Left side:

```
1  """
2      pipeline_framework: Foundational functions and
    classes that manage the ROMEO yaml-based pipeline
    s.
3  """
4  import abc
5  import gzip
6  import logging
7  import pathlib
8  import shutil
9  import sys
10 import typing
11 from dataclasses import dataclass, field
12
13 import yaml
14
15
16 @dataclass
17 class PipelineState:
18     """Encapsulates the runtime state of a pipelin
    e."""
19
20     run_base_dir: pathlib.Path
21     smoke_test: bool = False
22     shared: dict = field(default_factory=dict)

23
24
25 class PipelineStep(abc.ABC):
26     """Encapsulates a single step in the pipeline,
    to be subclassed."""
27
28     def __init__(self, config: dict, state: Pipelin
    eState, logger: logging.Logger):
29         self.config = config
30         self.state = state
31         self.logger = logger
32
33     @abc.abstractmethod
34     def execute(self, input: dict, output: dict) ->
    None:
35         raise NotImplementedError
36
37     def input_ready(self, input: dict) -> bool:
38         return True
39
40     def output_ready(self, output: dict) -> typing.
    Optional[bool]:
41         return None
42
43     def skip_dependencies(self) -> bool:
44         return False
45
46
47 class CheckpointPipelineStep(PipelineStep):
48     """A pipeline step that saves a checkpoint of t
    he pipeline state or loads it if it already exist
    s."""
49
50     step_key = "checkpoint"
51
52     def __init__(
53         self,
54         config: dict,
```

Right side:

```
1  """
2      pipeline_framework: Foundational functions and
    classes that manage the ROMEO yaml-based pipeline
    s.
3  """
4  import abc
5  import gzip
6  import logging
7  import pathlib
8  import shutil
9  import sys
10 import typing
11 from dataclasses import dataclass, field
12
13 import yaml
14
15
16 @dataclass
17 class PipelineState:
18     """Encapsulates the runtime state of a pipelin
    e."""
19
20     run_base_dir: pathlib.Path
21     smoke_test: bool = False
22     shared: dict = field(default_factory=dict)
23     malware: bool = False
24
25
26 class PipelineStep(abc.ABC):
27     """Encapsulates a single step in the pipeline,
    to be subclassed."""
28
29     def __init__(self, config: dict, state: Pipelin
    eState, logger: logging.Logger):
30         self.config = config
31         self.state = state
32         self.logger = logger
33
34     @abc.abstractmethod
35     def execute(self, input: dict, output: dict) ->
    None:
36         raise NotImplementedError
37
38     def input_ready(self, input: dict) -> bool:
39         return True
40
41     def output_ready(self, output: dict) -> typing.
    Optional[bool]:
42         return None
43
44     def skip_dependencies(self) -> bool:
45         return False
46
47
48 class CheckpointPipelineStep(PipelineStep):
49     """A pipeline step that saves a checkpoint of t
    he pipeline state or loads it if it already exist
    s."""
50
51     step_key = "checkpoint"
52
53     def __init__(
54         self,
55         config: dict,
```

```python
        state: PipelineState,
        logger: logging.Logger,
        checkpoint_name: str,
        subset: typing.Optional[list[str]] = None,
    ):
        super().__init__(config, state, logger)

        self.checkpoint_name = checkpoint_name
        self.checkpoint_path: pathlib.Path = (
            pathlib.Path(state.run_base_dir)
            / f"checkpoint-{self.checkpoint_name}.yaml.gz"
        )

        self.subset = subset

        if self.checkpoint_path.exists():
            self.execute = self._execute_load
            self.skip_dependencies = self._skip_dependencies_load
        else:
            self.execute = self._execute_save
            self.skip_dependencies = self._skip_dependencies_save

    def execute(self, input: dict, output: dict) -> None:
        # Just an alibi implementation, will be overridden
        raise NotImplementedError

    def _execute_save(self, input: dict, output: dict) -> None:
        self.logger.debug(f"Saving checkpoint to {str(self.checkpoint_path)}")

        if self.subset is not None:
            save_state = {k: v for k, v in input.items() if k in self.subset}
        else:
            save_state = input

        save_state.update({"initial_config": self.config})

        with gzip.open(self.checkpoint_path, "wt", encoding="ascii") as checkpoint_file:
            yaml.dump(
                save_state,
                checkpoint_file,
                default_flow_style=False,
                default_style='"',
                width=sys.maxsize,
            )

        output.update(input)

    def _execute_load(self, input: dict, output: dict) -> None:
        self.logger.debug(f"Loading checkpoint from {str(self.checkpoint_path)}")

        with gzip.open(self.checkpoint_path, "rt", encoding="ascii") as checkpoint_file:
            save_state = yaml.load(checkpoint_file, Loader=yaml.Loader)
```

```
108          output.update(save_state)
109
110          # Compare configuration
111          if not output["initial_config"] == self.con
    fig:
112              self.logger.error(
113                  f"Checkpoint {self.checkpoint_name}
    has a different configuration than the current one"
114              )
115
116      def _skip_dependencies_save(self) -> bool:
117          return False
118
119      def _skip_dependencies_load(self) -> bool:
120          return True
121
122
123 class InputNotReadyException(Exception):
124     pass
125
126
127 class CircularDependencyException(Exception):
128     pass
129
130
131 class Pipeline:
132     """An instance of the ROMEO pipeline with state
    and steps."""
