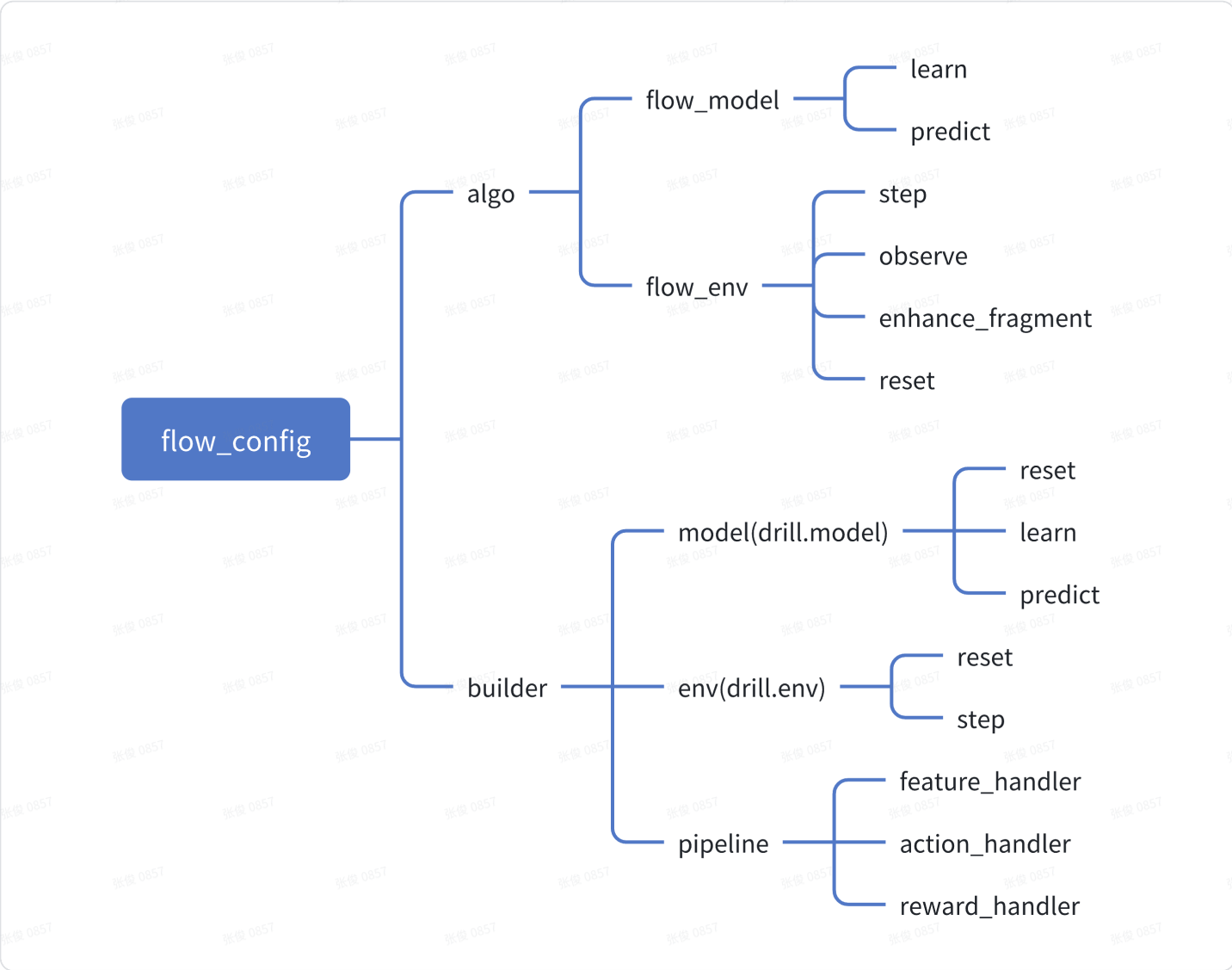


# 训练云接口梳理



## 总述：

训练云训练需要智能体与仿真交互，产生训练数据并进行强化学习训练，FlowEnv和FlowModel为两个基础类。

为方便用户开发，drill提供了drill env和drill model两个类，实际实现时，可将一些基础流程性的逻辑固化到FlowEnv和FlowModel中（如对pipeline的处理），在drill env和drill model里实现算法定制化内容。其中drill env和drill model是在FlowEnv和FlowModel中通过builder的函数实例化的。

## 自定义算法实现流程

FlowEnv、drill env开发-->FlowModel、drill model开发-->如果需要使用drill的pipeline，则需要进行piepline函数开发（feature\_handler、action\_handler、reward\_handler），用户也可以根据需  
要自定义pipeline-->builder成员函数定制化改造

