

# Reinforced Self-Training (ReST) for Language Modeling

Caglar Gulcehre<sup>\*†,1</sup>, Tom Le Paine<sup>\*†,1</sup>, Srivatsan Srinivasan<sup>\*†,1</sup>, Ksenia Konyushkova<sup>†,1</sup>, Lotte Weerts<sup>†,1</sup>

Abhishek Sharma<sup>†,1</sup>, Aditya Siddhant<sup>†,1</sup>, Alex Ahern<sup>1</sup>, Miaosen Wang<sup>1</sup>, Chenjie Gu<sup>1</sup>,

Wolfgang Macherey<sup>2</sup>, Arnaud Doucet<sup>1</sup>, Orhan Firat<sup>†,1</sup>, Nando de Freitas<sup>1</sup>

\* Contributed equally, † Core contributors

<sup>1</sup>Google DeepMind, <sup>2</sup>Google Research

Reinforcement learning from human feedback (RLHF) can improve the quality of large language model’s (LLM) outputs by aligning them with human preferences. We propose a simple algorithm for aligning LLMs with human preferences inspired by growing batch reinforcement learning (RL), which we call *Reinforced Self-Training (ReST)*. Given an initial LLM policy, *ReST* produces a dataset by generating samples from the policy, which are then used to improve the LLM policy using offline RL algorithms. *ReST* is more efficient than typical online RLHF methods because the training dataset is produced offline, which allows data reuse. While *ReST* is a general approach applicable to all generative learning settings, we focus on its application to machine translation. Our results show that *ReST* can substantially improve translation quality, as measured by automated metrics and human evaluation on machine translation benchmarks in a compute and sample-efficient manner.

*Keywords: Offline RL, reinforcement learning, RL from human feedback, language, natural language processing, machine translation*

## 1. Introduction

Large language models (LLMs) have demonstrated impressive abilities in generating high-quality text and in solving numerous language tasks (Brown et al., 2020; Bubeck et al., 2023; Rae et al., 2021). These models are trained to maximize the likelihood of the next token autoregressively using massive amounts of text and compute (Hoffmann et al., 2022; Srivastava et al., 2022). However, Perez et al. (2022) showed that producing text with high likelihood does not necessarily align well with human preferences on various tasks. Without proper alignment, the language models can also output unsafe contents with harmful consequences. Moreover, aligning LLMs helps to improve on other downstream tasks (Ouyang et al., 2022b). *Reinforcement learning from human feedback* (RLHF) aims to address the alignment problem by using human preferences (Glaese et al., 2022; Stiennon et al., 2020; Wu et al., 2021). Typically, human feedback is used to learn a reward model, which is then used to fine-tune LLM with a reinforcement learning (RL) objective.

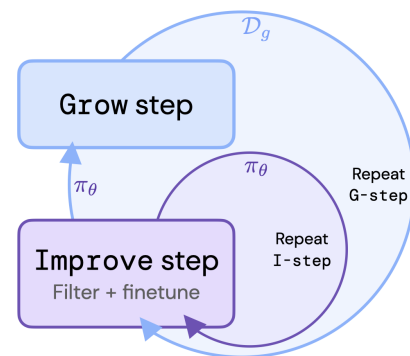


Figure 1 | **ReST method.** During Grow step, a policy generates a dataset. At Improve step, the filtered dataset is used to fine-tune the policy. Both steps are repeated, Improve step is repeated more frequently to amortise the dataset creation cost.

arXiv:2308.08998v2 [cs.CL] 21 Aug 2023

RLHF methods often rely on online RL methods such as PPO (Schulman et al., 2017) and A2C (Mnih et al., 2016). Online training requires sampling from the updated policy and scoring the samples with the reward model many times during training. The computational cost of dealing with a continual flow of new samples becomes a limitation of online methods, especially when the sizes of the policy and reward networks grow. Moreover, these methods are prone to reward “hacking” (Skalse et al., 2022), and prior works (Glaese et al., 2022) explored model regularization to mitigate this issue. Alternatively, offline RL methods learn from a fixed dataset of examples and, thus they are more computationally efficient and less prone to reward hacking. However, the quality of the policy learnt offline inevitably depends on the properties of the offline dataset (Fu et al., 2020; Gulcehre et al., 2021). As a result, carefully curated datasets become very important for the success of offline RL. Otherwise, the performance gains over supervised learning may be limited (Kumar et al., 2021). Concurrent to our work, (Rafailov et al., 2023) proposed a method called DPO (Direct Preference Optimization) that can make use of offline data to align an LM with human preferences.

We frame the alignment problem of a language model as a growing batch RL problem (Lange et al., 2012). Specifically, our Reinforced Self-Training (*ReST*) method includes two loops: in the inner loop (Improve), we improve the policy on a fixed dataset and in the outer loop (Grow), we grow the dataset by sampling from the latest policy (see Figure 1). In this work we consider conditional language modelling, then the steps of *ReST* are as follows:

1. **Grow (G):** The language model policy (initially, a supervised policy) is used to generate multiple output predictions for each context to augment the training dataset.
2. **Improve (I):** We rank and filter the augmented dataset with a scoring function. We use a learned reward model trained on human preferences as the scoring function in our experiments. Then, the language model is fine-tuned on the filtered dataset with an offline RL objective. This step can be repeated with an increasing filtering threshold. The final policy is then used in the next **Grow** step.

*ReST* is a general approach that allows different offline RL losses to be used in the inner loop when executing the **Improve** steps. To put it in practice, one only needs the ability to: i) sample from a model efficiently, ii) score the model’s samples. *ReST* provides several advantages over typical RLHF methods with online or offline RL:

- The computational burden is significantly reduced compared to online RL thanks to the output of **Grow** step being exploited across several **Improve** steps.
- The quality of the policy is not restricted by the quality of the original dataset (as in offline RL) since new training data is sampled from an improved policy during the **Grow** step.
- It is easy to inspect the data quality and potentially diagnose alignment issues, e.g., reward hacking, as the **Grow** and **Improve** steps are decoupled.
- The approach is simple, stable and has only a small number of hyperparameters to tune.

We explain the details of our proposed *ReST* approach in Section 3. Then, we present our experimental results on machine translation benchmarks in Section 4. Machine translation is a sequence-to-sequence learning problem (Sutskever et al., 2014), which is usually formulated as conditional language modelling where the context for conditioning is a sentence in a foreign language (source). We chose machine translation because i) it is an impactful application with strong baselines and a well-defined evaluation procedure, ii) several existing reliable scoring and evaluation methods are available for the use as a reward model (Freitag et al., 2022). In our experiments, we compare several offline RL algorithms on the IWSLT 2014 (Cettolo et al., 2014) and WMT 2020 benchmarks (Koehn et al., 2020) as well as more competitive, high-fidelity internal benchmarks on Web Domain. In our experiments *ReST* significantly improves reward model scores on test and validation sets. Furthermore, according to human raters, *ReST* generates higher quality translations compared to a supervised learning baseline.

## 2. Preliminaries

A conditional language model produces an output sequence  $\mathbf{y} = (y_1, y_2, \dots, y_T)$  given a context (or source input)  $\mathbf{x} = (x_1, x_2, \dots, x_L)$ , where the tokens  $x_l, y_t$  belong to a chosen vocabulary. A language generation policy  $\pi$  in an auto-regressive model characterized by a conditional probability distribution parameterized by  $\theta$  as

$$\pi_{\theta}(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^T \pi_{\theta}(y_t | \mathbf{y}_{1:t-1}, \mathbf{x}),$$

with the convention  $\mathbf{y}_{1:0} = \emptyset$  and  $\mathbf{y}_{1:t-1} = (y_1, y_2, \dots, y_{t-1})$ .

Let  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$  denote the data distribution. A given dataset  $\mathcal{D}$  consists of samples from this distribution:

$$\mathcal{D} = \{ (\mathbf{x}^i, \mathbf{y}^i)_{i=1}^N \text{ such that } \mathbf{x}^i \sim p(\mathbf{x}), \mathbf{y}^i \sim p(\mathbf{y}|\mathbf{x} = \mathbf{x}^i) \}.$$

Given this dataset, the supervised policy is trained by minimizing the negative log likelihood (NLL) loss:

$$\mathcal{L}_{\text{NLL}}(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} \left[ \sum_{t=1}^T \log \pi_{\theta}(y_t | \mathbf{y}_{1:t-1}, \mathbf{x}) \right]. \quad (1)$$

We refer to the model that is trained with the NLL loss as behavioral cloning (BC) (Pomerleau, 1989) following the RL literature nomenclature.

## 3. Reinforced Self-Training (*ReST*)

We present *ReST*, an RLHF algorithm that aligns the language model’s outputs with human preferences. Human preferences over sequences are modelled using a learned reward function (see Appendix A.4). In the underlying Markov decision process for conditional language modelling the states are the partial sequences, and the actions are the generated tokens (see Appendix A.1).

The *ReST* algorithm decouples the dataset growth and policy improvement of a typical RL pipeline into separate offline stages (Figure 1 and 2). We start by training an initial model  $\pi_{\theta}(\mathbf{y}|\mathbf{x})$  to map input sequences  $\mathbf{x}$  to output sequences  $\mathbf{y}$  on a given dataset of sequence pairs  $\mathcal{D}$  using the NLL loss from Equation (1). Next, the **Grow** step creates a new dataset  $\mathcal{D}_g$ , which augments the initial training dataset with samples from the model:

$$\mathcal{D}_g = \{ (\mathbf{x}^i, \mathbf{y}^i)_{i=1}^{N_g} \text{ such that } \mathbf{x}^i \sim \mathcal{D}, \mathbf{y}^i \sim \pi_{\theta}(\mathbf{y}|\mathbf{x}^i) \} \cup \mathcal{D}.$$

Here, the conditioning inputs are resampled from the original dataset  $\mathbf{x}^i \sim \mathcal{D}$ , as in self-training, but in situations where one has access to  $p(\mathbf{x})$  they could sample directly from it, i.e.,  $\mathbf{x}^i \sim p(\mathbf{x})$ . For example, consider a model that generates image from a textual description, in this case, the distribution of text inputs can be sampled from a language model  $p(\mathbf{x})$ .

