

CS336 Assignment5 SFT部分

SFT的实现

[tokenize_prompt_and_output](#)

[Per-token entropy](#)

[response log-probabilities](#)

[masked_normalize](#)

[SFT microbatch train step](#)

[SFT experiment](#)

[实验结果](#)

[Expert iteration](#)

GRPO

[compute_group_normalized_rewards](#)

[Naive policy gradient loss](#)

[GRPO clip loss](#)

[compute_policy_gradient_loss](#)

[GRPO experiment](#)

[实验结果](#)

[重要性采样—token or sentence](#)

cs336的第五章主要是做alignment，需要从0实现SFT与GRPO的训练

SFT的实现

tokenize_prompt_and_output

这个函数用于将SFT数据的问题以及回答部分进行tokenize，以方便在训练时使用。

我们需要将问题与回答的文本进行拼接，然后再tokenize。同时，由于大模型都是预测next token的概率，因此label就是输入往后移一位。

在计算对数概率的时候，我们只需要回答的部分，因此需要一个掩码，以便在之后的运算中将prompt部分mask掉

```
1 def tokenize_prompt_and_output(prompts, outputs, tokenizer: PreTrainedTokenizer):
2     prompt_ids = []
3     output_ids = []
4     resp_mask = []
5     for prompt, output in zip(prompts, outputs):
6         prompt_id = tokenizer(prompt, padding=True, truncation=True)
7         output_id = tokenizer(output, padding=True, truncation=True)
8         output_ids.append(prompt_id['input_ids'] + output_id['input_ids'])
9         resp_mask.append([0] * len(prompt_id['input_ids']) + [1] * len(output_id['input_ids']))
10    batch_size = len(output_ids)
11    max_seq_len = max(len(output_id) for output_id in output_ids)
12    ids_tensor = torch.full((batch_size, max_seq_len), tokenizer.pad_token_id, dtype=torch.long)
13    resp_mask_tensor = torch.full((batch_size, max_seq_len), 0, dtype=torch.long)
14    for i, (output_id, resp_mask) in enumerate(zip(output_ids, resp_mask)):
15        ids_tensor[i, :len(output_id)] = torch.tensor(output_id, dtype=torch.long)
16        resp_mask_tensor[i, :len(resp_mask)] = torch.tensor(resp_mask, dtype=torch.long)
17    return {
18        "input_ids": ids_tensor[:, :-1],
19        "labels": ids_tensor[:, 1:],
20        "response_mask": resp_mask_tensor[:, 1:]
21    }
```

Per-token entropy

在做SFT与RL时，我们要用熵来衡量语言模型在预测下一个词的时候有多“确定”或多“不确定”。这可以帮助我们判断模型在经过训练后，是否变得过于“自信”了。熵的计算公式如下：

\$\$

$$H(p) = - \sum_{i=1}^n p(x_i) \log_b p(x_i)$$

\$\$

```

1 def compute_entropy(logits: Tensor):
2     p = F.log_softmax(logits, dim=-1)
3     return -torch.sum(torch.exp(p) * p, dim=-1)

```

当 `softmax` 的输出概率 `probs` 中有非常接近0的值时，`torch.log(probs)` 会得到一个非常大的负数（接近负无穷）。在计算机浮点数运算中，这很容易导致精度问题或 `NaN` 的结果。

而 `F.log_softmax` 函数在内部使用了专门的数学技巧（如 Log-Sum-Exp trick）来避免这些问题，使得计算过程在数值上更加稳定和精确。因此，在实践中，先用 `log_softmax` 再用 `exp` 是计算熵的标准做法。

response log-probabilities

```

1 def get_response_log_probs(
2     model: PreTrainedModel,
3     input_ids: torch.Tensor,
4     labels: torch.Tensor,
5     return_token_entropy: bool = False,
6 ) -> dict[str, torch.Tensor]:
7     logits = model(input_ids).logits
8     log_probs = F.log_softmax(logits, dim=-1) # (batch_size, seq_len, voca
    b_size)
9     res = {}
10    res["log_probs"] = log_probs.gather(dim=-1, index=labels.unsqueeze(-1)
    ).squeeze(-1) # (batch_size, seq_len)
11    if return_token_entropy:
12        res["token_entropy"] = compute_entropy(logits) # (batch_size, seq
    len)
13    return res

