

# Deep Learning Homework 3

Björn Pettersson, bjpe900 2025-81022

01-12-2025

---

## Problem 1 : Continual Learning

How can we learn a sequence of tasks or data streams incrementally over time, without forgetting what they have learned previously?

### 1.1.a)

**Class incremental learning** The model receives new classes over time. Each training phase introduces additional labels that were not present before. The model must learn to classify both old and new classes using a single shared classifier, without forgetting the earlier ones.

### Domain incremental learning

The set of classes stays the same, but the input distribution changes over time (e.g., new environments, styles, sensors). The model must adapt to the new domain while keeping performance on the old domains, even though the labels themselves do not change.

### 1.1.b)

The 3 big categories of continual learning is:

1. **-Replay Methods:** They store and replay samples from previous tasks (or generate pseudo-samples) alongside new data to mitigate forgetting.
2. **-Regularization based methods:** They add regularization terms to the loss function to preserve important parameters learned from previous tasks. These constraints prevent large updates that could destroy old knowledge.
3. **-Parameter isolation methods:** They isolate model parameters for each task, either by using a fixed network with masks or dynamically growing the architecture, thereby avoiding interference between tasks.

## 1.2. Joint Training vs Offline CL vs Online CL (10 pts)

### 1.2.a)

#### Joint Training

We combine the data for learning all tasks into one dataset and train on this, Catastrophic forgetting is inherently avoided because the model repeatedly sees data from all tasks. This serves as an upper-bound on performance for the other two paradigms.

#### Offline continual learning (tasks incremental)

We sequentially train the model on tasks, ie first task A on dataset A, then task B on dataset B, etc. The model only sees the dataset for the current task during training but we train for multiple epochs ie the model sees a given example multiple times. When learning new things the model overwrites the weight values it had from learning previous tasks, leading to catastrophic forgetting. Problem is stability vs plasticity, how much we want to retain ability to perform old tasks, while keeping ability to learn new.

#### Online continual learning ()

The examples from the tasks are presented like a stream, each example is only shown once. The model trains on a single pass, like a,a,a,a,b,b,b,b,c,c,c,d,d,d,d,d for task from dataset A,B,C,D. This limits the time and data available to learn a task since we only see it once, but its closer to how humans learn.

### 1.2.b)

#### Most difficult is online continual learning

We still have the challenge of catastrophic forgetting of the previous tasks but we also have limited data for our current task risking underfitting. So in the stability plasticity tradeoff, we have problems with both, where as for offline CL we could at least get optimal parameter updates for each task since we had multiple passes to learn them.

## Problem 2 : Chain-of-Thought (CoT) Reasoning (20 points)

### 2.1

One shot prompt:

“

Q: A banker has 20 coins. He puts 5 coins into each of his 3 accounts. How many coins does he have left

A: We think about the problem step by step, i) the banker puts  $5*3=15$  coins into accounts. ii) he has  $20-15=5$  coins left. iii) Answer, the banker has 5 coins left.

Q: A baker has 20 cookies. He packs 4 cookies into each of his 3 boxes. How many cookies does he have left?

A: Lets think about this step by step

“

My one shot prompt, uses the step by step prompt reasoning helping the model see an example of how we want the problem broken down and reasoned about, ie how we want to arrive at a given answer. With zero shot, the model may attempt to just output the most likely number, while the reasoning steps here makes the answer more likely to be aligned with the type of reasoning we want.

## 2.2

Zero shot with CoT:

“

Q: A baker has 20 cookies. He packs 4 cookies into each of his 3 boxes. How many cookies does he have left ? A. Lets think about this step by step:

“

I added “think step by step” to encourage CoT behaviour, since its zero shot we dont show examples. But this line has been proven to make the model reason and effective for improving performance. There is no evidence of how or why exactly this works.

### 2.3.a)

One shot CoT triggers reasoning by providing the model with an example of how to solve a task, which improves generally performance. Another advantage is that we can choose the type of reasoning here after our own example.

