
  description: Ingests all data from sql source
  schedule_interval: "0 11 * * *"
  start_date: "2024-07-01"
  owner: date_team
  retries: 3
  retry_delay: 3
  pool: default_pool
  catchup: true

tasks:
  queries:
    - name: REDACTED_TABLE_NAME
```

DAG

Task

Operator

Hook

Other
Code

GENERATOR PATTERN

GENERATOR PATTERN

GENERATOR PATTERN

GENERATOR PATTERN

DAG

Task

Operator

Hook

Other
Code

```python
17  def create_dag_from_config(config_data):
18      default_args = {
19          "owner": config_data.get_dag_property("owner"),
20          "depends_on_past": config_data.get_dag_property("depends_on_past"),
21          "email_on_failure": config_data.get_dag_property("email_on_failure"),
22          "email_on_retry": config_data.get_dag_property("email_on_retry"),
23          "retries": config_data.get_dag_property("retries"),
24          "retry_delay": timedelta(minutes=config_data.get_dag_property("retry_del
25          "pool": "mis_pool",
26          "on_failure_callback": [task_failure_log_callback],
27      }
28
29      @dag(
30          dag_id=config_data.get_dag_property("name"),
31          schedule=config_data.get_dag_property("schedule_interval"),
32          default_args=default_args,
33          start_date=datetime.strptime(
34              config_data.get_dag_property("start_date"), "%Y-%m-%d"
35          ),
36          catchup=config_data.get_dag_property("catchup", default=False),
37          tags=config_data.get_tags(),
38          render_template_as_native_obj=True,
39      )
40      def generated_dag():
41          for query_info in config_data.queries:
42              mapped_tasks = config_data.get_mapped_task_intervals(query_info)
```

GENERATOR PATTERN

```
60
61  CONFIG_DIR = Path(__file__).parent / "config"
62  dag_configs = DAGConfig.from_directory(Path(CONFIG_DIR))
63
64  for config_data in dag_configs:
65      dag_id = config_data.get_dag_property("name")
66      globals()[dag_id] = create_dag_from_config(config_data)
67
```

DAG

Task

Operator

Hook

Other
Code

GENERATOR PATTERN

GENERATOR PATTERN

```python
def create_dag_from_config(config_data):
        catchup=False,
    )
    def generated_dag():
        sql_dir = str(Path(__file__).parent / "config" / "sql")
        oracle_hook = OracleHook(oracle_conn_id="oracle")
        wasb_hook = WasbHook(wasb_conn_id="azure")

        for query_info in config_data.queries:
            file_name = query_info["name"] + ".parquet"
            target = config_data.get_landing_zone_target(
                    generate_base_path(config_data),
                    file_name
            )


            execute_sql_task = CopyOperator(
                task_id=f"execute_{query_info['name']}_sql",
                query=query_info["sql"],
                oracle_hook=oracle_hook,
                wasb_hook=wasb_hook,
                sql_path=sql_dir,
                azure_container_name=AZURE_CONTAINER_NAME,
                target=target,
            )

            execute_sql_task
```

DAG

Task

Operator

Hook

Other Code

30 MPH

# "LOCAL" DEVELOPMENT

# deployToSandbox

git commit -m "I hope this works"

git tag -f deployToSandbox

git push -f deployToSandbox

50 MPH

# CHUNKING AND PARALLELIZATION

Source Table → Fetch → Query Results → Convert → Parquet Files → Upload → Landing Zone

**CHUNKING AND PARALLELIZATION**

1+ HOUR

Source Table → Fetch → Query Results → Convert → Parquet Files → Upload → Landing Zone

CHUNKING AND PARALLELIZATION

CHUNKING AND PARALLELIZATION