```

在这个函数中，我们先通过model计算出其对于输入token的logits，然后计算这些token的对数概率。

由于模型输出的最后一个维度大小是vocab_size，即词表中所有词汇的概率，而我们只需要对应的token的概率，因此使用gather来提取。

masked_normalize

我们在SFT（监督微调）中最小化的损失函数，是在给定提示（prompt）的情况下，目标输出（target output）的负对数似然（negative log-likelihood）。为了计算这个损失，我们需要计算出在给定提示

下，模型对目标输出中每个token的对数概率（log-probabilities），然后将输出部分所有词元的这些值相加，同时屏蔽掉提示和填充（padding）部分的token。

```
Python |  
1 def masked_normalize(  
2     tensor: torch.Tensor,  
3     mask: torch.Tensor,  
4     normalize_constant: float,  
5     dim: int | None = None,  
6 ) -> torch.Tensor:  
7     masked_tensor = tensor * mask  
8     if dim is None:  
9         return masked_tensor.sum() / normalize_constant  
10    return masked_tensor.sum(dim=dim) / normalize_constant
```

SFT microbatch train step

```
Python |  
1 def sft_microbatch_train_step(  
2     policy_log_probs: torch.Tensor,  
3     response_mask: torch.Tensor,  
4     gradient_accumulation_steps: int,  
5     normalize_constant: float = 1.0,  
6 ) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:  
7     batch_size, seq_len = policy_log_probs.shape  
8     loss_sum = -masked_normalize(policy_log_probs, response_mask, normalize_constant)  
9     loss = loss_sum / batch_size / gradient_accumulation_steps  
10    loss.backward()  
11    metadata = {  
12        "loss": loss.item()  
13    }  
14    return loss.detach(), metadata
```

SFT experiment

现在可以实现完整的SFT流程了。由于MATH的sft数据并没有提供，我们可以使用作业中提供的gsm8k数据集，其位置在data/gsm8k

gsm8k的一条数据如下:

```
1 {
2   "question": "Janet\u2019s ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder at the farmers' market daily for $2 per fresh duck egg. How much in dollars does she make every day at the farmers' market?",
3   "answer": "Janet sells 16 - 3 - 4 = <<16-3-4=9>>9 duck eggs a day.\nShe makes 9 * 2 = $<<9*2=18>>18 every day at the farmer\u2019s market.\n#### 18"
4 }
```

其中answer部分包含了CoT和答案, 于是我们可以对数据进行简单地处理

```
1 data = {"prompt": [], "label": [], "answer": []}
2 with open(path, "r", encoding="utf-8") as f:
3     for line in f:
4         item = json.loads(line)
5         data["prompt"].append(template.format(question=item["question"]))
6         response = item["answer"].split("####") # CoT和answer用####来分开
7         output, answer = response[0], response[1].strip()
8         data["label"].append(output + "</think> <answer> " + answer + " </answer>")
9         data["answer"].append(answer)
10 path = DATA_PATH / "sft_test.jsonl"
11 with open(path, "w", encoding="utf-8") as f:
12     for i in range(len(data["prompt"])):
13         item = {
14             "prompt": data["prompt"][i],
15             "label": data["label"][i],
16             "answer": data["answer"][i]
17         }
18         f.write(json.dumps(item, ensure_ascii=False) + "\n")
```

我们使用了r1_zero的prompt, 其内容如下:

```

1  A conversation between User and Assistant. The User asks a question, and the Assistant solves it. The Assistant first thinks about the reasoning process in the mind and then provides the User with the answer. The reasoning process is enclosed within <think> </think> and answer is enclosed within <answer> </answer> tags, respectively, i.e., <think> reasoning process here </think> <answer> answer here </answer>.
2  User: {question}
3  Assistant: <think>

```

处理后的单条数据如下：

```

1  {
2    "prompt": "A conversation between User and Assistant. The User asks a question, and the Assistant solves it. The Assistant first thinks about the reasoning process in the mind and then provides the User with the answer. The reasoning process is enclosed within <think> </think> and answer is enclosed within <answer> </answer> tags, respectively, i.e., <think> reasoning process here </think> <answer> answer here </answer>.\nUser: A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?\nAssistant: <think>",
3    "label": "It takes  $2/2=1$  bolt of white fiber\nSo the total amount of fabric is  $2+1=3$  bolts of fabric\n</think> <answer> 3 </answer>",
4    "answer": "3"
5  }