### 2.3.b)

I expect one shot to outperform zero shot, because there are less things that can go wrong through hallucination since i am providing the model with the reasoning steps it needs to solve this type of tasks, while for zero shot it has to make them up even though we can trigger CoT behaviour by the “lets think step by step” prompt.

### Problem 3 : Diffusion Models (DDPM vs. DDIM) (20 Points)

#### 3.1.a)

Markov definition of the forward process

Independence of gaussian noise

Show distribution of the  $x_t$  conditioned on  $x_0$  can be written in closed form as:

$$q(x_t | x_0) = N(x_t; \sqrt{\bar{a}_t} x_0, (1 - \bar{a}_t)I),$$

$$\bar{a}_t = \prod_{s=1}^t (1 - \beta_s)$$

We want to derive:  $x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$ , with  $\epsilon \sim N(0, I)$ .

From:  $x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon_t$ , and  $a_t = 1 - \beta_t$ :

We derive it as:

With definition of  $a_t$ :

$$x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon_t$$

$$x_t = \sqrt{a_t} x_{t-1} + \sqrt{1 - a_t} \epsilon_t$$

we apply the formula recursively for  $x_{t-1}$ :

$$x_{t-1} = \sqrt{a_{t-1}} x_{t-2} + \sqrt{1 - a_{t-1}} \epsilon_{t-1}$$

gives:

$$x_t = \sqrt{a_t} x_{t-1} + \sqrt{1 - a_t} \epsilon_t$$

$$x_t = \sqrt{a_t} (\sqrt{a_{t-1}} x_{t-2} + \sqrt{1 - a_{t-1}} \epsilon_{t-1}) + \sqrt{1 - a_t} \epsilon_t$$

$$x_t = \sqrt{a_t a_{t-1}} x_{t-2} + \sqrt{a_t (1 - a_{t-1})} \epsilon_{t-1} + \sqrt{1 - a_t} \epsilon_t$$

We know that for independent Normal functions  $\sqrt{a} \epsilon_1 + \sqrt{b} \epsilon_2 = \sqrt{a + b} \bar{\epsilon}$ :

Hence the two last terms becomes:

$$\begin{aligned} \sqrt{a_t (1 - a_{t-1})} \epsilon_{t-1} + \sqrt{1 - a_t} \epsilon_t &= \\ \sqrt{a_t (1 - a_{t-1}) + 1 - a_t} \bar{\epsilon}_{t-1} &= \\ \sqrt{a_t - a_t a_{t-1} + 1 - a_t} \bar{\epsilon}_{t-1} &= \\ \sqrt{1 - a_t a_{t-1}} \bar{\epsilon}_{t-1} &= \end{aligned}$$

hence:

$$x_t = \sqrt{a_t a_{t-1}} x_{t-2} + \sqrt{a_t(1 - a_{t-1})} \epsilon_{t-1} + \sqrt{1 - a_t} \epsilon_t$$

becomes

$$x_t = \sqrt{a_t a_{t-1}} x_{t-2} + \sqrt{1 - a_t a_{t-1}} \bar{\epsilon}_{t-1}$$

If we continue to recursively apply until we get back to  $x_0$ , we apply this process  $t$  times:

$$x_t = \sqrt{a_t a_{t-1} \dots a_1} x_0 + \sqrt{1 - a_t a_{t-1} \dots a_1} \epsilon$$

With our definition of  $\bar{a}_t = \prod_{s=1}^t (1 - \beta_s) = \prod_{s=1}^t a_s = a_1 a_2 \dots a_t$ :

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$$

So per our process of continuously noising our sample from a clean  $x_0$  to a completely random  $x_t$  we should now have arrived at a gaussian random noise distribution after  $t$  applications of the forward process, hence:

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$$

$$\epsilon \sim N(0, 1)$$

### 3.1.b.1)

For  $T=1000$ , generating a single image requires 1000 forward (denoising) passes through a large NN. This is computationally expensive because 1. 1000 is a large number of forward passes, 2. the NN is large, 3. We have to pass it through sequentially since it is a Markov process each step depends on the previous.

### 3.1.b.2)

What DDPM does: It generates images by gradually removing noise through many small sequential denoising steps. Each step is learning the reverse process of adding a slight bit of Gaussian noise to a given input image. Repeating this enough times takes us from complete noise to clear image.