```
57         return generated_dag()
34                 config_data.get_dag_property("start_date"), "%Y-%m-%d"
35             ),
36         catchup=config_data.get_dag_property("catchup", default=False),
37         tags=config_data.get_tags(),
38         render_template_as_native_obj=True,
39     )
40     def generated_dag():
41         for query_info in config_data.queries:
42             mapped_tasks = config_data.get_mapped_task_intervals(query_info)
43
44             execute_sql_task = CopyOperator.partial(
45                 task_id=f"execute_{query_info['name']}_sql",
46                 query_info=query_info,
47                 config_data=config_data,
48                 source_hook=config_data.get_batch_property("source"),
49                 target_hook=config_data.get_batch_property("target"),
50                 current_run_scheduled_date="{{ data_interval_end.to_date_string() }
51                 previous_run_scheduled_date="{{ data_interval_start.to_date_string(
52                 on_success_callback=[task_success_log_callback],
53             ).expand(query_clip_range=mapped_tasks)
54
55             execute_sql_task
56
57     return generated_dag()
58
```

DAG

Task

Operator

Hook

Other Code
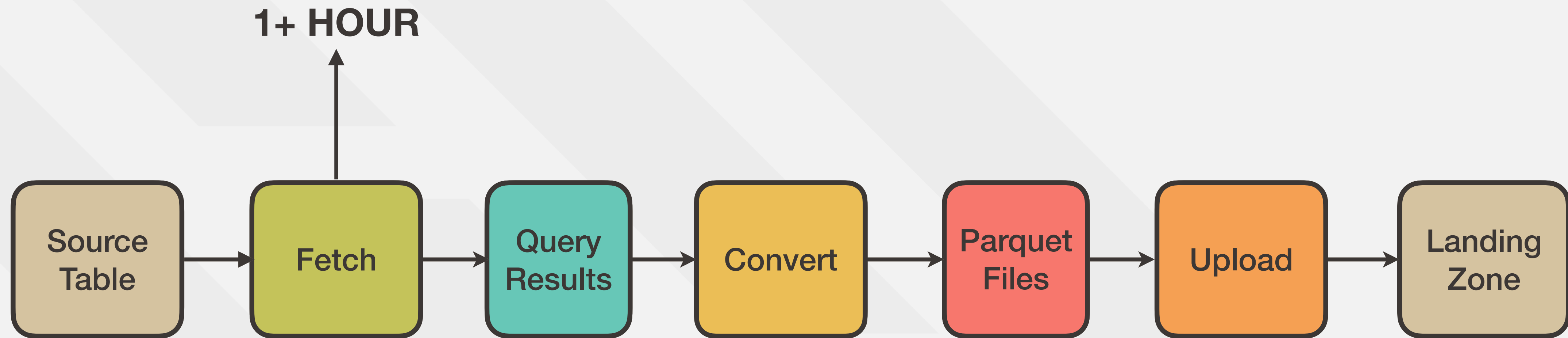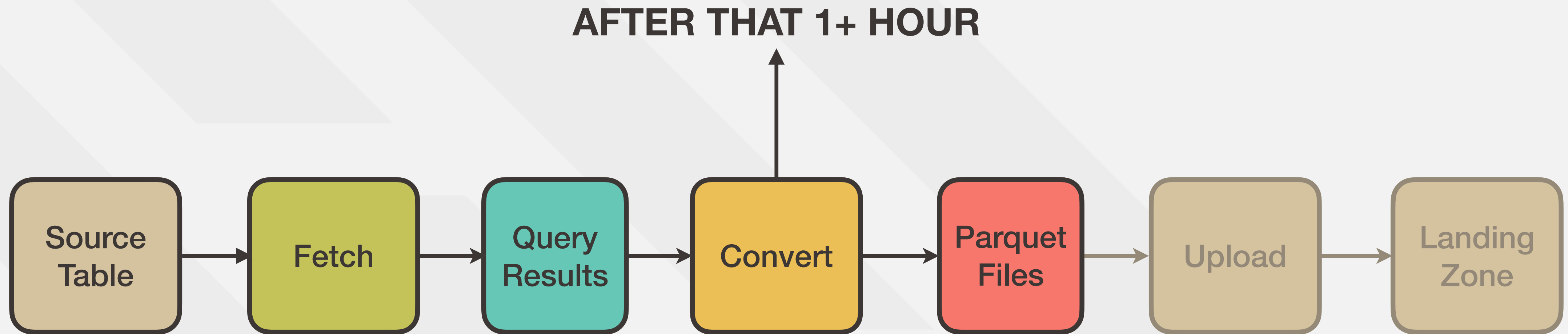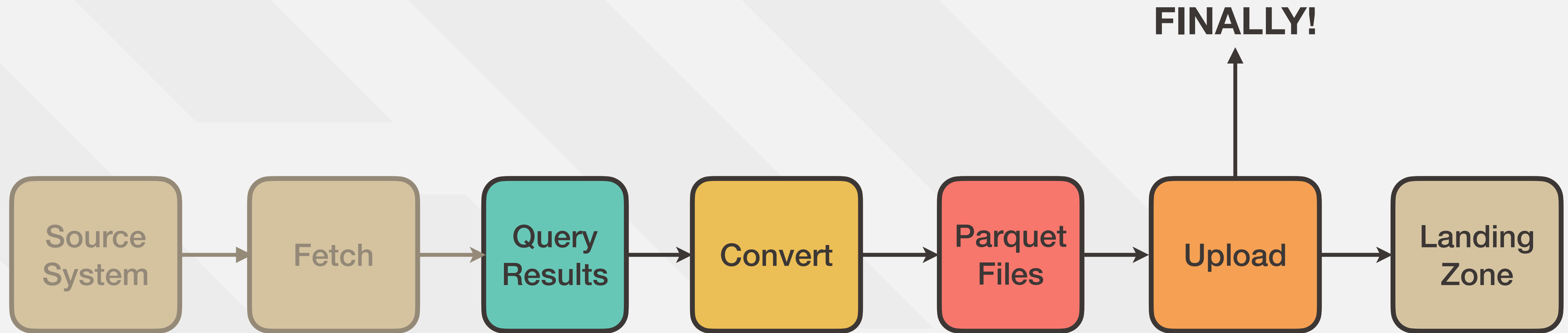
CHUNKING AND PARALLELIZATION

CHUNKING AND PARALLELIZATION

CHUNKING AND PARALLELIZATION

CHUNKING AND PARALLELIZATION

| depends_on_past | False |
| deps | frozenset({<TIDep(Trigger Rule)>, <TIDep(Task has been mapped)>, <TIDep(Previous Dagrun State)>, <TIDep(Not Previously Skipped)>, <TIDep(Not In Retry Period)>}) |
| downstream_task_ids | () |
| email | None |
| end_date | None |
| execution_timeout | None |
| executor_config | {} |

({'query_clip_range': [('0001-01-01', '2005-01-01'), ('2005-01-01', '2005-01-06'), ('2005-01-06', '2005-01-11'), ('2005-01-11', '2005-01-16'), ('2005-01-16', '2005-01-21'), ('2005-01-21', '2005-01-26'), ('2005-01-26', '2005-01-31'), ('2005-01-31', '2005-02-05'), ('2005-02-05', '2005-02-10'), ('2005-02-10', '2005-02-15'), ('2005-02-15', '2005-02-20'), ('2005-02-20', '2005-02-25'), ('2005-02-25', '2005-03-02'), ('2005-03-02', '2005-03-07'), ('2005-03-07', '2005-03-12'), ('2005-03-12', '2005-03-17'), ('2005-03-17', '2005-03-22'), ('2005-03-22', '2005-03-27'), ('2005-03-27', '2005-04-01'), ('2005-04-01', '2005-04-06'), ('2005-04-06', '2005-04-11'), ('2005-04-11', '2005-04-16'), ('2005-04-16', '2005-04-21'), ('2005-04-21', '2005-04-26'), ('2005-04-26', '2005-05-01'), ('2005-05-01', '2005-05-06'), ('2005-05-06', '2005-05-11'), ('2005-05-11', '2005-05-16'), ('2005-05-16', '2005-05-21'), ('2005-05-21', '2005-05-26'), ('2005-05-26', '2005-05-31'), ('2005-05-31', '2005-06-05'), ('2005-06-05', '2005-06