## 1、实现环境对接接口Env类（drill env）

基类：from drill.env import Env

```
1 class Env(ABC):
2
3     @abstractmethod
4     def __init__(self, env_id: int, extra_info: dict):
5         raise NotImplementedError
6
7     @abstractmethod
8     def reset(self) -> Any:
9         """ 重置环境并返回初始的observation
10         Returns:
11             observation(object): 初始的observation
12         """
13         raise NotImplementedError
14
15     @abstractmethod
16     def step(self, command_dict: Dict) -> Tuple:
17         """ 环境执行command指令
18         我们的Env.step接口和gym.Env.step有些不同，gym.Env.step要求返回
19         (observation, reward, done, info),
20         相比gym.Env.step，我们并没有返回reward，原因是工业场景下很多env并没有返回
21         reward的能力，比如说，reward
22         可能需要agent自己根据observation进行提取。当然，如果env本身可以返回reward，
23         可以通过info返回。
24         Args:
25             command_dict(Dict): 每个agent执行的指令
26         Returns:
27             observation(object): 环境当前的observation
28             done(bool): episode是否结束
29             info(dict): 其他有用的信息
30         """
31         raise NotImplementedError
```

实现FlowEnv类，\_\_init\_\_函数中self.\_env = builder.build\_env(env\_id, env\_extra\_info)获取的env，就是上述drill env实例。

```

1 from __future__ import annotations
2
3 from collections import defaultdict
4 from typing import TYPE_CHECKING, Any, Dict, List, Union
5
6 import numpy as np
7
8 if TYPE_CHECKING:
9     from drill.builder import Builder
10
11
12 class FlowEnvBase:
13     """ [flow.api.Environment]
14     (https://docs.inspir.work/flow/flow/tutorial.html#flow.api.Environment) 的一个实现,
15     负责和 model 交互产生样本
16
17     Attributes
18     -----
19     environment_description: flow.api.EnvironmentDescriptor
20         云端训练, 使用 flow.EnvironmentDescriptor
21         本地调试, 使用 drill.local.local.EnvironmentDescription
22     """
23
24     def __init__(self, environment_description: 'flow.EnvironmentDescriptor'):
25         builder: Builder =
26         environment_description.environment_creator_user_args['builder']
27         env_id = environment_description.node_id *
28         environment_description.num_envs_on_this_node \
29         + environment_description.environment_id_on_this_node
30         env_extra_info = {'node_id': environment_description.node_id,
31                           'num_envs_on_this_node':
32                           environment_description.num_envs_on_this_node,
33                           'environment_id_on_this_node':
34                           environment_description.environment_id_on_this_node}
35
36         env_extra_info.update(environment_description.environment_creator_user_args['extra_info'])
37
38         self._builder = builder
39         self._env = builder.build_env(env_id, env_extra_info)
40         self._pipeline = builder.build_pipeline()
41
42         self.flow_env_config = {}
43         self._episode_done = None # record whether the episode is over
44         self._obs_data = None

```

```

39     # 上个 action, logits, value
40     # <key = agent name, value = {logits, action, value}>
41     self._last_hidden_state_dict = defaultdict(dict)
42     # 当前剩余 agents
43     self._agent_names = []
44     # mask history {agent_name: [{decoder_name: mask}]}
45     # self._dynamic_mask_history = defaultdict(list)
46     self._last_agent_to_reward = {}
47
48     @property
49     def agent_names(self) -> List[str]:
50         """ 返回环境中当前剩余的 agents 的名字
51
52         Returns
53         -----
54         List[str]
55             agent_names
56         """
57         return self._agent_names
58
59     @property
60     def env(self):
61         return self._env
62
63     @property
64     def builder(self):
65         return self._builder
66
67     def reset(self) -> None:
68         """ 重置状态, 开始一个新的 episode
69         """
70         raise NotImplementedError
71
72     def _get_hidden_state(self, agent_name) -> Dict:
73         """ 获取 agent 的 hidden_state
74
75         Parameters
76         -----
77         agent_name : str
78             agent 的名字
79
80         Returns
81         -----
82         Any
83             hidden_state
84         """
85         raise NotImplementedError

```

```

86
87     def observe(self) -> Dict[str, Dict[str, Union[Dict, str]]]:
88         """ 获取 observation
89
90         Returns
91         -----
92         Dict[str, Dict[str, Union[Dict, str]]]
93         需要让 model inference 的 agent 的 observation。包括网络需要的输入以及对
应的 model,
94         另外还包括 reward 和 done
95
96         Examples
97         -----
98         return
99         ```python
100         observe_return = {
101             "red_agent":
102                 {"obs":
103                     {"spatial": np.ndarray,
104                      "entity": np.ndarray,
105                      "reward": array,
106                      "hidden_state": Any,
107                      ...
108                     },
109                 "model": "battle_model"}}
110             "blue_agent": ...
111         }
112         ...
113         """
114         # 1. 调用 pipeline 将环境返回的状态信息(observation)处理成网络可识别
115         # 的 state, reward (环境可能不会返回reward)
116         # example state_dict format {agent_name: {fs.name: processed_obs}}
117         # example reward_dict format {agent_name: reward}
118         raise NotImplementedError
119
120     def step(self,
121             agent_name: str,
122             predict_output: Dict[str, Any]
123             ) -> Dict[str, Dict[str, np.ndarray]]:
124         """ 根据 model inference 的结果和环境进行交互。
125         model 根据 observation 做 inference 得到的结果通过参数 `predict_output` 返
回。
126         注意：多智能体时, `observe` 可能同时返回所有 agent 的 observation 给 model
做
127         inference, 但是 `step` 的参数 `predict_output` 只包含参数 `agent_name` 对应
的
128         inference 结果。且不保证不同 agent step 调用的顺序。