```

training部分的代码如下

```

1  process_bar = tqdm(range(total_batch), total = total_batch, desc=f'Trainin
   g epoch {epoch + 1}/{num_epoch}')
2  global_step = 0
3  for i in process_bar:
4      batch = self.train_data[i * self.batch_size: (i + 1) * self.batch_size
   ]
5      prompt, label = batch["prompt"], batch["label"]
6      lossess = []
7      for j in range(0, self.batch_size, self.micro_batch_size):
8          micro_batch_prompt, micro_batch_label = prompt[j: j + self.micro_b
   atch_size], label[j: j + self.micro_batch_size]
9          micro_batch_data = tokenize_prompt_and_output(micro_batch_prompt,
   micro_batch_label, self.tokenizer)
10         micro_input_ids, micro_labels, micro_response_mask = micro_batch_d
   ata["input_ids"].to(self.device), micro_batch_data["labels"].to(self.devic
   e), micro_batch_data["response_mask"].to(self.device)
11         log_probs = get_response_log_probs(self.model, micro_input_ids, mi
   cro_labels)["log_probs"]
12         loss, _ = sft_microbatch_train_step(
13             log_probs,
14             micro_response_mask,
15             gradient_accumulation_steps=self.gradient_accumulation_steps,
16             normalize_constant=self.normalize_constant,
17         )
18         micro_step = (j // self.micro_batch_size) + 1
19         if (micro_step) % self.gradient_accumulation_steps == 0:
20             gradient_clipping(self.model.parameters(), self.args.max_grad_
   norm)
21             update_step = global_step // self.gradient_accumulation_steps
22             lr = self.scheduler(update_step)
23             self.optimizer.set_lr(lr)
24             self.optimizer.step()
25             self.optimizer.zero_grad()
26             global_step += 1
27             self._log_loss(process_bar, sum(lossess) / len(lossess), i)
28         if (i + 1) % self.eval_interval == 0:
29             eval_metrics = self.eval()
30             pprint(eval_metrics)

```

由于我们设置了gradient_accumulation_steps，因此在相应步的微batch之后更新一次参数

使用microbatch可以使得训练时的峰值显存降低，在gradient_accumulation_steps步进行一次更新，用时间换空间，通过多次计算来模拟一个大的批次，从而在有限的显存下实现大批量训练的效果

实验结果

由于只有3090，单卡无法进行全量微调，因此我在这使用了lora进行微调，效果与全量微调差异不大

Expert iteration

Algorithm 2 Expert iteration (EI)

Input initial policy model $\pi_{\theta_{\text{init}}}$; reward function R ; task questions \mathcal{D}

```
1: policy model  $\pi_{\theta} \leftarrow \pi_{\theta_{\text{init}}}$ 
2: for step = 1, ..., n_ei_steps do
3:   Sample a batch of questions  $\mathcal{D}_b$  from  $\mathcal{D}$ 
4:   Set the old policy model  $\pi_{\theta_{\text{old}}} \leftarrow \pi_{\theta}$ 
5:   Sample  $G$  outputs  $\{o^{(i)}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)$  for each question  $q \in \mathcal{D}_b$ 
6:   Compute rewards  $\{r^{(i)}\}_{i=1}^G$  for each sampled output  $o^{(i)}$  by running reward function  $R(q, o^{(i)})$ 
7:   Filter out wrong outputs (i.e.,  $o^{(i)}$  with  $r^{(i)} = 0$ ) to obtain a dataset  $\mathcal{D}_{\text{sft}}$  of correct question-response pairs
8:    $\pi_{\theta} \leftarrow \text{SFT}(\pi_{\theta}, \mathcal{D}_{\text{sft}})$  (Algorithm 1)
9: end for
```

Output π_{θ}

对于每一个question，采样G条输出，过滤掉错误的输出，然后用得到的问答对来做sft

GRPO

Algorithm 3 Group Relative Policy Optimization (GRPO)

