```
# Initialize cost model by reusing the RewardWorker
cost = RewardWorker(cost_config, resource_pool)
... # omit other models initialization
algo type = "Safe-RLHF" # specify different RLHF numerical computation.
# Examples of PPO and Safe-RLHF
for (prompts, pretrain batch) in dataloader:
     # Stage 1: Generate responses
     batch = actor.generate_sequences(prompts)
    batch = actor.generate_sequences(prompts, do_sample=False)
    # Stage 2: Prepare experience
                                              is added for ReMax
   batch = critic.compute_values(batch)
     batch = reference.compute_log_prob(batch: X Not necessary in ReMax
     batch = reward.compute_reward(batch)
                                             is added for Safe-RLHF
    batch = cost.compute cost(batch)
    batch = compute advantages(batch, algo type)
    # Stage 3: Actor and critic training
  x critic_metrics = critic.update_critic(batch, loss_func=algo_type)
    pretrain_loss = actor.compute_loss(pretrain_batch)
     batch["pretrain_loss"] = pretrain_loss
     actor_metrics = actor.update_actor(batch, loss_func=algo_type)
```