Why we can't jump through steps: 1 The denoising process is a Markov chain which means that each step depends on the previous step. 2 the model is only trained to partially denoise a given input. 3 we choose more steps in training this partial denoising over fewer for stability as very few steps in a much more complex task and may lead to worse results.

### 3.2. DDIM : Deterministic Reverse Process and Comparison (8 pts)

#### 3.2.a)

Solve:  $\hat{x}_0(x_t, t) = \frac{x_t - \sqrt{1 - \bar{a}_t} \epsilon_\theta(x_t, t)}{\sqrt{\bar{a}_t}}$  in terms of  $x_t$  and  $\epsilon_\theta(x_t, t)$ , and write down estimator:

We have the forward process from 1a:

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$$

We want to solve for  $x_0$  in terms of  $x_t$  and predict the noise  $\epsilon_\theta(x_t, t)$ :

In DDIM we get a networks noise prediction  $\epsilon \approx \epsilon_\theta(x_t, t)$ , hence we have:

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon_\theta(x_t, t)$$

we isolate  $x_0$  term:

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon_\theta(x_t, t)$$

$$x_t - \sqrt{1 - \bar{a}_t} \epsilon_\theta(x_t, t) = \sqrt{\bar{a}_t} x_0$$

$$x_t - \sqrt{1 - \bar{a}_t} \epsilon_\theta(x_t, t) = \sqrt{\bar{a}_t} x_0$$

$$\frac{x_t - \sqrt{1 - \bar{a}_t} \epsilon_\theta(x_t, t)}{\sqrt{\bar{a}_t}} = x_0$$

#### 3.2.b)

Show that we can choose a determinist update of the form  $x_{t-1} = \sqrt{\bar{a}_{t-1}} \hat{x}_0(x_t, t) + \sqrt{1 - \bar{a}_{t-1}} \epsilon_\theta(x_t, t)$  :

from 1a the forward process to create noisy image at time t from clean image  $x_0$  and noise  $\epsilon$ :

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$$

creating  $x_{t-1}$  from  $x_0$ :

$$x_{t-1} = \sqrt{\bar{a}_{t-1}} x_0 + \sqrt{1 - \bar{a}_{t-1}} \epsilon$$

From 3.2a we have estimate of clean image  $\hat{x}_0(x_t, t)$ , which and predicted noise  $\epsilon \approx \epsilon_\theta(x_t, t)$ :

$$x_{t-1} = \sqrt{\bar{a}_{t-1}} x_0 + \sqrt{1 - \bar{a}_{t-1}} \epsilon$$

becomes:

$$x_{t-1} = \sqrt{\bar{a}_{t-1}} \hat{x}_0(x_t, t) + \sqrt{1 - \bar{a}_{t-1}} \epsilon_\theta(x_t, t)$$

From  $x_t$  (noisy) we estimate the clean image  $\hat{x}_0$ . We then recreate  $x_{t-1}$  as a less noisy version of the same clean image. If we assume the image is the same throughout the process, the noise stays the same and produces the same output. We don't have any random sampling and use the same noise direction  $\epsilon_\theta(x_t, t)$  for both timesteps just scaled differently. This is why the process is deterministic.

### 3.2.c)

#### How DDIM can use only a small subset of timesteps

DDIM is deterministic, i.e. there is not randomness introduced in between steps unlike in DDPM where the noise introduced is new random sampling each step. DDIM instead uses the same noise scaled differently, and moves towards the same goal end image.

Each timestep in DDIM, we estimate the image that the noise is added to as well  $x_0$ . Early in the process this is a bad estimation but this doesn't matter so much because we add a lot of noise to it anyway. We can't really see what the noise is added to anyways. But we can use this "behind the scenes" peek at the final image from any process step, to jump in the process. In essence we always get an estimated final picture each step we just decide how many steps we want to run the denoising process to make it good.

#### Why this typically makes sampling much faster than DDPM

Sampling in the context of DDIM is the process of generating a new image (sample) from the original distribution of clean images from the random noise vector.