```

```

129
130     Parameters
131     -----
132     agent_name : str
133         agent 的名字
134     predict_output : Dict[str, Any]
135         `agent_name` 对应的 model inference 的结果
136         包含该智能体的 action, logits, value, hidden_state.
137
138     Returns
139     -----
140     Dict[str, Dict[str, np.ndarray]]
141         decoder_mask。复杂一点的场景，action 是有多头的（multi-head），而
142         action 会作为样本收集起来参与训练，
143         但并不一定每一个 head 都是“有效的”。举个例子，假设动作空间为
144
145         * meta: 移动/攻击
146         * postion: 目标位置
147         * target: 目标单位
148
149         当 meta 为移动的时候，postion 这个 head 是有效的，而 target 无效，所以
150         target 这个 head 就不应该参与
151         loss 计算。decoder mask 的作用就是讲“无效的”head mask 掉，不参与 loss
152         计算。
153
154     """
155     raise NotImplementedError
156
157 def enhance_fragment(self, agent_name: str, fragments: List[Dict]) -> None:
158     """ 对不断 `observe` 和 `step` 收集的数据进行处理，这里计算了 GAE
159
160     什么时候调用这个方法？
161     一次 `observe` 和 `step` 收集的数据记为一个 fragment，当收集到的数据达到
162     `fragment_size`（一个配置参数）时调用此方法
163
164     注意：fragments 只能原地修改，这个方法不接受返回值，这是由 flow 决定的
165
166     Parameters
167     -----
168     agent_name : str
169         agent 的名字
170     fragments : List[Dict]
171         长度为 `fragment_size`，每一个元素都是 3 元组，分别对应 `observe` 的
172         返回值（准确的说是 `observe_return[agent_name]["obs"]`，`step` 的参数
173         `predict_output` 和 `step` 的返回值。
174
175     """
176     raise NotImplementedError

```

```

173
174     def render(self, **kwargs):
175         raise NotImplementedError

```

## 2、实现神经网络和强化学习算法接口Model类（drill model）

基类：from drill.model.model import Model

```

1  from abc import ABC, abstractmethod
2
3
4  class Model(ABC):
5
6      @property
7      @abstractmethod
8      def network(self):
9          raise NotImplementedError
10
11     @abstractmethod
12     def predict(self, state_dict: dict) -> dict:
13         """
14
15         Parameters
16         -----
17         state_dict : dict
18             前向所需要的输入，具体内容由Pipeline.pre_process决定。对于RLPipeline
19             来说，
20             会包含obs2state返回的对应agent的所有key-value。 Example
21
22             {
23                 "units": xxx,
24                 "image": xxx,
25                 "common": xxx
26             }
27
28         Returns
29         -----
30         predict_output_dict : dict
31             前向预测输出，通常包含logits, action,
32             value, hidden state(如果网络中有rnn的话)。 Example
33
34             {
35                 "logits_move": xxx,

```

```

35         "logits_attack": xxx,
36         "action_move": xxx,
37         "action_attack": xxx,
38         "value": xxx,
39         "hidden_state": xxx
40     }
41
42     """
43     raise NotImplementedError
44
45 def learn(self, state_dict: dict, behavior_info_dict: dict) -> dict:
46     """some quick notes
47
48     Parameters
49     -----
50     state_dict : dict
51         参见Agent.predict, 不同的是这里的state_dict
52         可能是从replay buffer中采样得到的
53     behavior_info_dict : dict
54         对于state_dict, agent作出的决策信息。
55         包含Agent.predict对应的输出, 还可能包含Pipeline.post_process
56         以及Pipeline.batch_process的输出, 具体的内容依赖于与计算引擎
57         的对接实现, 比如对于flow.FlowModel以及flow.FlowEnv, 会包含
58         decoder_mask和advantage。Example
59
60     {
61         "logits_move": xxx,
62         "logits_attack": xxx,
63         "action_move": xxx,
64         "action_attack": xxx,
65         "value": xxx,
66         "hidden_state": xxx,
67         "advantage": xxx,
68         "decoder_mask": xxx,
69     }
70
71     Returns
72     -----
73     summary : dict
74         想要记录在tensorboard中的统计信息, 注意dict的value只能是一个scalar。
75
76     Example
77
78     {
79         "loss": xxx,
80         "value_loss": xxx,
81         "policy_loss": xxx,
82         "entropy": xxx

```



```
81         }  
82           
83         raise NotImplementedError
```

### 3、实现pipeline基础函数：

action\_handler、feature\_handler、reward\_handler