Input initial policy model $\pi_{\theta_{\text{init}}}$; reward function R ; task questions \mathcal{D}

```
1: policy model  $\pi_{\theta} \leftarrow \pi_{\theta_{\text{init}}}$ 
2: for step = 1, ..., n_grpo_steps do
3:   Sample a batch of questions  $\mathcal{D}_b$  from  $\mathcal{D}$ 
4:   Set the old policy model  $\pi_{\theta_{\text{old}}} \leftarrow \pi_{\theta}$ 
5:   Sample  $G$  outputs  $\{o^{(i)}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)$  for each question  $q \in \mathcal{D}_b$ 
6:   Compute rewards  $\{r^{(i)}\}_{i=1}^G$  for each sampled output  $o^{(i)}$  by running reward function  $R(q, o^{(i)})$ 
7:   Compute  $A^{(i)}$  with group normalization (Eq. 28)
8:   for train step = 1, ..., n_train_steps_per_rollout_batch do
9:     Update the policy model  $\pi_{\theta}$  by maximizing the GRPO-Clip objective (to be discussed, Eq. 29)
10:   end for
11: end for
```

Output π_{θ}

GRPO算法是将PPO算法中的GAE优势替换为群体相对优势，具体来说：

1. 先对一个问题采样G条responses
2. 对这G条response计算奖励
3. 优势的计算公式 $A_i = (r_i - \text{mean}(r)) / \text{std}(r)$

compute_group_normalized_rewards

这个函数要求给定一组响应以及其对应的答案，计算群体相对优势

其中rollout_response的长度是 `n_prompt*group_size`

```
Python |
1  def compute_group_normalized_rewards(
2      reward_fn: Callable[[str, str], dict[str, float]],
3      rollout_responses: list[str],
4      repeated_ground_truths: list[str],
5      group_size: int,
6      advantage_eps: float,
7      normalize_by_std: bool,
8  ) -> Tuple[torch.Tensor, torch.Tensor, dict[str, float]]:
9      rollout_batch_size = len(repeated_ground_truths)
10     assert rollout_batch_size % group_size == 0
11     advantages = torch.zeros(rollout_batch_size, dtype=torch.float32)
12     rewards = torch.zeros(rollout_batch_size, dtype=torch.float32)
13     for i in range(0, rollout_batch_size, group_size):
14         batch_resp = rollout_responses[i : i + group_size]
15         batch_gt = repeated_ground_truths[i : i + group_size]
16         batch_advantages = []
17         for resp, gt in zip(batch_resp, batch_gt):
18             reward_dict = reward_fn(resp, gt)
19             reward = reward_dict.get("reward", 0)
20             batch_advantages.append(reward)
21         rewards[i : i + group_size] = torch.tensor(batch_advantages, dtype=
=torch.float32)
22         denom = advantage_eps + rewards[i : i + group_size].std() if norma
lize_by_std else 1
23         advantages[i : i + group_size] = (rewards[i : i + group_size] - re
wards[i : i + group_size].mean()) / denom
24
25     metadata = {
26         "reward_mean": rewards.mean().item(),
27         "reward_std": rewards.std().item(),
28         "advantage_mean": advantages.mean().item(),
29         "advantage_std": advantages.std().item(),
30         "response_length": sum([len(resp) for resp in rollout_responses])
// len(rollout_responses),
31     }
32     return advantages, rewards.detach(), metadata
```

Naive policy gradient loss

这个函数不计算重要性采样部分，只是简单地将优势与对数概率相乘即可

$$-A_t \cdot \log p_{\theta}(o_t|q, o_{<t}).$$

```
1 def compute_naive_policy_gradient_loss(  
2     raw_rewards_or_advantages: torch.Tensor,  
3     policy_log_probs: torch.Tensor,  
4 ) -> torch.Tensor:  
5     return -raw_rewards_or_advantages * policy_log_probs
```

GRPO clip loss

这一部分是真正的GRPO的loss，我们需要引入重要性采样，同时引入clip，使得新策略与旧策略不会相差过大

```
1 def compute_grpo_clip_loss(  
2     advantages: torch.Tensor, # (batch_size, 1)  
3     policy_log_probs: torch.Tensor, # (batch_size, sequence_length)  
4     old_log_probs: torch.Tensor, # (batch_size, sequence_length)  
5     cliprange: float,  
6 ) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:  
7     ratio = torch.exp(policy_log_probs - old_log_probs)  
8     clipped_ratio = torch.clamp(ratio, 1 - cliprange, 1 + cliprange)  
9     loss = -torch.min(ratio, clipped_ratio) * advantages  
10    return loss, {"ratio": ratio.detach(), "clipped_ratio": clipped_ratio.  
    detach()}
```

