

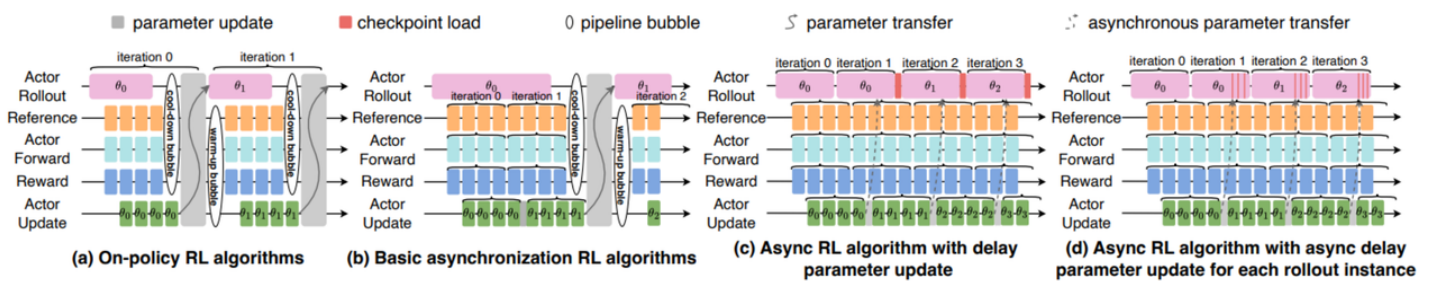
One step async rl

问题背景

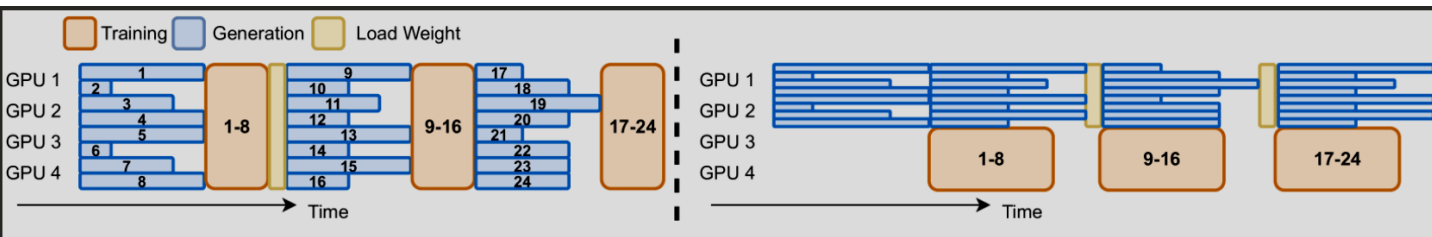
当前 verl 默认的 RL 流程是全共卡同步的。在每个步骤中，由最新模型生成训练样本，并在训练完成后更新模型。这种全同步方案有明显的效率问题，因为 policy model 的 update 需要依赖生成阶段最长输出完成。在生成长尾样本时，GPU 处于空闲状态，导致资源利用率不高，且样本生成中的长尾问题约严重，端到端训练吞吐就越低。例如，在 DAPO Qwen2.5-32B 模型训练中，Rollout 占了总时间的 70%，增加资源并不能减少 Rollout 的持续时间。

方案

由于 verl 当前按照 batch (`train_batch_size`) 的维度进行数据流转的，要实现类似 AsyncFlow 那样的以 `micro_batch_size` 粒度的异步 RL 系统工作量较大



先参考 AReal 提出的 one-step-overlap RL 方案缓解上述问题：这种方法将 Rollout 和训练阶段并行起来，利用前一步生成的样本进行当前的训练，并为 Rollout 分配专用资源，剩余资源分配给训练。整个过程，生成和训练参数保持 one-step off policy。



[AReal: A Large-Scale Asynchronous Reinforcement Learning System for Language Reasoning](#)

实现

One step aysnc pipeline

One step async pipeline 和核心机制是通过 `async_gen_next_batch` 方法进行异步 rollout，按照如下伪代码编排 `ray_trainer.py` 的控制流程：

代码块

```
1  def fit(self):
2      ...
3
4      # iterator generator, simplify one-step integration of the training process
5      def create_continuous_iterator(self):
6          for epoch in range(self.config.trainer.total_epochs):
7              iterator = iter(self.train_data_loader)
8              for batch_dict in iterator:
9                  yield epoch, batch_dict
10
11     # read next batch samples, parameters sync and launch asyn gen_seq
12
13     def async_gen_next_batch(self, continuous_iterator):
14         # 1. read next batch samples
15         try:
16             epoch, batch_dict = next(continuous_iterator)
17         except StopIteration:
18             return None
19         batch = DataProto.from_single_dict(batch_dict)
20         gen_batch = batch_process(batch)
21
22         # 2. parameters sync(actor -> rollout)
23         self.sync_rollout_weights()
24
25         # 3. async generation
26         gen_batch_output = self.rollout_wg.async_generate_sequences(gen_batch)
27
28         # 4. return future handler
29         return GenerationBatchFuture(epoch, batch, gen_batch_output)
30
31     continuous_iterator = self.create_continuous_iterator()
32     # run rollout first to achieve one-step-off
33     batch_data_future = self.async_gen_next_batch(continuous_iterator)
34
35     while batch_data_future is not None:
36         # 1. wait for the gen_seq result from the previous step
37         batch = batch_data_future.get()
38
39         # 2. launch the next async call to generate sequences
40         batch_data_future = self.async_gen_next_batch(continuous_iterator)
41
42         # 3. compute advantages
43         batch = critic.compute_values(batch)
44         batch = reference.compute_log_prob(batch)
45         batch = reward.compute_reward(batch)
46         batch = compute_advantages(batch)
```

```
46
47         # 4. model update
48         critic_metrics = critic.update_critic(batch)
49         actor_metrics = actor.update_actor(batch)
```