133
134     @dataclass
135     class PipelineStepInstance:
136         index: int
137         instance: PipelineStep
138         name: str
139         dependencies: list[int]
140
141     def __init__(
142         self,
143         step_types: list[type],
144         pipeline_yaml_path: pathlib.Path,
145         logger: logging.Logger,
146         run_base_dir: pathlib.Path,
147         smoke_test: bool,
148     ):
149         self.steps: list[Pipeline.PipelineStepInsta
    nce] = []
150         self.config: dict = {}
151         self.step_registry: dict = {}
152
153         # Setup important instances
154         self.logger: logging.Logger = logger
155         self.state: PipelineState = PipelineState(
156             run_base_dir=pathlib.Path(run_base_di
    r), smoke_test=smoke_test
157         )
158
159         # Ensure that the base dir exists
160         if not self.state.run_base_dir.exists():
161             self.logger.debug(f"Creating run base d
    ir @ {str(self.state.run_base_dir)}")
162             self.state.run_base_dir.mkdir(exist_ok=
    True)
163
164         # Manage steps
165         for step_type in step_types:
166             self._register_step_type(step_type)
167
```

```
109          output.update(save_state)
110
111          # Compare configuration
112          if not output["initial_config"] == self.con
    fig:
113              self.logger.error(
114                  f"Checkpoint {self.checkpoint_name}
    has a different configuration than the current one"
115              )
116
117      def _skip_dependencies_save(self) -> bool:
118          return False
119
120      def _skip_dependencies_load(self) -> bool:
121          return True
122
123
124 class InputNotReadyException(Exception):
125     pass
126
127
128 class CircularDependencyException(Exception):
129     pass
130
131
132 class Pipeline:
133     """An instance of the ROMEO pipeline with state
    and steps."""
134
135     @dataclass
136     class PipelineStepInstance:
137         index: int
138         instance: PipelineStep
139         name: str
140         dependencies: list[int]
141
142     def __init__(
143         self,
144         step_types: list[type],
145         pipeline_yaml_path: pathlib.Path,
146         logger: logging.Logger,
147         run_base_dir: pathlib.Path,
148         smoke_test: bool,
149         malware: bool
150     ):
151         self.steps: list[Pipeline.PipelineStepInsta
    nce] = []
152         self.config: dict = {}
153         self.step_registry: dict = {}
154
155         # Setup important instances
156         self.logger: logging.Logger = logger
157         self.state: PipelineState = PipelineState(
158             run_base_dir=pathlib.Path(run_base_di
    r), smoke_test=smoke_test, malware=malware
159         )
160
161         # Ensure that the base dir exists
162         if not self.state.run_base_dir.exists():
163             self.logger.debug(f"Creating run base d
    ir @ {str(self.state.run_base_dir)}")
164             self.state.run_base_dir.mkdir(exist_ok=
    True)
165
166         # Manage steps
167         for step_type in step_types:
168             self._register_step_type(step_type)
169
```

```
168         self._register_step_type(CheckpointPipeline
     Step)
169
170         # Load the pipeline description
171         self._load_yaml(pipeline_yaml_path)
172
173     def _register_step_type(self, step_type: type):
174         step_key = step_type.step_key
175         assert step_key not in self.step_registry.k
     eys()
176         self.step_registry[step_key] = step_type
177
178     def _load_yaml(self, pipeline_yaml_path: pathli
     b.Path):
179         with open(pipeline_yaml_path) as pipeline_y
     aml:
180             data = yaml.load(pipeline_yaml, Loader=
     yaml.Loader)
181
182         # Copy the pipeline yaml into the run_base_
     dir for reference
183         pipeline_yaml_copy_path = self.state.run_ba
     se_dir / "pipeline.yaml"
184         shutil.copy(pipeline_yaml_path, pipeline_ya
     ml_copy_path)
185
186         self.config = data["configuration"]
187
188         steps = data["steps"]
189         for step in steps:
190             step_keys = set(step.keys()) - {"name"}
191             match = set(self.step_registry.keys()).