compute_policy_gradient_loss

这个函数用于在训练时作为接口，根据loss_type来区分使用哪个loss，方便后面做消融实验

```

1 def compute_policy_gradient_loss(
2     policy_log_probs: torch.Tensor,
3     loss_type: Literal["no_baseline", "reinforce_with_baseline", "grpo_clip"],
4     raw_rewards: torch.Tensor | None = None,
5     advantages: torch.Tensor | None = None,
6     old_log_probs: torch.Tensor | None = None,
7     cliprange: float | None = None,
8 ) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:
9     if loss_type == "no_baseline":
10         assert raw_rewards is not None
11         loss = compute_naive_policy_gradient_loss(raw_rewards, policy_log_probs)
12         return loss, {}
13     elif loss_type == "reinforce_with_baseline":
14         assert advantages is not None
15         loss = compute_naive_policy_gradient_loss(advantages, policy_log_probs)
16         return loss, {}
17     elif loss_type == "grpo_clip":
18         assert advantages is not None
19         assert old_log_probs is not None
20         assert cliprange is not None
21         loss, metadata = compute_grpo_clip_loss(advantages, policy_log_probs, old_log_probs, cliprange,)
22         return loss, metadata

```

GRPO experiment

现在我们可以完整地实现GRPO训练的流程了，GRPO训练的流程大致如下：

1. 先采样一批问题 D
2. 从这批问题中，对每一个问题都采样G条答案

```

1 prompts = batch['prompts']
2 ground_truths = batch['ground_truths']
3
4 responses = self._generate_responses(prompts, G)

```

3. 计算奖励与优势

```

1  advantages, raw_rewards, reward_metadata = compute_group_normalized_rewards
  (
2      self.reward_fn,
3      responses,
4      repeated_ground_truths,
5      self.args.group_size,
6      self.args.advantage_eps,
7      self.args.use_std_normalization
8  )

```

responses是长度为(n_prompt*group_size)的list

最后得到的advantages形状是(n_prompt*group_size,)

4. 使用当前的模型，计算所有response的对数概率

```

1  policy_log_probs, response_mask = self._get_log_probs(repeated_prompts, res
  ponses)

```

将prompt和response拼接，再次输入给大模型进行一次前向传播，得到每个token的logits

将当前的log_probs设置为old policy的log_probs

5. 将batch分成micro_batch，计算损失并更新参数

```

1  micro_policy_log_probs = policy_log_probs[i:end_idx]
2  micro_response_mask = response_mask[i:end_idx]
3  micro_advantages = advantages[i:end_idx].unsqueeze(-1) # 添加维度以匹配
4  micro_raw_rewards = raw_rewards[i:end_idx].unsqueeze(-1)
5  micro_old_log_probs = old_log_probs[i:end_idx] if old_log_probs is not None
   else None

```

每一个micro_batch，计算一次损失：

```

1  loss, metadata = grpo_microbatch_train_step(
2      micro_policy_log_probs,
3      micro_response_mask,
4      self.args.gradient_accumulation_steps,
5      self.args.loss_type,
6      raw_rewards=micro_raw_rewards if self.args.loss_type == "n
   o_baseline" else None,
7      advantages=micro_advantages if self.args.loss_type in ["re
   inforce_with_baseline", "grpo_clip"] else None,
8      old_log_probs=micro_old_log_probs,
9      cliprange=self.args.cliprange if self.args.loss_type == "g
   rpo_clip" else None
10 )