It is faster in DDIM than DDPM because we can skip steps because of the deterministic nature of DDIM, see above explanation.

## Problem 4 : Basic Methods of Reinforcement Learning (20 points)

### 4.1.a)

A process satisfies the Markov property if the future state of the process depends only on the present state and not on the sequence of events preceding it.

Why does it hold in the student's MDP?

Arrows are possible actions, states are circles with rewards  $R$ . So at a given state we have the option to move along the connected arrows (actions). If the actions we can take and hence where we can end up at state  $x_{t+1}$  is a function of the current state  $x_t$ . The actions in  $x_t$  does not depend on how one arrived there i.e.  $x_{t-1}$ . Once we are at  $x_t$  the history before that does not matter. Hence the Markov property is fulfilled.

#### 4.1.b)

##### i. Cumulative reward of 1,2,3,4,5

$$R_1 + R_2 + R_3 + R_4 + R_5 = -0.5 + -1.2 + 0 + 2 + 50 = 50.3$$

##### ii. Probability of sequence

$$P(\text{sequence}) = P_{\text{relax}} * P_{\text{youtube}} * P_{\text{study}} * P_{\text{submit}} = 0.6 * 0.5 * 0.6 * 0.4 = 0.072$$

Probability of the whole seequence occurring is  $P(\text{sequence}) = 7.20000 * 10^{-2}$

#### 4.1.c)

##### Policy

The agents behaviour, ie a mapping from all states to what actions to take for those states. It can be deterministic or stchastic probability of taking soem action given a state. For the student a policy could be to always study for state 1,2,3 and to submit if studying state 4. :

Policy: Learning agent's way of behaving at a given time • A mapping from perceived states of the environment to actions to be taken when in those states

##### Reward

The reward is a scalar feedback signal recieved for moving form one state to another. This can be used to train a RL agent. In the student MDP its is the R signal. :

Reward (signal): Agent's benefit at each time step • On each time step, a single number that the environment sends to the reinforcement learning agent

##### Value Function / Q funciton

ESEtiamtes ling term expectedred cumuative reward ffrom a gvein state or state action pair. Value funciton is a funciton only of state, and Q a funciton of state action pairs. For the student MDP it will let the student move to study even though immediate reward is negative, in order to get the higher cumujlative reward including the submitting at +50.:

Value function: Agent's benefit at the end of the episode • The value of a state is the total amount of reward an agent can expect to accumulate over the future, starting from that stat

##### Model

(optional) Model (of an environment): Mimicking the behavior of the environment

## 4.2 Policy Gradient Methods PGM

### 4.2.a)

Given policy gradient theorem:

$$\nabla_{\theta} J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi_{\theta}}(s, a) \nabla_{\theta} \pi_{\theta}(a | s)$$

### Role of $\mu(s)$ on policy state distribution in the policy gradient theorem

A policy is deciding on the action to take from some set of states visited. If we want to evaluate how “good” a given policy is we need to look at the expected return from visiting the states for the policy and taking the policy's actions from those states. Hence states that are not visited given our policy has no impact on how good the policy is and by the same logic states visited often has a larger impact. This is why we must weight the states by  $\mu(s)$  when calculating the gradient of the objective.

### Why can we not simply weight all states equally

Irrelevant states problem: Some states may be completely unreachable under the current policy  $\pi_{\theta}$ . For example, in the student MDP, if the policy never chooses "Youtube" from state 2, then certain states might never be visited. Optimizing the policy for unreachable states wastes computational resources and provides no benefit to actual performance. Misaligned optimization objective: The goal  $J(\theta) = v_{\pi_{\theta}}(s_0)$  measures expected return from the start state  $s_0$  under policy  $\pi_{\theta}$ . States that are never or rarely visited contribute negligibly to this objective, so equal weighting would optimize for the wrong goal. Incorrect gradient direction: Equal weighting could lead to gradient updates that improve performance in rarely-visited states at the expense of frequently-visited states, potentially decreasing overall performance  $J(\theta)$ . Violates the causal structure: The value of improving the policy at state  $s$  depends on how often  $s$  is reached. A state visited 100 times per episode should naturally have  $100\times$  more influence than a state visited once.