Subsequently, the **Improve** steps use  $\mathcal{D}_g$  to fine-tune the policy  $\pi_{\theta}$ . Note that we keep the original dataset in the training mixture to ensure that the policies do not diverge. Below, we describe **Grow** and **Improve** steps in more details.

## Grow

The Grow step corresponds to the acting or data-generation step in RL. We create an augmented dataset of trajectories  $\mathcal{D}_g$  by sampling many output sequences from the current policy  $\pi_\theta$ , i.e.,  $y \sim \pi_\theta(y|x)$  for  $x \sim \mathcal{D}$ . The new dataset of sequences is then scored with a reward function  $R(x, y)$ . The datapoints with the reward above a threshold score are used to update the policy (see next). Once the policy is improved, a new dataset of better quality samples can be created once again (Figure 2, bottom).

## Improve

At the Improve step (exploitation or policy improvement in RL terminology), the goal is to use the new dataset  $\mathcal{D}_g$  to fine-tune the policy  $\pi_\theta$ . We start by defining a filtering function that includes only samples with rewards higher than a certain threshold  $\tau$ :

$$F(\mathbf{x}, \mathbf{y}; \tau) = \mathbb{1}_{R(\mathbf{x}, \mathbf{y}) > \tau}.$$

Let us note that the threshold based filtering function may result into learning suboptimal behaviors that favors outcomes with high variance in the environments with stochastic dynamics (Brandfonbrener et al., 2022). However, in this work we formulate the language modeling and translation tasks as deterministic RL problems (Appendix A.1.)

Next, we finetune the current best policy typically trained with either the supervised learning loss  $\mathcal{L}_{\text{NLL}}$  from equation 1 or an offline RL loss  $\mathcal{L}(\mathbf{x}, \mathbf{y}; \theta)$  on the filtered data such as V-MPO (Song et al., 2020) or offline actor-critic (Mathieu et al., 2021). To sum up, we use the following reward weighted loss  $J$ :

$$J(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_g} [F(\mathbf{x}, \mathbf{y}; \tau) \mathcal{L}(\mathbf{x}, \mathbf{y}; \theta)]. \quad (2)$$

Standard imitation learning approaches, such as BC (Pomerleau (1989), equation 1) and one-step RL methods like Behavior Value Estimation (BVE) (Gulcehre et al., 2021) perform one-step of Improve on the fixed dataset  $\mathcal{D}$ . In contrast, the basic version of ReST additionally includes a Grow step that allows the model to gather multiple new output sequences (potential translations) for contexts  $\mathbf{x}$  from the original dataset (source sentences to translate).

When iterating over Improve steps, we increase the filtering thresholds:  $\tau_1 < \dots < \tau_{N-1} < \tau_N$  (Figure 2). This filtering with the growing threshold results in data subsets of increasing quality but of decreasing size. As LLMs overfit to small datasets quickly, we fine-tune every new policy from the previous policy with a lower learning rate. Consecutive fine-tuning of policies  $\{\pi_{\theta_k}\}_{k \geq 1}$  on higher quality data subsets ensures policy improvement with a fixed dataset  $\mathcal{D}_g$ . If we were to sample from policies  $\{\pi_{\theta_k}\}_{k \geq 1}$ , the average reward of the generated samples would be increasing (shown in grey in Figure 2). As sampling from a policy in the Grow step is computationally expensive, after each such step we perform several Improve steps. Thus, the cost of a single dataset generation is amortised over multiple Improve steps. Algorithm 1 outlines the full ReST algorithm with multiple dataset growth and policy improvement steps.

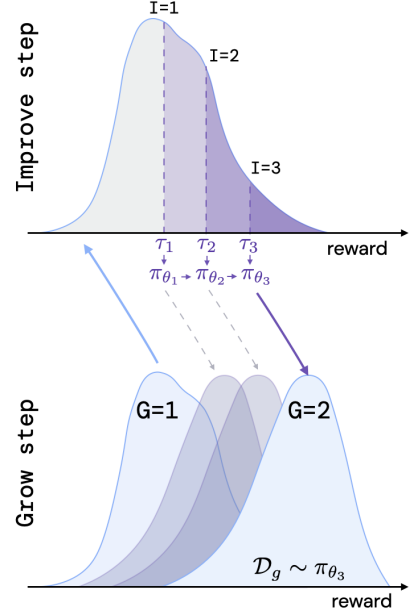


Figure 2 | **ReST algorithm.** *Top:* At Improve steps  $I=1, I=2, I=3$ , the dataset from the initial policy is filtered with thresholds  $\tau_1 < \tau_2 < \tau_3$  and a sequence of policies  $\pi_{\theta_1}, \pi_{\theta_2}, \pi_{\theta_3}$  are fine-tuned. *Bottom:* If we were to sample from those policies (grey), the quality of samples would increase. In practice, only the final policy  $\pi_{\theta_3}$  is used to generate the next dataset  $\mathcal{D}_g$ .

---

**Algorithm 1: ReST algorithm.** *ReST* is a growing-batch RL algorithm. Given an initial policy of reasonable quality (for example, pre-trained using BC) iteratively applies **Grow** and **Improve** steps to update the policy. Here  $F$  is a filtering function, and  $\mathcal{L}$  is an loss function.

---

**Input:**  $\mathcal{D}$ : Dataset,  $\mathcal{D}_{eval}$ : Evaluation dataset,  $\mathcal{L}(\mathbf{x}, \mathbf{y}; \theta)$ : loss,  $R(\mathbf{x}, \mathbf{y})$ : reward model,  $G$ : number of grow steps,  $I$ : number of improve steps,  $N$ : number of samples per context

Train  $\pi_\theta$  on  $\mathcal{D}$  using loss  $\mathcal{L}$ .

**for**  $g = 1$  to  $G$  **do**

    // **Grow**

    Generate dataset  $\mathcal{D}_g$  by sampling:  $\mathcal{D}_g = \{ (\mathbf{x}^i, \mathbf{y}^i) \}_{i=1}^{N_g}$  s.t.  $\mathbf{x}^i \sim \mathcal{D}$ ,  $\mathbf{y}^i \sim \pi_\theta(\mathbf{y}|\mathbf{x}^i)$  }  $\cup \mathcal{D}$ .

    Annotate  $\mathcal{D}_g$  with the reward model  $R(\mathbf{x}, \mathbf{y})$ .

**for**  $i = 1$  to  $I$  **do**

        // **Improve**

        Choose threshold s.t.  $\tau_1 > V_{\pi_\theta}$  for  $V_{\pi_\theta} = \mathbb{E}_{\mathcal{D}_g} [R(\mathbf{x}, \mathbf{y})]$  and  $\tau_{i+1} > \tau_i$ .

**while** *reward improves on  $\mathcal{D}_{eval}$*  **do**

            Optimise  $\theta$  on objective:  $J(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_g} [F(\mathbf{x}, \mathbf{y}; \tau_i) \mathcal{L}(\mathbf{x}, \mathbf{y}; \theta)]$

**end**

**end**

**end**

**Output:** Policy  $\pi_\theta$

---

**Probabilistic interpretation of the Improve step** Let us consider the particular choice  $\mathcal{L} = \mathcal{L}_{\text{NLL}}$ , with  $\theta'$  being the parameters of the model from the last **Grow** step,  $\lambda$  the proportion of data sampled from this model in  $\mathcal{D}_g$  and a single step of growth. The expression for the gradient in this case takes the following form:

$$\nabla J(\theta) = -\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \lambda \mathbb{E}_{\mathbf{y} \sim \pi_{\theta'}(\mathbf{y}|\mathbf{x})} [F(\mathbf{x}, \mathbf{y}; \tau) \nabla \log \pi_\theta(\mathbf{y} | \mathbf{x})] + (1 - \lambda) \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [F(\mathbf{x}, \mathbf{y}; \tau) \nabla \log \pi_\theta(\mathbf{y} | \mathbf{x})] \right]. \quad (3)$$

The first term on the RHS of (3) is similar to an online policy gradient term at the beginning of training when  $\theta \approx \theta'$  with  $F(\mathbf{x}, \mathbf{y}; \tau)$  replacing the state-action value function  $Q^\pi(\mathbf{x}, \mathbf{y})$ , when starting in state  $\mathbf{x}$  and taking sequential actions  $\mathbf{y}$ , that is generating synthetic data  $\mathbf{y}$  using policy  $\pi_\theta$  in our context. For the second term on the RHS of (3), we consider the original data  $\mathcal{D}$ , but we still ensure that it passes the threshold  $F(\mathbf{x}, \mathbf{y}; \tau)$ . Intuitively, people choose  $\mathcal{D}$  for training according to some possibly unknown criteria. In this work, we make the criterion  $F(\mathbf{x}, \mathbf{y}; \tau)$  explicit. The last term is therefore a form of offline policy gradients which prevents  $\pi_\theta(\mathbf{y} | \mathbf{x})$  to move too far from  $p(\mathbf{y} | \mathbf{x})$  which could lead to *model collapse* (Shumailov et al., 2023). Finally, note the similarity of this approach with self-training (Clark et al., 2003; Scudder, 1965; Xie et al., 2020) techniques. We provide a population interpretation (i.e., as  $N, N_g \rightarrow \infty$ ) of *ReST* in Appendix A.9.

In the following section, we explore how the choice of loss, filtering function and threshold, and synthetic data generated by language policy via sampling (exploration data) empirically affect the performance of the resulting policies  $\pi_\theta$ .