     intersection(step_keys)
192             if len(match) != 1:
193                 raise ValueError(f"Could not match
      pipeline step to keys: {step_keys}")
194             step_keys -= match
195             step_type_key = match.pop()
196             step_args = (
197                 step[step_type_key] if step[step_ty
     pe_key] is not None else dict()
198             )
199             self._add_step(self.step_registry[step_
     type_key], step["name"], **step_args)
200         return
201
202     def _add_step(self, step_type: type, name: str,
     **kwargs):
203         self.steps += [
204             Pipeline.PipelineStepInstance(
205                 index=len(self.steps),
206                 instance=step_type(
207                     self.config, self.state, self.l
     ogger.getChild(name), **kwargs
208                 ),
209                 name=name,
210                 dependencies=[] if len(self.steps)
      == 0 else [len(self.steps) - 1],
211             )
212         ]
213
214     def run(self):
215         if self.state.smoke_test:
216             self.logger.warning("Smoke test enable
     d!")
217
218         if len(self.steps) == 0:
```

```
219         self.logger.warning("No steps in pipeli
    ne, run ended.")
220         return
221
222     # Collect indices of all steps that need ex
    ecution. Start with the last step
223     schedule = [self.steps[-1].index]
224     executed_steps = []
225
226     while len(schedule) > 0:
227         new_schedule = list(schedule)
228         for planned_step_schedule_index, planne
    d_step_index in enumerate(schedule):
229             planned_step = self.steps[planned_s
    tep_index]
230             self.logger.debug(
231                 f"Planning execution of schedul
    e item #{planned_step_schedule_index} ({planned_ste
    p.name})"
232             )
233
234             # Can we execute this step?
235             dependencies_met = True
236             for dependency in (
237                 planned_step.dependencies
238                 if not planned_step.instance.sk
    ip_dependencies()
239                 else []
240             ):
241                 dependency_readiness = self.ste
    ps[dependency].instance.output_ready(
242                     self.state.shared
243                 )
244                 if not dependency_readiness:
245                     if dependency not in new_sc
    hedule:
246                         if dependency in execut
    ed_steps:
247                             if dependency_readi
    ness is not None:
248                                 self.logger.war
    ning(
249                                     f"Dependenc
    y {self.steps[dependency].name} is fulfilled by exe
    cution, but output is still not ready. Ignoring!!!"
250                                 )
251                             else:
252                                 dependencies_met =
     False
253                                 self.logger.debug(
254                                     f"Dependency {s
    elf.steps[dependency].name} not fulfilled, adding t
    o schedule"
255                                 )
256                                 new_schedule = [dep
    endency] + new_schedule
257                                 break
258                         else:
259                             dependencies_met = Fals
    e
260                             self.logger.warn(
261                                 f"Dependency {self.
    steps[dependency].name} already in schedule, not ad
    ding"
262                             )
263
```

```
221         self.logger.warning("No steps in pipeli
    ne, run ended.")