```

在训练过程中，用到了vllm来加快模型的推理，因此训练的时候需要用到两张GPU，一张用来装载训练模型，一张用于装载vllm模型。

完整的训练代码如下:

```

1  global_step = 0
2  process_bar = tqdm(self.train_dataloader, total=min(len(self.train_dataloader), self.args.n_grpo_steps), desc="Training")
3  for grpo_step, batch in enumerate(process_bar):
4      if grpo_step >= self.args.n_grpo_steps:
5          break
6      self.policy.eval()
7      if self.vllm_model is not None:
8          self._load_policy_to_vllm()
9      prompts = batch["prompts"]
10     ground_truths = batch["ground_truths"]
11
12     responses = self._generate_responses(prompts, n_samples=self.args.group_size)
13
14     rollout_prompts = [p for p in prompts for _ in range(self.args.group_size)]
15     repeated_ground_truths = [gt for gt in ground_truths for _ in range(self.args.group_size)]
16
17     advantages, raw_rewards, reward_metadata = compute_group_normalized_rewards(
18         self.reward_fn,
19         responses,
20         repeated_ground_truths,
21         self.args.group_size,
22         self.args.advantage_eps,
23         self.args.use_std_normalization,
24     )
25     advantages = advantages.unsqueeze(-1)
26     raw_rewards = raw_rewards.unsqueeze(-1)
27     self._log(reward_metadata)
28     if reward_metadata["reward_mean"] <= 0:
29         print("Reward mean is zero, stop train.")
30         break
31     with torch.no_grad():
32         if self.args.loss_type == "grpo_clip":
33             old_log_probs = []
34             for i in range(0, len(responses), self.micro_train_batch_size):
35                 micro_batch_prompts = rollout_prompts[i : i + self.micro_train_batch_size]
36                 micro_batch_responses = responses[i : i + self.micro_train_batch_size]
37                 masked_log_probs, _ = self._get_log_probs(micro_batch_prompts, micro_batch_responses, self.args.importance_sample_level)
38                 old_log_probs.append(masked_log_probs.detach())

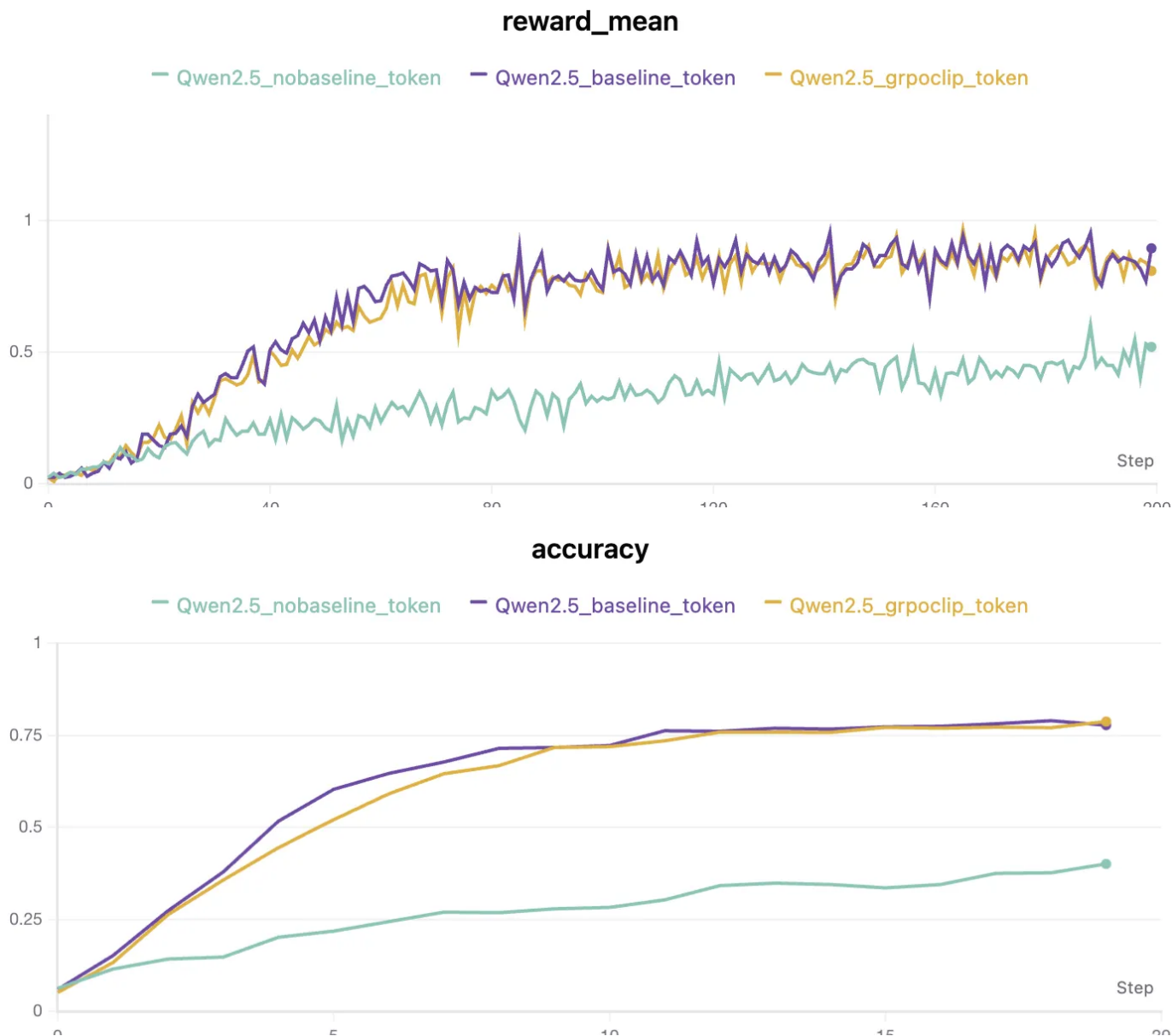
```

```

39     else:
40         old_log_probs = None
41         # === Training Phase ===