## 4. Experiments and analysis

We chose machine translation as a testbed for *ReST* as it is an impactful application of conditional language modeling where established reward models are available, for example, Metric X (Freitag et al., 2022), BLEURT (Sellam et al., 2020) and COMET (Rei et al., 2020). We ran experiments on two common benchmarks: IWSLT 2014 (Cettolo et al., 2014), and WMT 2020 (Koehn et al., 2020),

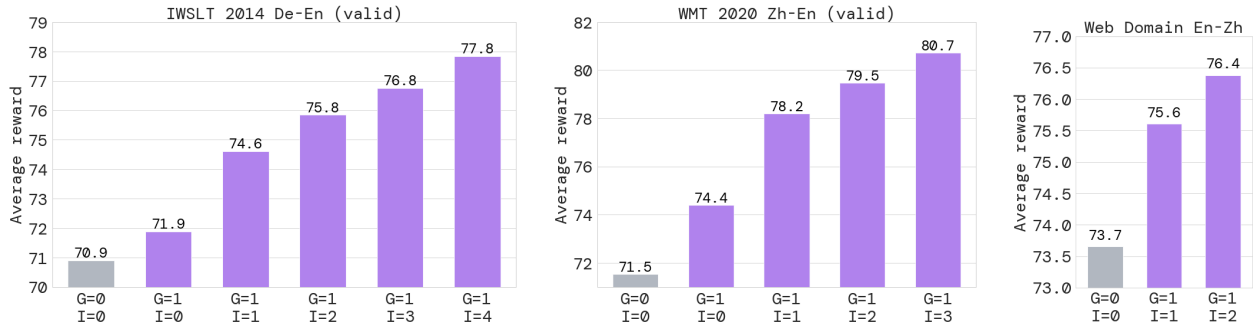


Figure 3 | **ReST with multiple Improve steps.** Average reward model scores on IWSLT 2014 De-En, WMT 2020 Zh-En, and Web Domain En-Zh validation sets. On each dataset, we report results with BC ( $G = 0, I = 0$ ) and ReST with a single Grow step and several Improve steps with an increasing reward threshold. Each Improve step increases the reward model score in all three validation datasets. We found the suitable number of Improve steps to be a dataset-dependent hyperparameter.

as well as an internal benchmark dataset which we call *Web Domain* (a version of this dataset was previously used by Ghorbani et al. (2021)). These datasets contain a set of sentences in the source language and the corresponding human “reference” translation. We selected a different language pair for each dataset to test the generality of the results. We kept a separate validation and test sets with unseen source sentences for the evaluation purposes.

We used Metric X in our experiments, a state-of-art reference-free reward model (Freitag et al., 2022) which, for a given source text and a proposed translation, outputs a numerical score. We report results in terms of average rewards on samples generated by a policy on the validation set <sup>1</sup>. For the details of the datasets and models, we refer to Appendix A.3. Also, Table 2 indicates the size of the datasets by reporting the number of samples per source sentence generated at each Grow step.

**Nomenclature** We named variants of *ReST* by the loss type, number of Grow steps, and number of Improve steps, for example GOLD  $G=1, I=2$ . With this convention, BC  $G=0, I=0$  refers to standard supervised learning, which is trained only on the original dataset  $\mathcal{D}$  and performs neither Grow nor Improve steps. When the loss type is not specified, the BC loss is used, i.e., the model is trained with auto-regressive supervised learning with the NLL loss as typical in training language models. In all plots, we colored supervised learning in grey and *ReST* variants in shades of purple.

**Baselines** We reported the results with several different offline RL method, including Offline Actor Critic (OAC) (Mathieu et al., 2021), Behavior VMPO (BVMP0), Generation by Off-policy Learning from Demonstrations (GOLD) (Pang and He, 2021), and BC (Pomerleau, 1989) <sup>2</sup>.

**Do multiple Improve steps in *ReST* increase the reward model scores?** We evaluated *ReST* on three different datasets by fixing the loss function to BC and increasing the number of Improve steps. The range of rewards for training was normalized between 0 and 1 <sup>3</sup>. For our experiments, we

<sup>1</sup>Performance on the test set follows the same trends (see Appendix A.4). We also experimented with BLEURT and BLEU scores, and *ReST* improved those scores as well. We noticed that online PPO algorithm can learn to exploit the weaknesses and biases of these two metrics quickly, which can cause the model to generate samples that maximize the rewards but deteriorate the quality of the model’s output.

<sup>2</sup>Details on our baselines are in Appendix A.8 and the experiments with additional losses are in Appendix A.2

<sup>3</sup>Note that the plots show rewards between 0 and 100.

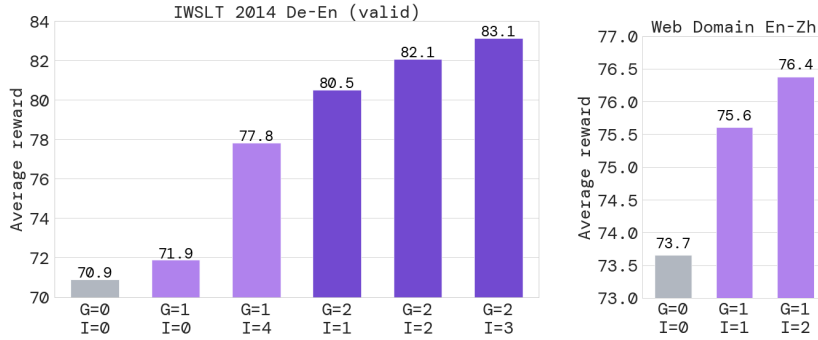
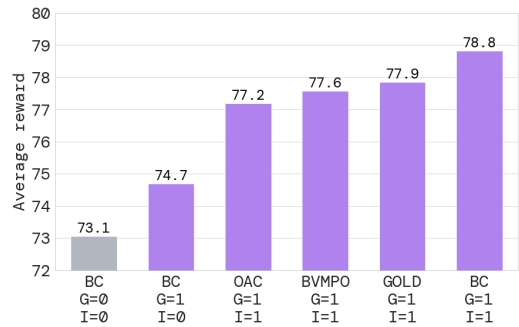


Figure 4 | **ReST with two Grow steps.** The second Grow step with subsequent Improve steps improves the performance by 5.3 points on IWSLT 2014 De-En and 0.8 points on Web Domain En-Zh task over the first Grow step.

picked the filtering thresholds  $\tau_i$  from a sequence of increasing values [0.0, 0.7, 0.8, 0.9, 0.95, 0.99]<sup>4</sup>. The  $\tau_0 = 0.0$  case corresponds to using the full dataset. We did five Improve steps on IWSLT 2014, four on WMT-2020, and two on Web Domain. In Figure 3 we plotted the average reward of different variants of *ReST*. We see that each subsequent Improve step improves the performance of the translation model significantly across all three datasets.

**Do additional Grow steps improve reward model scores?** We performed a second Grow step with successive Improve steps to measure the effect of the extra Grow step on the performance. In Figure 4, a method with an additional Grow step achieves further improvement on the IWSLT 2014 and Web Domain datasets. We noticed a 5.3 point improvement between the end of the first and the second Grow step.



**Does *ReST* improve over supervised training?** To answer this question, in Figure 5 we plotted the average reward achieved by the supervised learning model as well as several variants of *ReST* with different losses and the number of Grow and Improve steps. Different variants of *ReST* (purple) significantly outperform supervised learning (gray) even after just the first grow step. This observation was consistent across different datasets and language pairs that we tested.

Figure 5 | **WMT 2020 zh-en (test):** BC (in grey,  $G = 0$   $I = 0$ ) and *ReST* trained with different offline RL losses. *ReST* is trained with one Grow and Improve step except  $G = 1$   $I = 0$ , which is trained on the entire dataset generated after the first Grow step without any Improve (all in purple). All variants of *ReST* outperform the initial BC baseline, with BC loss resulting in the best performance.

**Which loss is the best for a single step of *ReST*?** Figure 5 depicts variants of *ReST* with different offline RL losses  $\mathcal{L}(x, y; \theta)$ . We find that BC loss outperforms other loss functions. Note that normally BC algorithm does not depend on the reward, but in *ReST*, the reward is taken into account through the reward filtering stage for  $I \geq 1$  (with  $\tau_1 = 0.8$  for WMT 2020.) Results with multiple Grow and Improve steps are displayed in Figure 4 (see also Appendix A.6).

<sup>4</sup>If we used less than 6 steps, we skipped the initial smaller thresholds. Details are given in Appendix A.3. We ensured that in all our datasets  $\tau_1 > 0.7 \geq V_{\pi_\theta}$  by empirically measuring  $V_{\pi_\theta}$  over the dataset.

Algorithm	Average Reward	Distinct samples
BC (G=0, I=0)	70.9	16 000 000
ReST (G=1, I=0)	71.9	16 000 000
ReST (G=1, I=4)	77.8	16 000 000
ReST (G=2, I=3)	<b>83.1</b>	32 000 000
Online RL	71.6	24 000 000

Table 1 | **Online RL for IWSLT 2014:** Online RL performs as well as *ReST* (G=1, I=0) and *ReST* (G=1, I=4) is significantly better.

**Can *ReST* be improved further with Best-of-N sampling at inference time?** Best-of-N sampling technique at inference time generates  $N$  samples which are then ranked by the reward model. Then, the top ranked candidate is selected (Gao et al., 2022). We show results with Best-of-N sampling on top of BC (G=0 I=0) and *ReST* variants in Figure 6. The performance of *ReST* improves both with  $N$  and with the number of Improve steps. The best *ReST* variant with  $N < 10$  matches the performance of the BC model with  $N = 200$ . Even though RL is known to limit the diversity of samples, this experiment shows that *ReST* can still benefit from Best-of-N sampling. After three Improve steps with  $N = 200$ , *ReST* achieves the highest possible reward of 1, outperforming the “reference” translations in  $\mathcal{D}$ .