-10'), ('2005-06-10', '2005-06-15'), ('2005-06-15', '2005-06-20'), ('2005-06-20', '2005-06-25'), ('2005-06-25', '2005-06-30'), ('2005-06-30', '2005-07-05'), ('2005-07-05', '2005-07-10'), ('2005-07-10', '2005-07-15'), ('2005-07-15', '2005-07-20'), ('2005-07-20', '2005-07-25'), ('2005-07-25', '2005-07-30'), ('2005-07-30', '2005-08-04'), ('2005-08-04', '2005-08-09'), ('2005-08-09', '2005-08-14'), ('2005-08-14', '2005-08-19'), ('2005-08-19', '2005-08-24'), ('2005-08-24', '2005-08-29'), ('2005-08-29', '2005-09-03'), ('2005-09-03', '2005-09-08'), ('2005-09-08', '2005-09-13'), ('2005-09-13', '2005-09-18'), ('2005-09-18', '2005-09-23'), ('2005-09-23', '2005-09-28'), ('2005-09-28', '2005-10-03'), ('2005-10-03', '2005-10-08'), ('2005-10-08', '2005-10-13'), ('2005-10-13', '2005-10-18'), ('2005-10-18', '2005-10-23'), ('2005-10-23', '2005-10-28'), ('2005-10-28', '2005-11-02'), ('2005-11-02', '2005-11-07'), ('2005-11-07', '2005-11-12'), ('2005-11-12', '2005-11-17'), ('2005-11-17', '2005-11-22'), ('2005-11-22', '2005-11-27'), ('2005-11-27', '2005-12-02'), ('2005-12-02', '2005-12-07'), ('2005-12-07', '2005-12-12'), ('2005-12-12', '2005-12-17'), ('2005-12-17', '2005-12-22'), ('2005-12-22', '2005-12-27'), ('2005-12-27', '2006-01-01'), ('2006-01-01', '2006-01-06'), ('2006-01-06', '2006-01-11'), ('2006-01-11', '2006-01-16'), ('2006-01-16', '2006-01-21'), ('2006-01-21', '2006-01-26'), ('2006-01-26', '2006-01-31'), ('2006-01-31', '2006-02-05'), ('2006-02-05', '2006-02-10'), ('2006-02-10', '2006-02-15'), ('2006-02-15', '2006-02-20'), ('2006-02-20', '2006-02-25'), ('2006-02-25', '2006-03-02'), ('2006-03-02', '2006-03-07'), ('2006-03-07', '2006-03-12'), ('2006-03-12', '2006-03-17'), ('2006-03-17', '2006-03-22'), ('2006-03-22', '2006-03-27'), ('2006-03-27', '2006-04-01'), ('2006-04-01', '2006-04-06'), ('2006-04-06', '2006-04-11'), ('2006-04-11', '2006-04-16'), ('2006-04-16', '2006-04-21'), ('2006-04-21', '2006-04-26'), ('2006-04-26', '2006-05-01'), ('2006-05-01', '2006-05-06'), ('2006-05-06', '2006-05-11'), ('2006-05-11', '2006-05-16'), ('2006-05-16', '2006-05-21'), ('2006-05-21', '2006-05-26'), ('2006-05-26', '2006-05-31'), ('2006-05-31', '2006-06-05'), ('2006-06-05', '2006-06-10'), ('2006-06-10', '2006-06-15'), ('2006-06-15', '2006-06-20'), ('2006-06-20', '2006-06-25'), ('2006-06-25', '2006-06-30'), ('2006-06-30', '2006-07-05'), ('2006-07-05', '2006-07-10'), ('2006-07-10', '2006-07-15'), ('2006-07-15', '2006-07-20'), ('2006-07-20', '2006-07-25'), ('2006-07-25', '2006-07-30'), ('2006-07-30', '2006-08-04'), ('2006-08-04', '2006-08-09'), ('2006-08-09', '2006-08-14'), ('2006-08-14', '2006-08-19'), ('2006-08-19', '2006-08-24'), ('2006-08-24', '2006-08-29'), ('2006-08-29', '2006-09-03'), ('2006-09-03', '2006-09-08'), ('2006-09-08', '2006-09-13'), ('2006-09-13', '2006-09-18