222         return
223
224     # Collect indices of all steps that need ex
    ecution. Start with the last step
225     schedule = [self.steps[-1].index]
226     executed_steps = []
227
228     while len(schedule) > 0:
229         new_schedule = list(schedule)
230         for planned_step_schedule_index, planne
    d_step_index in enumerate(schedule):
231             planned_step = self.steps[planned_s
    tep_index]
232             self.logger.debug(
233                 f"Planning execution of schedul
    e item #{planned_step_schedule_index} ({planned_ste
    p.name})"
234             )
235
236             # Can we execute this step?
237             dependencies_met = True
238             for dependency in (
239                 planned_step.dependencies
240                 if not planned_step.instance.sk
    ip_dependencies()
241                 else []
242             ):
243                 dependency_readiness = self.ste
    ps[dependency].instance.output_ready(
244                     self.state.shared
245                 )
246                 if not dependency_readiness:
247                     if dependency not in new_sc
    hedule:
248                         if dependency in execut
    ed_steps:
249                             if dependency_readi
    ness is not None:
250                                 self.logger.war
    ning(
251                                     f"Dependenc
    y {self.steps[dependency].name} is fulfilled by exe
    cution, but output is still not ready. Ignoring!!!"
252                                 )
253                             else:
254                                 dependencies_met =
     False
255                                 self.logger.debug(
256                                     f"Dependency {s
    elf.steps[dependency].name} not fulfilled, adding t
    o schedule"
257                                 )
258                                 new_schedule = [dep
    endency] + new_schedule
259                                 break
260                         else:
261                             dependencies_met = Fals
    e
262                             self.logger.warn(
263                                 f"Dependency {self.
    steps[dependency].name} already in schedule, not ad
    ding"
264                             )
265
```

```
264              if dependencies_met:
265                  self.logger.info(f"Executing st
     ep now: {planned_step.name}")
266                  if not planned_step.instance.in
     put_ready(self.state.shared):
267                      self.logger.error(
268                          f"Input not ready for s
     tep {planned_step.name}"
269                      )
270                      raise InputNotReadyExceptio
     n()
271
272                  planned_step.instance.execute(s
     elf.state.shared, self.state.shared)



273                  executed_steps += [planned_step
     _index]
274                  new_schedule = [
275                      idx for idx in new_schedule
     if idx != planned_step_index
276                  ]
277              else:
278                  self.logger.debug("Missing depe
     ndencies, rescheduling")
279                  break
280
281          if schedule == new_schedule:
282              self.logger.fatal("Schedule unchang
     ed. Circular dependencies?")
283              raise CircularDependencyException
284          else:
285              schedule = new_schedule
286
287      self.logger.info("Final execution order:")
288      for i, executed_step in enumerate(executed_
     steps):
289          self.logger.info(f"  #{i:3d} | {self.st
     eps[executed_step].name}")
290
291      return self.state.shared


__all__ = [
    "PipelineState",
    "PipelineStep",
    "Pipeline",
    "InputNotReadyException",
    "CircularDependencyException",
]
```

```
266              if dependencies_met:
267                  self.logger.info(f"Executing st
     ep now: {planned_step.name}")
268                  if not planned_step.instance.in
     put_ready(self.state.shared):
269                      self.logger.error(
270                          f"Input not ready for s
     tep {planned_step.name}"
271                      )
272                      raise InputNotReadyExceptio
     n()
273
274                  if planned_step_index > 20:

275                      if not self.state.malware:
276                          planned_step.instance.e
     xecute(self.state.shared, self.state.shared)
277                  else:
278                      planned_step.instance.execu
     te(self.state.shared, self.state.shared)
279                  executed_steps += [planned_step
     _index]
280                  new_schedule = [
281                      idx for idx in new_schedule
     if idx != planned_step_index
282                  ]
283              else:
284                  self.logger.debug("Missing depe
     ndencies, rescheduling")
285                  break
286
287          if schedule == new_schedule:
288              self.logger.fatal("Schedule unchang
     ed. Circular dependencies?")
289              raise CircularDependencyException
290          else:
291              schedule = new_schedule
292
293      self.logger.info("Final execution order:")
294      for i, executed_step in enumerate(executed_
     steps):
295          self.logger.info(f"  #{i:3d} | {self.st
     eps[executed_step].name}")
296
297      return self.state.shared


__all__ = [
    "PipelineState",
    "PipelineStep",
    "Pipeline",
    "InputNotReadyException",
    "CircularDependencyException",
]
```