**How does *ReST* compare with Online RL?** We compared *ReST* with PPO (Schulman et al., 2017), an online RL algorithm widely used for RLHF (Glaese et al., 2022; Ouyang et al., 2022a). For our online RL experiments, we used the setup of Donato et al. (2022) where PPO had access to a similar amount of training data as *ReST* with 1 Grow step. The results are summarized in Table 1. Online RL performs as well as *ReST* with one Grow and no Improve steps which is equivalent to BC on the  $\mathcal{D}_g$  dataset. With the same amount of training data, *ReST* with multiple Improve steps achieves significantly higher rewards. Furthermore, we noticed that the BLEU score for the online RL policy on the validation set dropped by nearly 8 points (BLEU score of *ReST* did not change) which indicates a potential reward hacking behaviour. *ReST*’s ability to improve the reward model score without deteriorating the performance on other metrics suggests that the “alignment tax” it pays is lower than for online RL approaches.

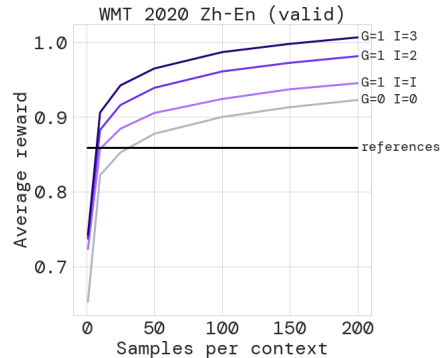


Figure 6 | **Best-of-N sampling at inference time.** All variants of *ReST* benefit as much from Best-of-N sampling as supervised models.

**Does *ReST* improve human preferences?** We evaluated the *ReST* models by human raters to investigate if *ReST* can outperform BC in human evaluations as well. We displayed a sentence in the source language and two generated translations: one by BC model (G=0 I=0) and one by a *ReST* variant. Human raters scored each translation on a scale from 0 to 6, and we measured the difference between the average score of *ReST* method and of BC which we refer as “Human eval diff”. In Figure 7 (right), we see that all variants of *ReST* outperform the BC baseline significantly. However, if we compare the human score gains with the gains in the learned reward (Figure 7, left), the rankings do not match. We hypothesise that the difference is due to the fact that the reward



models cannot generalize well on OOD data since the learned reward models are an imperfect proxy of human preferences. In particular, we found that the reward models generalise worse as our policy moves away from the behaviour model which can happen as the number of Grow and Improve steps increases at which point *ReST* can start overfitting to the reward model. Thus, in our analysis, we focused on evaluating models based on how well they align with a reward signal and we treat reward model generalisation as an independent issue that could be mitigated by, for example, finetuning the reward model between the consecutive Grow steps on the human-annotated data from the most recent policy. *ReST* with the B-VMPO loss utilises a learned value function and KL regularisation term to prevent over-fitting and thus attains high human preference scores.

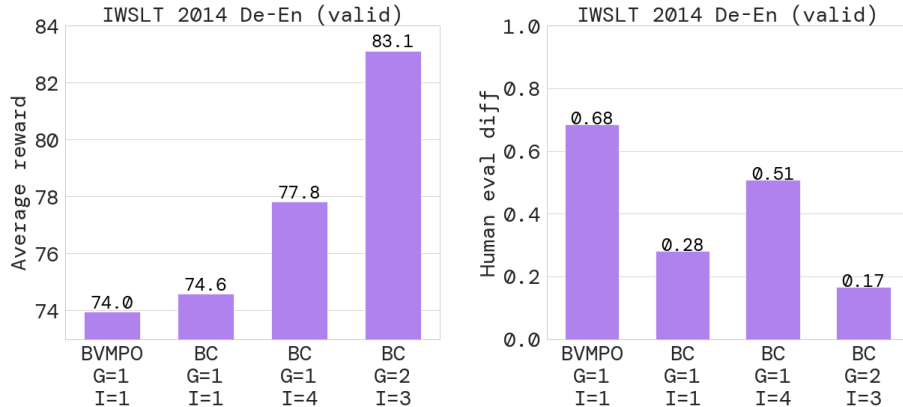


Figure 7 | **Comparison of performance based on learned reward and on human evaluation.** All *ReST* variants outperform BC in terms of human ratings, but the rankings of the methods based on reward model scores and human scores are different.

## 5. Related works

There has been large number of works recently on self-improving alignment algorithms for language modelling. In Figure 8, we also compare *ReST* against different approaches: supervised learning, self-training, online and offline RL. We conclude from this comparison that *ReST* is the only approach that is compute efficient, but also can leverage exploration data and rewards. Next, we describe some of the particular works related to *ReST*.

	Exploration	Rewards	Compute Efficient
Supervised Learning	✗	✗	✓
Self-training	✓	✗	✓
Online RL	✓	✓	✗
Offline RL	✗	✓	✓
<i>ReST</i>	✓	✓	✓

Figure 8 | ***ReST* vs alternatives:** *ReST* is the only approach that can leverage the exploration data and rewards, but is also computationally efficient.

**Self-training** Self-training is an established semi-supervised learning approach which utilizes unlabeled data to improve a model (Scudder, 1965). Since its introduction, self-training has been successfully applied to many tasks, including image classification and recognition (Xie et al., 2020), protein folding (Jumper et al., 2021), as well as several language tasks (He et al., 2019; Sun et al., 2021; Yarowsky, 1995; Zhang and Zong, 2016). He et al. (2019) empirically demonstrated that noisy self-training improves the performance of translation models. The Improve step of *ReST* resembles self-training. *ReST*’s main difference from self-training is that the Grow step of *ReST* generates synthetic exploration data for training with RL.

**Expert Iteration (EI)** Anthony (2021) proposed an RL framework which is a form of policy iteration approach that makes use of a planning mechanism. EI explicitly decomposes the RL problem into two parts: planning and generalisation. Similar to *ReST*, EI uses the policy to generate data and exploit it to learn a policy with RL. Unlike EI, *ReST* does not require any planning mechanism and it makes use of iterative Improve steps that enables it to leverage the gathered data more effectively.

**Reasoning with language models** The EI method inspired several related approaches (Uesato et al., 2022; Zelikman et al., 2022). Zelikman et al. (2022) proposed a technique called STAR that iteratively leverages a small number of rationales to fine-tune the model. Then, they sample rationales with the answers from the model and filter the generated answers by their correctness. Uesato et al. (2022) proposed a similar method which learns to solve math problems using learned reward models. Recently, Jung et al. (2023) proposed a method called “Impossible Distillation” similar to ours, which generates datasets for large language models by sampling from a suboptimal model and filters the low-quality examples with a filtering mechanism. This approach corresponds to *ReST* with a single Grow and Improve steps. In contrast to these methods, *ReST* can be trained with any offline RL losses, with or without planning and various filtering mechanisms. It can also deal with continuous-valued reward scores. Furthermore, *ReST* can perform iterative policy improvement on a fixed dataset with the Improve steps.

**Iterated Learning (IL)** IL is the process where an agent learns its behavior by being exposed to another agent’s behavior, which itself learned it in the same way (Kirby et al., 2014). Recently, this approach was adopted for interactive language learning using deep learning (Lu et al., 2020a,b). IL differs from *ReST* as it operates in a multi-agent setting and does not use RL.

**Self Imitation Learning (SIL)** SIL learns a policy for an off-policy actor-critic algorithm where the policy tries to reproduce the good behaviour demonstrated by the agent (Oh et al., 2018). SIL achieves it by filtering out the unsuccessful trajectories from the replay buffer and training the agent only on the high-reward trajectories. In that sense, *ReST* can be considered to be closely related to the SIL-based approaches. The main difference is that *ReST* is agnostic to the underlying RL algorithm used to train the policy and, unlike SIL, does not necessitate a value function to filter out the unsuccessful trajectories. Also, *ReST* is applied to generative AI settings, which do not require interactions with an environment in an online fashion.

**Reward ranked Fine-Tuning (RAFT)** Concurrently to our work, Dong et al. (2023) proposed RAFT. RAFT can be interpreted as a particular case of *ReST* which uses only one Improve step for each Grow step, and relies on a filtering threshold which is a fixed quantile of the empirical distribution of the rewards of the current samples. The authors reported improvements over BC and PPO on a variety of language modeling and image generation tasks. Our experiments with *ReST* showed that multiple Improve steps with an increasing filtering threshold for one Grow step lead to further performance improvements.

## 6. Discussion

In this paper, we proposed an algorithm called *ReST* which is simple, has minimal hyper-parameters to tune, and is flexible to work with many designs of Grow and Improve steps. We studied the performance of *ReST* in machine translation as robust and established reward models are available for

this task. We experimented with different offline RL losses in the *ReST* loop, but found BC to perform the best for improving the reward model scores. Multiple steps of NLL training with progressively increasing filtering thresholds in the Improve step lead to continuous improvements in the model’s reward on the holdout set. However, improvements in reward model scores do not necessarily reflect human preferences since the reward model is only a proxy for human preferences. The results indicate that one Grow step is the best option when considering human evaluation scores, even though rewards continue to grow with more Grow steps. To overcome this limitation, the reward models could be fine-tuned on the subset of  $D_g$  annotated with human preferences similar to Bai et al. (2022) and Glaese et al. (2022), which we leave as future work. Let us note that the risk of overfitting to the reward model increases with the repeated iterations of the Grow steps; thus we believe it is essential to address this issue, especially in cases where multiple Grow steps are needed to train the model.

As we have seen, simple BC loss still outperforms many offline RL losses in terms of aligning the model with the reward model scores. However, we found that BC can overfit to the reward model so we explain it by the fact that learning value functions in RL is challenging due to sparse rewards, credit assignment problems, sensitivity to hyperparameters, and limited exploration in the Grow step. *ReST* could benefit from better RL exploration strategies at Grow step, such as MCTS (Leblond et al., 2021). The ability to exploit the generated data during the Grow step could result in a broader exploration of the state-action space and better generalization. Additionally, the determinism of the environment does not allow for large gains over BC for offline RL losses.