'), ('2006-09-18', '2006-09-23'), ('2006-09-23', '2006-09-28'), ('2006-09-28', '2006-10-03'), ('2006-10-03', '2006-10-08'), ('2006-10-08', '2006-10-13'), ('2006-10-13', '2006-10-18'), ('2006-10-18', '2006-10-23'), ('2006-10-23', '2006-10-28'), ('2006-10-28', '2006-11-02'), ('2006-11-02', '2006-11-07'), ('2006-11-07', '2006-11-12'), ('2006-11-12', '2006-11-17'), ('2006-11-17', '2006-11-22'), ('2006-11-22', '2006-11-27'), ('2006-11-27', '2006-12-02'), ('2006-12-02', '2006-12-07'), ('2006-12-07', '2006-12-12'), ('2006-12-12', '2006-12-17'), ('2006-12-17', '2006-12-22'), ('2006-12-22', '2006-12-27'), ('2006-12-27', '2007-01-01'), ('2007-01-01', '2007-01-06'), ('2007-01-06', '2007-01-11'), ('2007-01-11', '2007-01-16'), ('2007-01-16', '2007-01-21'), ('2007-01-21', '2007-01-26'), ('2007-01-26', '2007-01-31'), ('2007-01-31', '2007-02-05'), ('2007-02-05', '2007-02-10'), ('2007-02-10', '2007-02-15'), ('2007-02-15', '2007-02-20'), ('2007-02-20', '2007-02-25'), ('2007-02-25', '2007-03-02'), ('2007-03-02',

**CHUNKING AND PARALLELIZATION**

CHUNKING AND PARALLELIZATION

# STRATEGY PATTERN

**Context**

- strategy

+ set(strategy)
+ do_a_thing(a, b)

<< interface >>
**Strategy**

+ do_a_thing(a, b)

strategy = strategy.do_a_thing(a, b)

**Concrete Strategy**

**Different Concreate Strategy**

**Different Concrete Strategy Again**

**STRATEGY PATTERN**

**Context**

- strategy

+ set(strategy)
+ do_a_thing(a, b)

<< interface >>
**Strategy**

+ do_a_thing(a, b)

strategy = strategy.do_a_thing(a, b)

Concrete Strategy

Different Concreate Strategy

Different Concrete Strategy Again

**STRATEGY PATTERN**

**Context**

- strategy

+ set(strategy)
+ do_a_thing(a, b)

<< interface >>
**Strategy**

+ do_a_thing(a, b)

strategy = strategy.do_a_thing(a, b)

Concrete Strategy

Different Concreate Strategy

Different Concrete Strategy Again

**STRATEGY PATTERN**

Context

- strategy

+ set(strategy)
+ do_a_thing(a, b)

strategy = strategy.do_a_thing(a, b)

<< interface >>
Strategy

+ do_a_thing(a, b)

Concrete Strategy

Different Concreate Strategy

Different Concrete Strategy Again

STRATEGY PATTERN

**Context**

- strategy

+ set(strategy)
+ do_a_thing(a, b)

strategy = strategy.do_a_thing(a, b)

<< interface >>
**Strategy**

+ do_a_thing(a, b)

Concrete Strategy

Different Concreate Strategy

Different Concrete Strategy Again

Adapted from
https://refactoring.guru/design-patterns/strategy

**STRATEGY PATTERN**

**Retriever**

- fetch_strategy

+ set(strategy)
+ fetch_data(a, b)

**<< interface >>**
**FetchStrategy**

+ fetch_data(a, b)

strategy = fetch_strategy.fetch_data(a, b)

**SQL Strategy**

**Date Interval SQL Strategy**

**Incremental SQL Strategy**

**STRATEGY PATTERN**

STRATEGY PATTERN

**Interface**

**Context**

**Strategy**

**Config**

**Utilize**