### 4.2.b)

PGM are stochastic gradient ascent realization of generalised policy iteration GPI

### Why REINFORCE can be viewed as a policy improvement step:

1. REINFORCE updates parameters. Moving  $\theta$  in the direction that increases  $J(\theta)$  the expected return from the start state.
2. The update pushes probability mass towards actions that yield higher return and away from actions that yielded low return. This “improves” the policy.
3. Instead of computing the full gradient REINFORCE uses sampled trajectories to get an unbiased stochastic estimate of this gradient. Each sample provides a noisy but correct

in expectation directly for improvement. 4. Unlike traditional policy iteration which makes the policy greedy wrt  $q_\pi$  REINFORCE performs soft improvement by gradually shifting action probabilities the stochastic policy naturally explore while still moving towards better actions.

### **Why using sampled return $G_t$ corresponds to improving an approximation of the action value function**

REINFORCE is a policy improvement algorithm that determines the gradient of a policy's performance with respect to its parameters  $\theta$  and uses gradient ascent to update them. It acts as an update step because it maximizes the expected return. The update rule is based on the likelihood ratio method and involves the use of the sampled accumulated return  $G_t$  as an unbiased estimate of the true action-value function. Since  $G_t$  is directly observed from a trajectory, it is used instead of approximating  $Q$ , which would require a separate learning process. The core identity leveraged is that the expected return is equal to  $Q$ . Therefore, using  $G_t$  in the gradient calculation provides a known, true value estimate, which improves the stability and practicality of the policy gradient estimate.

## **Problem 5 : PPO vs GRPO for Training Language Models (20 points)**

### **5.a.i**

Excessive large policy updates can harm training due to several fundamental RL related issues. In PPO advantage estimates are computed under the old policy  $\pi_{\theta_{old}}$  hence when policy changes dramatically these advantage estimates become invalid guides for policy optimization. Making gradient estimates unreliable. Furthermore in RL the training data is itself dependent on the current policy, i.e. take action get feedback depends on what action. Hence data distribution over observations and rewards are constantly changing, and large policy updates makes this worse by causing dramatic shifts in the data distribution. Lastly RL is very sensitive to hyperparameter tuning, if updates are too large the policy network can be pushed into a region of parameter space where it collects the next batch of data under very bad policy causing performance collapse from which it may never recover. In essence large updates break the assumption that collected experience remain relevant for optimization leading to poor data efficiency and training instability. This is because PPO is on policy assuming data used for training was collected by or is close to the current policy being optimized.

## 5.a.ii

### Describe PPOs clipping idea for preventing such instability

One idea to combat this type of instability is to set a bound for how far away from the current policy we can move when updating our policy. Idea introduced in Trust Region Policy Optimisation TRPO which is the basis on which PPO was built.

We move from objective in GPM:

$$\hat{g} = \hat{E}_t[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t]$$

In Trust region methods:

$$\max_{\theta} \hat{E}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t \right]$$

To make sure updated policy does not move too far away from current policy PPO adds a KL divergence constraint. The problem is that this constraint adds additional overhead and sometimes leads to bad training behaviour. PPO tries to mitigate this by incorporating it directly into the optimisation objective.

PPO Clipped surrogate Objective:

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$$

$$L^{CPI}(\theta) = \hat{E}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t \right] = \hat{E}_t [r_t(\theta) \hat{A}_t]$$

Main objective function in PPO

$$L^{CLIP}(\theta) = \hat{E}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)]$$

This is normal policy gradient objective  $\min(r_t(\theta) \hat{A}_t$

Clipped version of normal policy gradients objective:  $\text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)$

### How does this mechanism limit the size of the policy update?

How clip works:

Advantage can be positive or negative affecting the value of the min operator:

- Positive advantage: Selected action had better than expected effect on outcome. Loss function flattens out when r gets to high. This happens when action is a lot more likely under current policy than under old policy, in this case we don't want to overdo the policy update so the objective function gets clipped here to limit the effect of the gradient update.