To conclude, *ReST* is a general and efficient approach. It can be applied when 1) a robust reward model of human preferences is available and 2) we are able to generate samples from the model at scale. Thus, it can be applied to many tasks within the language domain, such as summarization, turn-based dialogue, and other generative audio and video models. With several avenues for future exploration and applications, we believe that *ReST* is a useful growing batch RL methodology for RLHF.

**Acknowledgements** We would like to thank the members from the machine translation teams at Google and Google DeepMind for their inputs to the project during the brainstorming phases and for setting up the codebase that this project was developed upon (Yu et al., 2020). We would like to thank Matt Hoffman, Bobak Shahriari, Taylan Cemgil and Chris Dyer for the discussions about this project. We are grateful for the feedback provided by Bilal Piot for an early draft of this paper. We would also like to thank those responsible for various different frameworks that we used during the project such as the DeepMind JAX ecosystem (Babuschkin et al., 2020) and Launchpad (Yang et al., 2021).

## References

- A. Abdolmaleki, S. Huang, G. Vezzani, B. Shahriari, J. T. Springenberg, S. Mishra, D. Tirumala, A. Byravan, K. Bousmalis, A. György, et al. On multi-objective policy optimization as a tool for reinforcement learning. *arXiv preprint arXiv:2106.08199*, 2021.
- R. Agarwal, M. Schwarzer, P. S. Castro, A. Courville, and M. G. Bellemare. Beyond tabula rasa: Reincarnating reinforcement learning. *arXiv preprint arXiv:2206.01626*, 2022.
- T. W. Anthony. *Expert iteration*. PhD thesis, UCL (University College London), 2021.
- I. Babuschkin, K. Baumli, A. Bell, S. Bhupatiraju, J. Bruce, P. Buchlovsky, D. Budden, T. Cai, A. Clark, I. Danihelka, C. Fantacci, J. Godwin, C. Jones, R. Hemsley, T. Hennigan, M. Hessel, S. Hou,

- S. Kapturowski, T. Keck, I. Kemaev, M. King, M. Kunesch, L. Martens, H. Merzic, V. Mikulik, T. Norman, J. Quan, G. Papamakarios, R. Ring, F. Ruiz, A. Sanchez, R. Schneider, E. Sezener, S. Spencer, S. Srinivasan, L. Wang, W. Stokowiec, and F. Viola. The DeepMind JAX Ecosystem, 2020. URL <http://github.com/deepmind>.
- Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- R. Bellman. A Markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684, 1957.
- D. Brandfonbrener, A. Bietti, J. Buckman, R. Laroche, and J. Bruna. When does return-conditioned supervised learning work for offline reinforcement learning? *Advances in Neural Information Processing Systems*, 35:1542–1553, 2022.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 2020.
- S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, et al. Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*, 2023.
- M. Cettolo, J. Niehues, S. Stüker, L. Bentivogli, and M. Federico. Report on the 11th iwslt evaluation campaign. In *Proceedings of the 11th International Workshop on Spoken Language Translation: Evaluation Campaign*, 2014.
- L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In *Advances in Neural Information Processing Systems*, 2021.
- S. Clark, J. R. Curran, and M. Osborne. Bootstrapping pos-taggers using unlabelled data. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003*, pages 49–55, 2003.
- D. Donato, L. Yu, W. Ling, and C. Dyer. Mad for robust reinforcement learning in machine translation. *arXiv preprint arXiv:2207.08583*, 2022.
- H. Dong, W. Xiong, D. Goyal, R. Pan, S. Diao, J. Zhang, K. Shum, and T. Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.
- L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, S. Legg, and K. Kavukcuoglu. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International Conference on Machine Learning*, 2018.
- M. Freitag, R. Rei, N. Mathur, C.-k. Lo, C. Stewart, G. Foster, A. Lavie, and O. Bojar. Results of the wmt21 metrics shared task: Evaluating metrics with expert-based human evaluations on ted and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, 2021.
- M. Freitag, R. Rei, N. Mathur, C.-k. Lo, C. Stewart, E. Avramidis, T. Kocmi, G. Foster, A. Lavie, and A. F. T. Martins. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), Dec. 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.wmt-1.2>.

- J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4RL: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- L. Gao, J. Schulman, and J. Hilton. Scaling laws for reward model overoptimization. *arXiv preprint arXiv:2210.10760*, 2022.
- B. Ghorbani, O. Firat, M. Freitag, A. Bapna, M. Krikun, X. Garcia, C. Chelba, and C. Cherry. Scaling laws for neural machine translation. *arXiv preprint arXiv:2109.07740*, 2021.
- A. Glaese, N. McAleese, M. Trębacz, J. Aslanides, V. Firoiu, T. Ewalds, M. Rauh, L. Weidinger, M. Chadwick, P. Thacker, et al. Improving alignment of dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- C. Gulcehre, S. G. Colmenarejo, Z. Wang, J. Sygnowski, T. Paine, K. Zolna, Y. Chen, M. Hoffman, R. Pascanu, and N. de Freitas. Regularized behavior value estimation. *arXiv preprint arXiv:2103.09575*, 2021.
- J. He, J. Gu, J. Shen, and M. Ranzato. Revisiting self-training for neural sequence generation. *arXiv preprint arXiv:1909.13788*, 2019.
- J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, J. Welbl, A. Clark, et al. An empirical analysis of compute-optimal large language model training. In *Advances in Neural Information Processing Systems*, 2022.
- J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- J. Jung, P. West, L. Jiang, F. Brahma, X. Lu, J. Fisher, T. Sorensen, and Y. Choi. Impossible distillation: from low-quality model to high-quality dataset & model for summarization and paraphrasing. *arXiv preprint arXiv:2305.16635*, 2023.
- S. Kirby, T. Griffiths, and K. Smith. Iterated learning and the evolution of language. *Current Opinion in Neurobiology*, 28:108–114, 2014.
- P. Koehn, V. Chaudhary, A. El-Kishky, N. Goyal, P.-J. Chen, and F. Guzmán. Findings of the wmt 2020 shared task on parallel corpus filtering and alignment. In *Proceedings of the Fifth Conference on Machine Translation*, pages 726–742, 2020.
- T. Kudo. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 66–75, 2018.
- A. Kumar, J. Hong, A. Singh, and S. Levine. Should i run offline reinforcement learning or behavioral cloning? In *International Conference on Learning Representations*, 2021.
- S. Lange, T. Gabel, and M. Riedmiller. Batch reinforcement learning. In M. Wiering and M. van Otterlo, editors, *Reinforcement Learning: State-of-the-Art*, pages 45–73. Springer Berlin Heidelberg, 2012.
- R. Leblond, J.-B. Alayrac, L. Sifre, M. Pislari, L. Jean-Baptiste, I. Antonoglou, K. Simonyan, and O. Vinyals. Machine translation decoding beyond beam search. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2021.

- Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- Y. Lu, S. Singhal, F. Strub, A. Courville, and O. Pietquin. Countering language drift with seeded iterated learning. In *International Conference on Machine Learning*, 2020a.
- Y. Lu, S. Singhal, F. Strub, O. Pietquin, and A. Courville. Supervised seeded iterated learning for interactive language learning. *arXiv preprint arXiv:2010.02975*, 2020b.
- M. Mathieu, S. Ozair, S. Srinivasan, C. Gulcehre, S. Zhang, R. Jiang, T. Le Paine, K. Zolna, R. Powell, J. Schrittwieser, et al. Starcraft ii unplugged: Large scale offline reinforcement learning. In *Deep RL Workshop NeurIPS 2021*, 2021.
- V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Harley, T. P. Lillicrap, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Learning Representations*, 2016.
- J. Oh, Y. Guo, S. Singh, and H. Lee. Self-imitation learning. In *International Conference on Machine Learning*, pages 3878–3887. PMLR, 2018.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with human feedback, 2022a.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022b.
- R. Y. Pang and H. He. Text generation by learning from demonstrations. In *International Conference on Learning Representations*, 2021.
- E. Perez, S. Huang, F. Song, T. Cai, R. Ring, J. Aslanides, A. Glaese, N. McAleese, and G. Irving. Red teaming language models with language models. *arXiv preprint arXiv:2202.03286*, 2022.
- D. A. Pomerleau. ALVINN: An autonomous land vehicle in a neural network. In *Advances in Neural Information Processing Systems*, pages 305–313, 1989.
- J. W. Rae, S. Borgeaud, T. Cai, K. Millican, J. Hoffmann, F. Song, J. Aslanides, S. Henderson, R. Ring, S. Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.
- R. Rei, C. Stewart, A. C. Farinha, and A. Lavie. COMET: A neural framework for MT evaluation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2020.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- H. Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371, 1965.

- T. Sellam, D. Das, and A. Parikh. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- I. Shumailov, Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, and R. Anderson. The curse of recursion: Training on generated data makes models forget. *arXiv preprint arxiv:2305.17493*, 2023.
- J. Skalse, N. H. R. Howe, D. Krasheninnikov, and D. Krueger. Defining and characterizing reward hacking. In *Advances in Neural Information Processing Systems*, 2022.
- H. F. Song, A. Abdolmaleki, J. T. Springenberg, A. Clark, H. Soyer, J. W. Rae, S. Noury, A. Ahuja, S. Liu, D. Tirumala, N. Heess, D. Belov, M. Riedmiller, and M. M. Botvinick. V-mpo: On-policy maximum a posteriori policy optimization for discrete and continuous control. In *International Conference of Learning Representations*, 2020.
- A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shoeb, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022.
- N. Stiennon, L. Ouyang, J. Wu, D. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and P. F. Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 2020.
- H. Sun, R. Wang, K. Chen, M. Utiyama, E. Sumita, and T. Zhao. Self-training for unsupervised neural machine translation in unbalanced training data scenarios. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021.
- I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, 2014.
- J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- J. Wu, L. Ouyang, D. M. Ziegler, N. Stiennon, R. Lowe, J. Leike, and P. Christiano. Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862*, 2021.
- Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10687–10698, 2020.
- F. Yang, G. Barth-Maron, P. Stańczyk, M. Hoffman, S. Liu, M. Kroiss, A. Pope, and A. Rustemi. Launchpad: A programming model for distributed machine learning research. *arXiv preprint arXiv:2106.04516*, 2021.
- D. Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *33rd Annual Meeting of the Association for Computational Linguistics*, 1995.
- L. Yu, L. Sartran, P.-S. Huang, W. Stokoweic, D. Donato, S. Srinivasan, A. Andreev, W. Ling, S. Mokra, A. D. Lago, Y. Doron, S. Young, P. Blunsom, and C. Dyer. The DeepMind chinese–english document translation system at wmt2020. In *Proceedings of the Fifth Conference on Machine Translation*, 2020.