42         self.policy.train()
43         for _ in range(self.args.epochs_per_rollout_batch):
44             losses = []
45             for i in range(0, len(responses), self.micro_train_batch_size):
46                 global_step += 1
47                 micro_batch_slice = slice(i, i + self.micro_train_batch_size)
48                 batch_prompts = rollout_prompts[micro_batch_slice]
49                 batch_responses = responses[micro_batch_slice]
50                 batch_advantages = advantages[micro_batch_slice].to(self.policy
51 y.device)
52                 batch_raw_rewards = raw_rewards[micro_batch_slice].to(self.pol
53 icy.device)
54                 batch_old_log_probs = old_log_probs[i // self.micro_train_batch
55 h_size].to(self.policy.device) if old_log_probs else None
56                 policy_log_probs, response_mask = self._get_log_probs(batch_pr
57 ompts, batch_responses, self.args.importance_sample_level)
58                 loss, metadata = grpo_microbatch_train_step(
59                     policy_log_probs=policy_log_probs,
60                     response_mask=response_mask,
61                     gradient_accumulation_steps=self.args.gradient_accumulation
62 n_steps,
63                     loss_type=self.args.loss_type,
64                     raw_rewards=batch_raw_rewards,
65                     advantages=batch_advantages,
66                     old_log_probs=batch_old_log_probs,
67                     cliprange=self.args.cliprange,
68 )
69                 if ((i // self.micro_train_batch_size) + 1) % self.args.gradient
70 accumulation_steps == 0:
71                     torch.nn.utils.clip_grad_norm_(self.trainable_parameters,
72 self.args.max_grad_norm)
73                     self.optimizer.step()
74                     self.optimizer.zero_grad()
75             if (grpo_step + 1) % self.args.eval_interval == 0:
76                 self.policy.eval()
77                 eval_metrics = self.evaluate(output_dir=self.output_dir / f"checkp
78 oint-step-{grpo_step + 1}")
79                 self._log(eval_metrics)