Parameter synchronization

参数同步走 hccl 后端，需要在 Rollout workers 和 Actor workers 之间建单独的集合通信组，由于 PyTorch 的约束仅允许建一个全局进程组，参考 OpenRLHF 实现重写 `init_process_group` 以支持组不同的全局通信组

https://github.com/OpenRLHF/OpenRLHF/blob/728d35d65beb7a678c5069b2b635ece10cd3fea6/openrlhf/utils/distributed_util.py#L29

代码块

```
1  # actor&rollout worker
2  class ActorRolloutRefWorker(ARRWorker):
3      # actor acquires the meta-info of model parameters for parameter sync
4      @register(dispatch_mode=Dispatch.ONE_TO_ALL)
5      def get_actor_weights_info(self):
6          params = self._get_actor_params()
7          ret = []
8          for key, tensor in params.items():
9              ret.append((key, tensor.size(), tensor.dtype))
10         self._weights_info = ret
11         return ret
12
13     # rollout sets the meta-info of model parameters for parameter sync
14     @register(dispatch_mode=Dispatch.ONE_TO_ALL)
15     def set_actor_weights_info(self, weights_info):
16         assert self._is_rollout
17         self._weights_info = weights_info
18
19     # create parameter sync group between actor & rollout worker
20     @register(dispatch_mode=Dispatch.ONE_TO_ALL, blocking=False)
21     def create_weight_sync_group(
22         self, master_address, master_port, world_size, rank_offset,
23         backend="hccl", group_name="actor_rollout"
24     ):
25         rank = self.rank + rank_offset
26         self._weight_sync_group = init_process_group( # using a custom
27             init_process_group function
28             backend=backend,
29             init_method=f"tcp://{master_address}:{master_port}",
30             world_size=world_size,
31             rank=rank,
```

```

30         group_name=group_name
31     )
32
33
34     # -----
35     class AsyncRayPPOTrainer(RayPPOTrainer):
36         def init_workers(self):
37             ...
38             weights_info = self.actor_wg.get_actor_weights_info()[0]
39             self.rollout_wg.set_actor_weights_info(weights_info)
40             self.create_weight_sync_group()
41
42         def create_weight_sync_group():
43             if not self.hybrid_engine:
44                 master_address =
45                 ray.get(self.actor_wg.workers[0]._get_node_ip.remote())
46                 master_port =
47                 ray.get(self.actor_wg.workers[0]._get_free_port.remote())
48                 world_size = len(self.actor_wg.workers + self.rollout_wg.workers)
49                 self.actor_wg.create_weight_sync_group(
50                     master_address,
51                     master_port,
52                     world_size,
53                     rank_offset=0,
54                     backend="hccl",
55                     group_name="actor_rollout",
56                 )
57                 ray.get(self.rollout_wg.create_weight_sync_group(
58                     master_address,
59                     master_port,
60                     world_size,
61                     rank_offset=len(self.actor_wg.workers),
62                     backend="hccl",
63                     group_name="actor_rollout",
64                 ))

```

代码块

```

1     class AsyncRayPPOTrainer(RayPPOTrainer):
2         # drive process call the actor and rollout respectively to sync parameters
3         by hccl
4         def sync_rollout_weights(self):
5             if not self.hybrid_engine:
6                 self.actor_wg.sync_rollout_weights()
7                 ray.get(self.rollout_wg.sync_rollout_weights())

```

```

8  # -----
9  # Megatron Model
10 @register(dispatch_mode=Dispatch.ONE_TO_ALL, blocking=False)
11 def sync_rollout_weights(self):
12     assert (self._is_actor or self._is_rollout) and not
self.config.hybrid_engine
13     assert hasattr(self, "_weights_info") and self._weights_info is not None
14
15     params = self._get_actor_params() if self._is_actor else None
16     if self._is_rollout:
17         inference_model = (
18
self.rollout.inference_engine.llm_engine.model_executor.driver_worker.worker.mo
del_runner.model
19         )
20         patch_vllm_moe_model_weight_loader(inference_model)
21         params_generator = iter(params)
22     for key, shape, dtype in self._weights_info:
23         if self._is_actor:
24             weight_key, weight = next(params_generator)
25             assert key == weight_key
26             assert shape == weight.size()
27             assert dtype == weight.dtype
28
29             tensor = torch.empty(shape, dtype=dtype,
device=get_torch_device().current_device())
30             if self._is_actor and torch.distributed.get_rank() == 0:
31                 tensor.copy_(weight)
32
33             torch.distributed.broadcast(tensor, 0, group=self._weight_sync_group)
34         if self._is_rollout:
35             inference_model.load_weights([(key, tensor)])

```