E. Zelikman, Y. Wu, J. Mu, and N. Goodman. Star: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*, 2022.

J. Zhang and C. Zong. Exploiting source-side monolingual data in neural machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2016.

## A. Appendix

### A.1. RLHF for conditional language modeling as MDP

We can formulate conditional language modeling as a sequence to sequence problem. The goal is to map a source sequence  $\mathbf{x} = (x_1, x_2, \dots, x_L)$  into a target sequence  $\mathbf{y} = (y_1, y_2, \dots, y_T)$ , that is to learn a mapping from  $\mathbf{x}$  to  $\mathbf{y}$ . Machine translation is a classic example of a sequence to sequence problem (Sutskever et al., 2014).

The model can be described as a Markov Decision Process (MDP) (Bellman, 1957). An MDP,  $\mathcal{M} \stackrel{\text{def}}{=} (\mathcal{S}, \mathcal{A}, T, r, d)$ , consists of finite sets of states  $\mathcal{S}$  and actions  $\mathcal{A}$ , a transition distribution  $T(s'|s, a), s, s' \in \mathcal{S}, a \in \mathcal{A}$ , a reward function  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , and an initial state distribution  $d : \mathcal{S} \rightarrow [0, 1]$ . In the offline setting the agent learns from a dataset  $\mathcal{D}$  which contains sequences  $(s_n, a_n, r_{n+1})$ . The dataset  $\mathcal{D}$  is often assumed to be generated by an unknown *behaviour policy*  $\mu$  characterised by a distribution over actions conditioned on the state:  $\mu(a|s)$ . The elements of MDP have the following meaning in conditional language modeling:

- **States** ( $s$ ): The state  $s_n$  consists of the input sequence and the partially generated output up to the  $n$ -th token:  $s_n = [\mathbf{x}, y_{1:n-1}]$ .
- **Actions** ( $a$ ): The actions are the tokens  $a_n = y_n$  to be generated by policy ( $\pi$ ).
- **Rewards** ( $r$ ): Rewards can be given or learned. In our experiments we train a deep neural network on human preferences to assigns a score to the full generated sequence. Reward is produced after end-of-sequence (EOS) token, and it is zero at all other steps.
- **Transition operator** ( $T$ ): The transition operator defines the dynamics of an environment:  $T(s'|a, s) := T(y_n | s_n = [\mathbf{x}, y_{1:n-1}])$ . In our setup, it is a deterministic operator that concatenates newly produced token (action) to the current sequence (state).

This RLHF formulation of conditional language modeling can be also seen as a contextual bandit problem with a very large action space.

### A.2. Negative results with offline RL

In addition to the experiments from Section 4, we also run experiments with Q-learning with BVE type of one-step RL approaches and reward conditioned approaches, such as decision transformers (Chen et al., 2021). However, we could not obtain notable improvements with them over the supervised baseline. Q-function-based methods for learning a policy and value function performed worse than supervised learning when used as a policy even after training from initialization from a supervised checkpoint. This matches previous observations by Mathieu et al. (2021) who showed that offline Q-learning has difficulty in large action spaces (vocabulary size of 32 000 tokens) and state-value based methods are preferable. Furthermore, as Brandfonbrener et al. (2022) showed, return-conditioned methods tend to learn sub-optimal policies on tasks with sparse continuous rewards.



### A.3. Data and Model details

In all our experiments, the base policy architecture is a minor modification of a standard Transformer architecture (Vaswani et al., 2017). We use a vocabulary of 32 000 tokens and a max sequence length of 128 at decoding time. We use the `sentencepiece` tool (Kudo, 2018) for building the tokenizers.

**Grow step** During the `Grow` step, we sampled from the latest checkpoint of the policy with tempered softmax using temperature 0.8 following the procedure proposed by Li et al. (2022) to generate the dataset. Moreover, in our analysis, we found that temperature 0.8 often covers a broad range of rewards in the dataset.

**Thresholds in Improve step** If we set the threshold  $\tau$  to a value such that  $\tau \geq V^{\pi_0}$ , and learn a policy  $\pi^1$  on data filtered with this threshold, the updated average value  $V^{\pi_1}$  of this policy satisfies  $V^{\pi_1} > V^{\pi_0}$ . Then, iterative threshold ensures policy improvement in the `Improve` step. In our experiments, we do several improvement steps until we reach the maximum threshold. At each step, we trained the model for 500 000 SGD steps and pick the checkpoint with the best reward model score.

**IWSLT 2014 De-En** We use train, validation and test sets from IWSLT 2014 De-En dataset (Cettolo et al., 2014) which includes source sentences in German with human translations (references) in English. In each `Grow` step, we generate 100 candidate translation for each source sentence in the training set, effectively yielding us  $|\mathcal{D}_g| = 16\,000\,000$ . For this dataset, we use a tiny version of the Transformer encoder-decoder (Vaswani et al., 2017) architecture with the feedforward MLP layers of size 512, feedforward dimension of 1024, 4 attention heads and 6 encoder and decoder layers.

**WMT 2020 Zh-En** We use the source-reference pairs in Chinese and English from the work of Koehn et al. (2020) for our training, validation and test sets. Exact details on the datasets and preprocessing can be found in Yu et al. (2020). In each `Grow` step, we generate 25 candidates for each source sentence in the training set, effectively yielding us  $|\mathcal{D}_g| = 890\,000\,000$ . We choose to generate 25 candidates per source as the WMT dataset is significantly ( $\sim 100$  times) larger than the IWSLT dataset. We use an architecture that mimics the Transformer-base encoder-decoder (Vaswani et al., 2017) architecture with model dimension 512, feedforward dimension of 2048, 8 attention heads and 6 encoder and decoder layers.

**Web Domain En-Zh** Finally, we use the in-house dataset for English to Chinese translation with custom training, fine-tuning and test sets. This training corpus is our biggest dataset<sup>5</sup>, so we use a modified version of Transformer-big encoder-decoder (Vaswani et al., 2017) architecture with model dimension 1024, feedforward dimension of 8192, 16 attention heads and 6 encoder and decoder layers.

In Table 2, we list all the datasets with their sizes. In all the experiments, unless stated otherwise, we report the average reward scores on the validation set.

<sup>5</sup>For computational reasons, we run *ReST* on the fine-tuning corpus with a model warm-started from a supervised model trained on the entire massive training corpus.

Dataset	$ \mathcal{D}_g $	# Eval samples	# Candidates per source
IWSLT 2014 de-en	16 000 000	7466	100
WMT 2020 zh-en	890 000 000	2000	25
Web Domain Finetuning en-zh	3 000 000	6000	1000

Table 2 | Details of the datasets and their sizes used in *ReST* experiments.

#### A.4. Reward model

We used learned reward models that assign a score to the whole translation. We considered two types of reward models for translation (Freitag et al., 2022): i) *Reference-free* reward models estimate of how good a translation is based on the source and candidate sentences only, ii) *Reference-based* reward models additionally use the reference human translation to decide how good a candidate translation is. The reference-free reward models are more flexible since they do not require reference translations. Moreover, by their design, reference-based models assign higher scores for sentences similar to the reference, and this could happen even when the reference is erroneous or incomplete. Finally, reference-free models open possibilities to evaluate and discover candidate translations that supersede the quality of references. Thus, in our experiments, we chose to work with reference-free reward models. On the other hand, since the reference-free reward models are not grounded in a human reference translation, it is more vulnerable to the distribution shifts in the candidates generated by the model and reward hacking. In practice, we pre-computed and stored  $R(x, y)$  for the generated  $\mathcal{D}_g$ .

During the development of the model, we relied on the methodology of “unit tests” for reward models where we tested them in various hand-crafted settings, such as on various permutations and repetitions of sentences. The unit-tests ensured the high quality of rewards generated for *ReST*. Even though we used the most robust available reference-free reward model, we found that it was still not perfect, and sometimes showed signs of delusions. For example, the reward of a translation occasionally increased when we repeated the sentence, independent of the initial quality of the translation.

#### A.5. Alignment between human evaluation scores and the reward model scores

In Figure 9 and 10, we showed the distribution of human preferences and reward model scores for *ReST* with BC and  $G=1$   $I=4$ . Here, we found out that our reward model has a very high variance on samples with low human preference score from *ReST* (BC,  $G=1$   $I=4$ ) and our supervised learning baseline (BC,  $G=0$   $I=0$ ). Our hypothesis for this is that the training dataset for the reward model was dominated by translations of high quality<sup>6</sup>. Thus the model overfitted to the samples with high scores, and it has high variance in its prediction over the samples with lesser quality. This gap can be addressed with the incremental retraining of the reward model.

#### A.6. Results on test sets

Figures 11 and 12 present the results on the test sets in IWSLT 2014 De-En and WMT 2020 Zh-En datasets. Similar to the validations set, *ReST* outperforms other value-based offline RL baselines. Also, increasing the number of both Grow and Improve steps helps to further increase the rewards.