```python
class Retriever:

    def fetch_data(
        self, source_file_path, source_connection, query_range_interval=None
    ):
        """
        Returns a dataframe from the source system's response to the supplied
        SQL query.

        :param source_file_path: Path to the file containing information to
                                 call the downstream system.
        :param source_connection: Connection object to use when calling the
                                   source system
        :param query_clip_range: A tuple representing a range of dates to
                                  insert into a WHERE clause in the provided
                                  query.
        :return: Response from the source system
        """

        return self._strategy.fetch_data(
            source_file_path, source_connection, query_range_interval
        )
```

**STRATEGY PATTERN**

STRATEGY PATTERN

Interface

Context

Strategy

Config

Utilize

```yaml
12  dag:
13    name: sql_source_dag
14    description: Ingests all data from sql source
15    schedule_interval: "0 11 * * *"
16    start_date: "2024-07-01"
17    owner: date_team
18    retries: 3
19    retry_delay: 3
20    pool: default_pool
21    catchup: true
22
23  tasks:
24    queries:
25      - name: REDACTED_TABLE_NAME
26        predicate:
27          min_date: "2005-01-01"
28          max_date: "2024-09-01"
29          interval: "30 days"
30
```

STRATEGY PATTERN

80 MPH

ORCHESTRATE INGESTION

```
48              config_data=config_data,
49              source_hook=config_data.get_connection_property("source_conn_
50              target_hook=config_data.get_connection_property("target_conn_
51              current_run_scheduled_date="{{ data_interval_end | ds }}",
52              previous_run_scheduled_date="{{ data_interval_start | ds }}",
53          ).expand(query_clip_range=mapped_tasks)
54
55          sql_tasks.append(execute_sql_task)
56
57      databricks_run_task = DatabricksRunNowOperator(
58          databricks_conn_id="databricks",
59          task_id="databricks_run_task",
60          job_name="databricks_job_name",
61          notebook_params={
62              "storage_container"  "{{ get_variable('landing_zone') }}",
63              "storage_account":  {{ conn.azure.login }}",
64              "airflow_environment_path": "{{ reverse_domain(urlparse(conf.
65          },
66      )
67
68      sql_tasks >> databricks_run_task
69
```

DAG

Task

Operator

Hook

Other
Code

**ORCHESTRATE INGESTION**

ORCHESTRATE INGESTION

ORCHESTRATE INGESTION

100 MPH

Ludicrous Speed GO!

# ORCHESTRATE EVERYTHING

- **Config-driven API DAGs**

- **ExternalTaskSensors for silver jobs**

- **Improved DLT job integration**

ORCHESTRATE EVERYTHING

# FINISH LINE

## DATA AT BURNS & MCDONNELL

SCALEABLE

RELIABLE

EVOLVABLE

# THANK YOU!



Slide Deck and Resources