- Negative advantage: Action had an estimated negative value, objective flattens when  $r$  goes near zero. This corresponds to actions that are much less likely now than in the old policy, here we don't want to overdo the update and set the probability of these actions to 0. Advantage is noisy so we don't want to destroy policy based on a single estimate. If on the other hand the action is very probable  $r$  is big and at the same time the advantage is negative the objective tells us to make this action less probable by an amount proportional to how wrong we were. This is the only region where the unclipped version of the function has a lower value than the clipped one and gets returned by the min operator.

In essence PPO does the same as TRPO objective which is to force the policy updates to be conservative if they move very far away from the current policy. The difference is that TRPO pushes complicated KL divergence to achieve this while PPO uses a simpler min function. PPO often outperforms the more complex function anyway.

### 5.b.i

The per-token KL penalty serves as a regularizer that prevents the fine-tuned policy  $\pi_\theta$  from straying too far from the reference model  $\pi_{\text{ref}}$ . This prevents overfitting, RLHF is prone to overfitting to the specific reward model which is trained on limited human feedback data. KL term penalises large deviations from the reference model. RLHF want to tune answers to human preference not fundamentally change the model.

### 5.b.ii

Reward Hacking Prevention: The reward model is an imperfect approximation of human preferences learnt from limited data, without the KL penalty, the policy can discover inputs or generation patterns that give high rewards but don't represent desirable behaviour. The gradients for these kinds of hacks could be steep and start to drastically change the model. KL anchors policy to a reference making sure it does not go too far away, into some bizarre policy region.

Distributional shift: Policy may collapse into a narrow mode repeating a smaller "safer" high reward policy rather than letting us access coverage of the whole language space or abandon natural language completely when token probabilities become extreme, the model may only favour some tokens and not others.

In essence the reward model is imperfect based on human feedback. It is possible that humans or the model itself could take advantage of this and find bizarre behaviour that still give high rewards. Hence without tying updates to a somewhat "sane" reference model with a KL penalty the model could go off the rails.

## 5.c

### How GRPO provides baseline without learnt critic:

We group outputs and compute each output's reward relative position in the group of sampled outputs. For each prompt GRPO samples  $G$  outputs from the old policy and computes the group relative advantage as  $\hat{A}_t^{GRPO} = \frac{R - \mu_R}{\sigma_R + \delta}$ . Reward  $R$ , average reward in group  $\mu_R$ , std of rewards in the group  $\sigma_R$ . This is as opposed to in PPO (agent critic) where advantage estimation is  $A_t^{PPO} = R_t - V(s_t)$  where  $V(s)$  is achieved by training a separate critic network. For  $\hat{A}_t^{GRPO}$  advantage is derived purely statistically. It measures the values of outcome  $R$  relative to the average output of the current batch.

### Why this baseline mechanism is especially suitable for LLMs:

Having a critic network as in PPO has some advantages such as lower variance and stable learning due to more precise value estimation. However GRPO has other advantages such as memory and computation efficiency and significant memory reduction by cutting out the critic network and deriving advantage statistically and more imprecisely. While PPO is best suited for general RL tasks, especially those requiring precise value estimation and low variance GRPO is best suited for large scale models, particularly LMs in RLHF where memory is a bottleneck and rewards are often sparse/binary.

LLMs as the name suggests large, PPO requires a critic network often similar size of policy network, which requires storage of critic gradients and optimiser states as well as forward and backward propagation. Being able to essentially cut the memory footprint in half by cutting out the critic is great for making RLHF feasible for very large models.

RLHF often has binary rewards, i.e. was this response helpful or not, correct/incorrect etc. Hence learning  $V(s)$  for all intermediate states overkill and unnecessary. Group based output is more natural simply asking what outputs are better. Each output is "discrete" one output / word prediction not continuous. Furthermore LLMs naturally generate  $G$  response per prompt during inference (context window) it is natural to use these samples to construct advantage directly. GRPO mirrors how humans evaluate LLM outputs, comparing multiple responses to the same query and selecting the better ones which is often how RLHF is carried out in practice.