<sup>6</sup>The data came from WMT 2021 metrics shared task dataset (Freitag et al., 2021), which reviews the human evaluation of good translation models

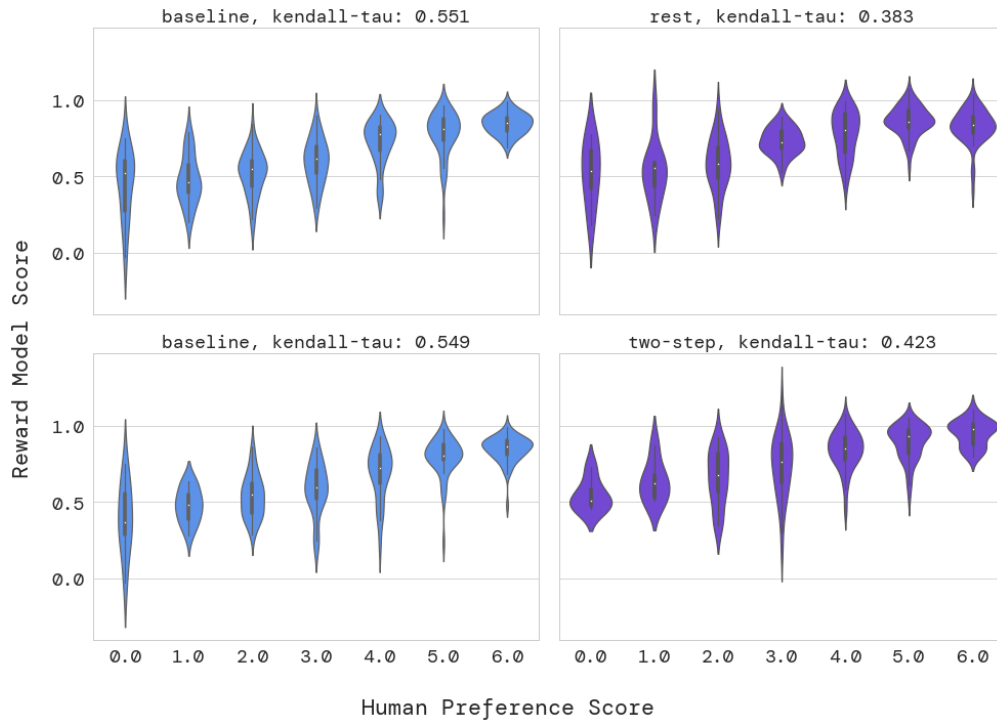


Figure 9 | [IWSLT 2014 De-En]: distribution of human preference and reward model scores for *ReST* (BC,  $I=4$ ,  $G=1$ ) in side by side evaluation with supervised model. The human preference scores less than or equal to 3 have significantly higher variance in rewards compared to the human preference scores above 3. The plots on the left-hand side are for the samples generated from a supervised baseline and the right-hand side are for the samples generated with *ReST*.

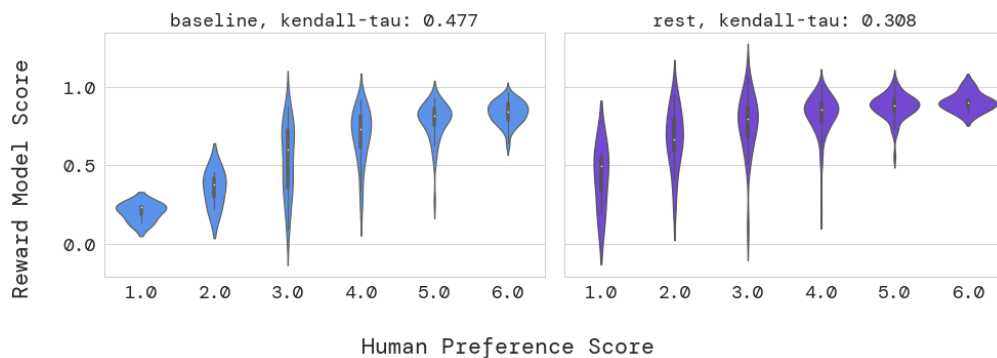


Figure 10 | [WMT 2020 Zh-En]: distribution of human preference and reward model scores for *ReST* (BC,  $I=4$ ,  $G=1$ ) in side by side evaluation with supervised model. The human preference scores lower than or equal to 3 has a variance in terms of reward model scores. The reward model has less certainty on the scores of the candidates having lower scores.

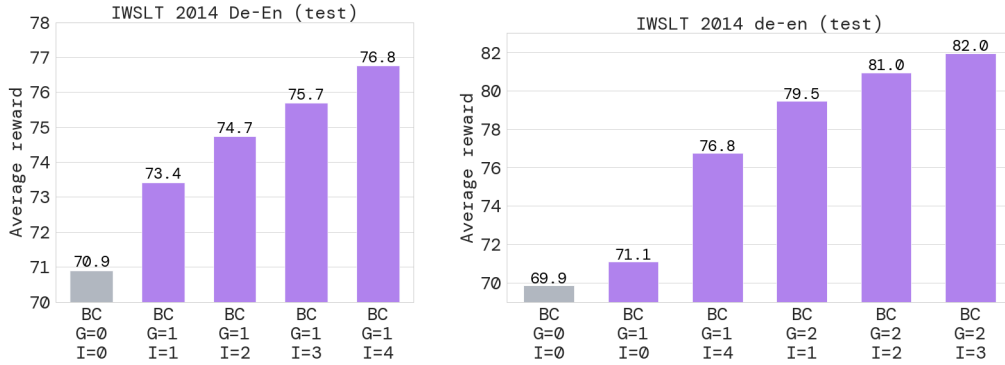


Figure 11 | [IWSLT 2014 De-En]: *ReST* results on test set. All variations of *ReST* (purple) outperform supervised learning baseline (grey). As we increase the number of Grow and Improve steps, the reward of the model on the test set continues to increase.

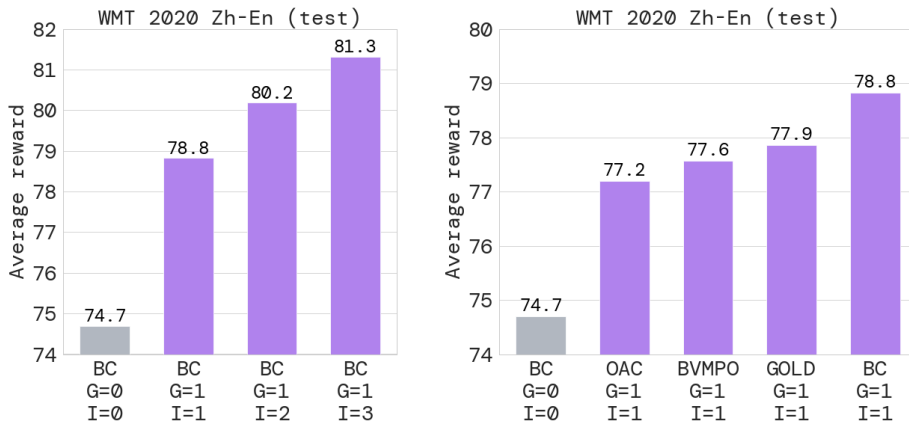


Figure 12 | [WMT 2020 Zh-En]: *ReST* results on test set. We see that *ReST* clearly outperforms the baselines (right) and the reward of the model on the test set increases with the number of Improve steps.

## A.7. Additional experiments

This section reports the results of several ablations in *ReST* which explain some of the design choices that we made.

**GOLD loss compared to BC loss in Improve step** In Figure 13, we compare GOLD against BC losses with the increasing number of *ReST* Improve steps on WMT 2020 Zh-En dataset. The performance of both BC and GOLD improves with the number of Improve steps, but BC performs better than GOLD in every setting.

**Loss comparison on IWSLT 2014 De-En** In Figure 14, we compare BVMPO, BC and GOLD losses in Improve step on IWSLT 2014. The results are consistent with WMT dataset: in short, *ReST* with BC loss and multiple Improve steps outperforms other approaches.

**Selecting threshold based on percentiles of reward model scores** Using a single threshold for all the source-candidate pairs may lead to a scenario with no training data for certain (harder) source

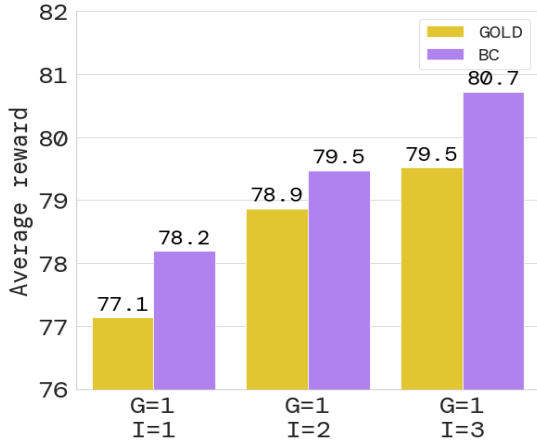


Figure 13 | [WMT 2020 Zh-En]: GOLD loss vs BC loss in multi-step improvement. The performance of all methods improves with more steps, but BC loss always performs better than other offline RL losses that we tried.

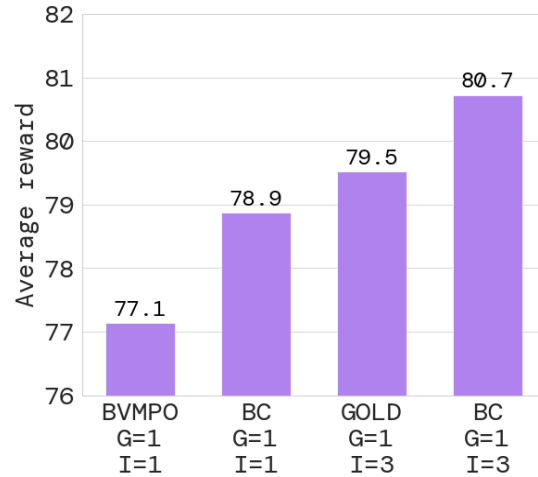


Figure 14 | [IWSLT 2014 De-En]: BVMPO, BC, GOLD losses used in ReST. BC loss with three Improve steps yields the best results.