```

实验结果



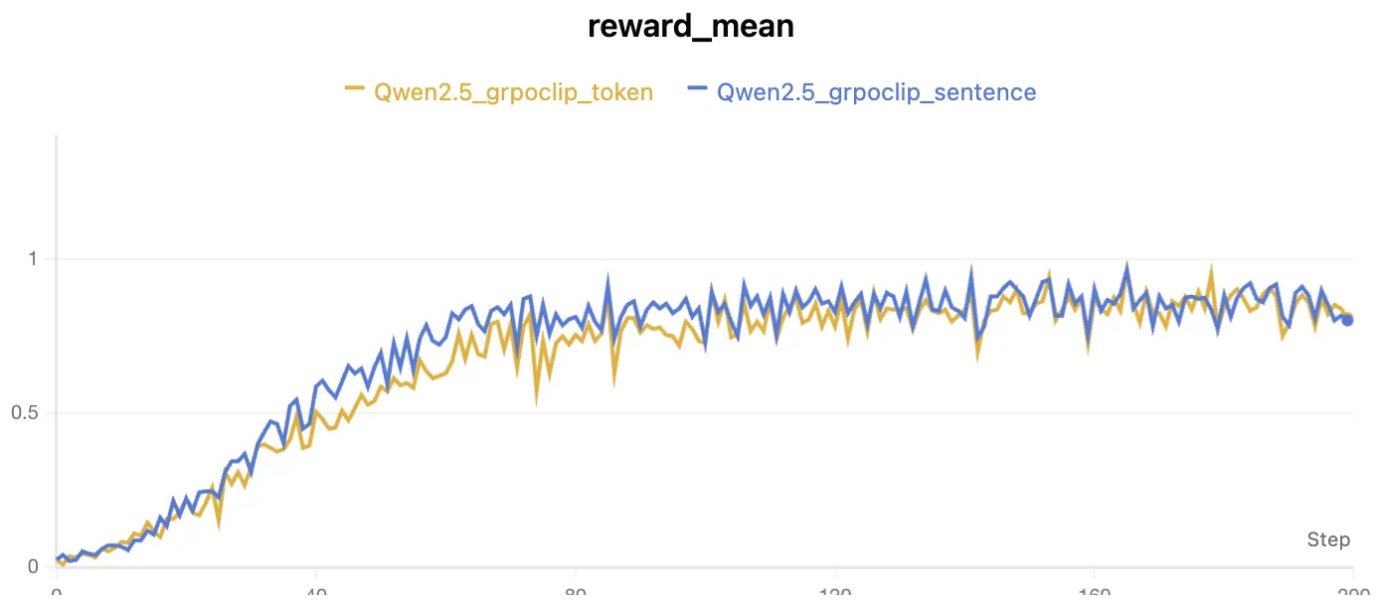
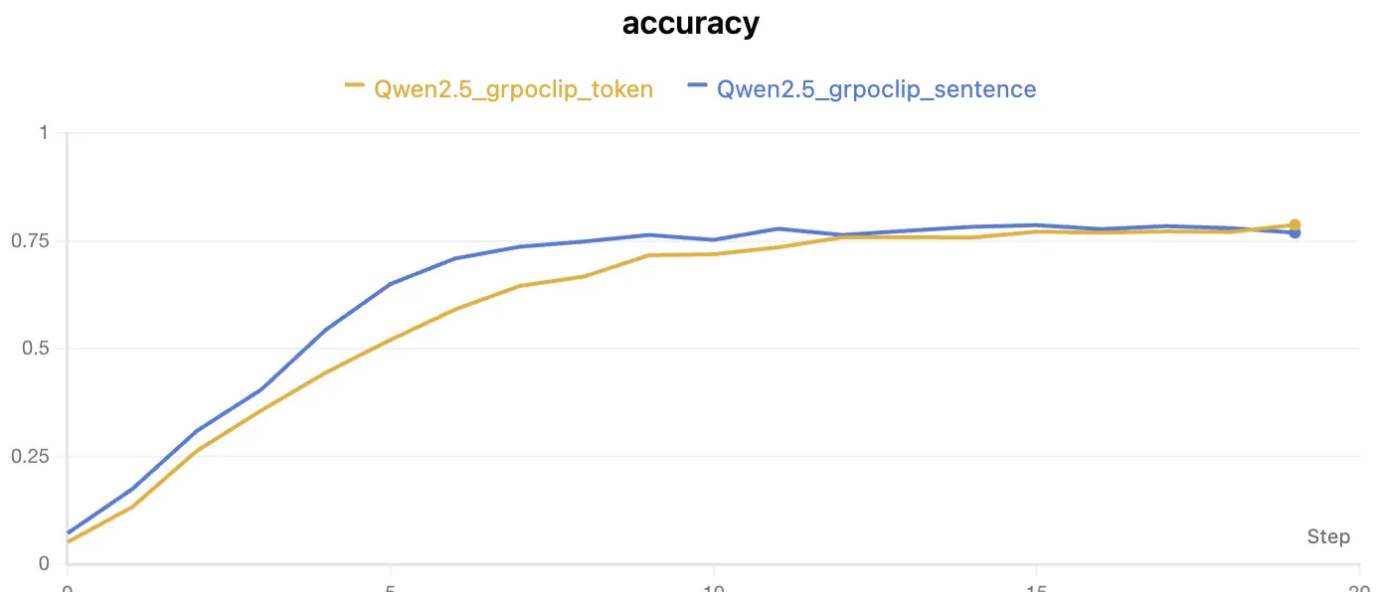
重要性采样—token or sentence

GRPO的重要性采样是token level，即每个token的对数概率都要乘以优势

GSPO则是用每条完整回答的概率（sequence likelihood）来做重要性采样

在实现的时候，在计算log-probabilities的时候就进行mask_mean，这样就得到了完整回答的概率

对比token level 与 sentence level，实验结果如下：



可以看到sentence级别的重要性采样相比较于token级的重要性采样，收敛会更快一些