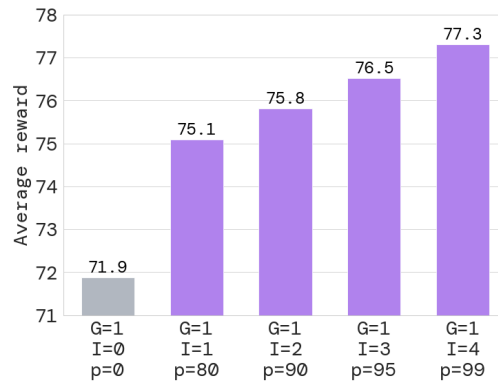
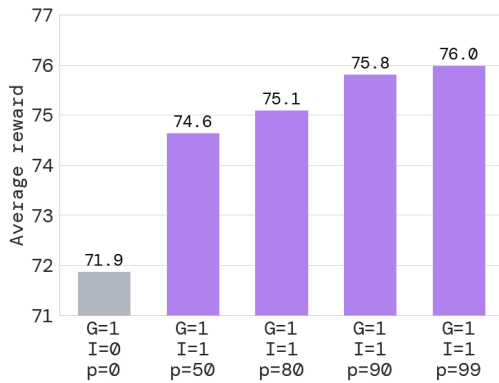


Figure 15 | [IWSLT 2014 De-En] percentiles: Filtering data based on percentiles of reward scores for each source candidate. ReST performs better with increasing percentile and the performance saturates after  $p = 90$ .

sentences or contexts. Alternatively, we can filter the candidate sentences based on the percentile of the scores given their source. This way of filtering based on percentiles effectively results in changing the thresholds for each source sentence. The results are presented in Figure 15. As the filtering percentile increases, the performance increases, but it starts saturating, and the gains become smaller after  $p = 90$ . The results are generally similar to those with global threshold-based filtering.

**Selecting threshold based on linear interpolation between the mean and max scores of candidates** An alternative way of specifying filters based on a source sentence is to compute the filter based on a max reward score ( $V^{\max}(\text{source})$ ) and the mean candidate score ( $V^{\text{mean}}(\text{source})$ ) for a give source sentence. Then, we select the threshold for each source sentence as  $\tau_{\gamma}(\text{source}) = \gamma V^{\max}(\text{source}) + (1 - \gamma)V^{\text{mean}}(\text{source})$ . At each Improve step, we increase the parameter  $\gamma$ . We report the results with respect to different  $\gamma$  values in Figure 16. In a nutshell, computing the threshold by interpolating the max and mean scores for a given candidate gives results similar to the percentile-based way of computing the thresholds per source. Also we can see that the schedule of thresholds

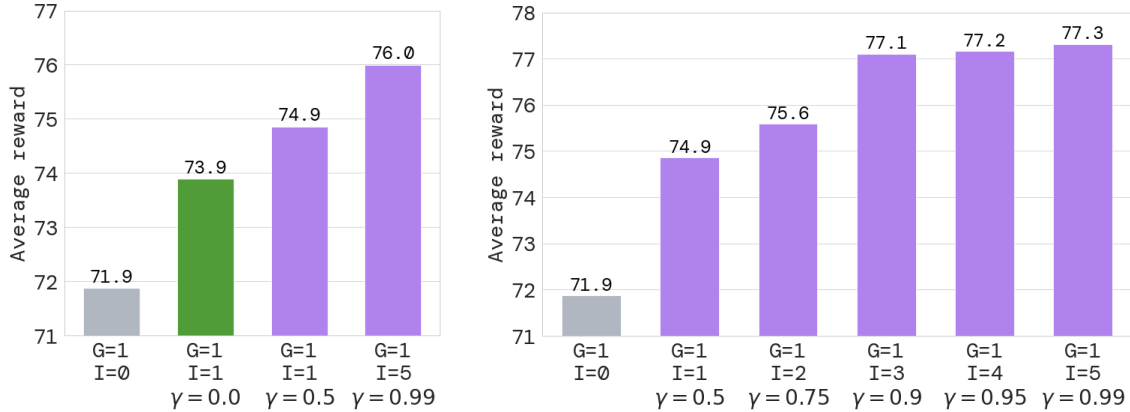


Figure 16 | [IWSLT 2014 De-En] interpolation experiments: Threshold in Improve step is chosen for each candidate sentence as  $\gamma V^{\max}(\text{source}) + (1 - \gamma)V^{\text{mean}}(\text{source})$ . On the left, experiments with  $I=1$  provide improvements and increasing  $\gamma$ , which also helps substantially. On the right, we present the results with  $\gamma$  starting with 0.5 instead of 0.5. As it can be seen from this figure, the number of Improve steps and how the thresholds are selected greatly influence the final performance of the model.

and the number of Improve steps have a large impact on the final performance of the model.

## A.8. Baseline Algorithms

This section provides details on the baseline losses used in Figure 5.

### A.8.1. BVMPO

The BVMPO approach is similar to DIME proposed by Abdolmaleki et al. (2021). The main difference is that we use a state-value function instead of Q function with v-trace (Espeholt et al., 2018), similarly to V-MPO (Song et al., 2020). We train separate neural networks for policy and value function. The policy is pre-trained with BC and the value network is pre-trained with BVE to estimate the behavior value  $V^\beta$ . BVMPO uses the behavioral cloning policy  $\pi^\beta(a|s)$  as a behaviour prior:

$$\mathcal{L}_{BVMPO} = \mathcal{L}_{VMPO} + \lambda \text{KL}(\pi^\beta | \pi). \quad (4)$$

As opposed to DIME, BVMPO starts with a fixed coefficient  $\lambda = 1$  and anneals it to  $1e - 5$  during the training. We found that starting with a large  $\lambda$  and annealing it stabilizes the training without any noticeable drop in performance. Similarly, annealing the behavior regularization co-efficient was also shown to be effective in offline Q-learning (Agarwal et al., 2022).

### A.8.2. GOLD

The GOLD loss was introduced by Pang and He (2021). When using it in Improve step, we start with a BC policy. In all GOLD experiments, we used  $k = 5$  steps to unroll the MLE model into the future for reward computations. Increasing  $k$  did not improve the performance.

### A.8.3. Offline Actor Critic (OAC)

Offline actor critic is an actor critic approach proposed by Mathieu et al. (2021) where the state-value function is trained to estimate the behavior value in the dataset (Gulcehre et al., 2021). Unlike

Mathieu et al. (2021) we did not use the V-trace as it did not improve over the vanilla advantage function in our initial experiments.

#### A.8.4. Value Functions

We trained state-value functions on IWSLT 2014 De-En and WMT 2020 Zh-En datasets after the initial grow step. The value functions estimate the behavior value of the policy on the training dataset. The value functions predict the sum of discounted rewards for each word until the end of the sentence. We used a transformer-based encoder-decoder architecture for the value functions. The encoder gets the source sentence as the input, and the decoder predicts Monte Carlo discounted returns conditioned on the source and history of generated sequences of tokens in the target.

On the training sets of IWSLT and WMT, our value functions achieved Kendall-tau correlation of 92.4% and Spearman correlation above 99% with respect to the reward model scores in the dataset. As we do RL only offline, the value functions do not need to generalize beyond the training dataset.

#### A.9. ReST: Population interpretation

Here, we provide a population interpretation of the *ReST* algorithm. In this setup, the initial BC training step can be thought of as identifying  $\theta$  by minimizing

$$\begin{aligned} \text{KL}(p(\mathbf{x}, \mathbf{y}) || \pi_{\theta}(\mathbf{x}, \mathbf{y})) &= \sum_{\mathbf{x}, \mathbf{y}} p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{x}, \mathbf{y})}{\pi_{\theta}(\mathbf{x}, \mathbf{y})} \\ &= \sum_{\mathbf{x}, \mathbf{y}} p(\mathbf{x}, \mathbf{y}) \log \frac{p(\mathbf{y}|\mathbf{x})}{\pi_{\theta}(\mathbf{y}|\mathbf{x})} \\ &= \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\text{KL}(p(\mathbf{y}|\mathbf{x}) || \pi_{\theta}(\mathbf{y}|\mathbf{x}))], \end{aligned}$$

where KL is the Kullback–Leibler discrepancy,  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$  is the (training) data distribution and  $\pi_{\theta}(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})\pi_{\theta}(\mathbf{y}|\mathbf{x})$ . Let us call the minimizer  $\theta_0$ .

Assume now that one has the parameter  $\theta_k$  at the end of Improve steps of the  $k$ th Grow step. Then at the next  $(k + 1)$ th Grow step, we consider Improve steps using a sequence of increasing thresholds  $(\tau_i)_{i=1}^I$  which are used to filter the available data distribution (mixture of original data distribution and synthetic data distribution) to obtain a new distribution

$$\pi_{\theta_k}^{\tau_i}(\mathbf{x}, \mathbf{y}) \propto \{(1 - \lambda)p(\mathbf{x}, \mathbf{y}) + \lambda\pi_{\theta_k}(\mathbf{x}, \mathbf{y})\} \mathbb{1}(R(\mathbf{x}, \mathbf{y}) > \tau_i).$$

We then obtain the next parameter  $\theta$  by minimizing the KL

$$\text{KL}(\pi_{\theta_k}^{\tau_i}(\mathbf{x}, \mathbf{y}) || \pi_{\theta}(\mathbf{x}, \mathbf{y}))$$

and keep increasing the threshold until the reward  $\mathbb{E}_{\pi_{\theta}(\mathbf{x}, \mathbf{y})} [R(\mathbf{x}, \mathbf{y})]$  stops increasing. This yields  $\theta_{k+1}$ .