

---

# Few-steps [Offline] Reinforcement Learning with Large Language Models

---

**Antonio Mari**  
antonio.mari@epfl.ch

**Matteo Santelmo**  
matteo.santelmo@epfl.ch

**Stefano Viel**  
stefano.viel@epfl.ch

## Abstract

1 We explore the application of offline reinforcement learning (RL) algorithms in  
2 few-step environments using large language models (LLMs). We target algorithms  
3 that optimize immediate reward. We employ two environments: the game of  
4 Wordle and a simpler word replacement task. Our experiments aim to assess how  
5 RL algorithms, specifically behavioral cloning (BC), filtered behavioral cloning,  
6 and our proposed reward-weighted behavioral cloning (Weighted-BC), perform  
7 in these settings. We find that Weighted-BC generalizes both BC and filtered BC,  
8 and our empirical results align with theoretical expectations. Additionally, we  
9 investigate the impact of performance conditioning on these models.

## 10 1 Introduction

11 Reinforcement Learning serves as a powerful tool for large language models training and there has  
12 been extensive research on Reinforcement Learning for Human Feedback (RLHF) (Ouyang et al.  
13 2022), whose aim is aligning the output of a language model with human preferences. More generally,  
14 RL is useful to make LLMs accomplish tasks in a goal-directed manner.

15 However, most of the recent work on RL applied to LLMs has focused on “single-step” RL problems,  
16 where a single response is optimized with respect to some reward model. This is not surprising:  
17 it is likely much easier to evaluate improvements to algorithms for single-step text generation as  
18 compared to multi-step generation, with multi-turn dialogue requiring time-consuming studies with  
19 human participants.

20 In our work, we focus on understanding how some RL algorithms can be applied and extended  
21 in a more general multi-step settings. We are interested in few-steps environments to assess how  
22 algorithms suited for the single-step case (which optimize the immediate reward) perform on short  
23 horizon scenarios.

24 Our few-steps environment is Wordle (Wardle 2021), a game in which a player has to guess a secret  
25 five-letters word and has in total six guesses to do so. After every guess, the game provides feedback  
26 on which letters of the guess are present in the secret word and if they are correctly positioned.  
27 Moreover, we also experiment with a simple single-turn toy-environment, namely word replacement,  
28 which consists of replacing some specific words in a piece of text.

29 The reasons behind the choice of these environments: i) we have access to the exact reward, ii) we  
30 can generate a dataset with a given-policy, having access to an environment simulator. This way, we  
31 generate datasets with both expert policies and noisy policies. In the former case, we would like to  
32 clone the expert behavior, while in the latter case we are interested in learning the behavior that leads  
33 to high reward. To this end, we test which methods can well exploit good samples in low quality data.

34 Contributions:

- 35 1. We formalize the Wordle environment in 2.1.2, explaining why it differs from the common  
36 settings where RL is usually applied on LLMs.
- 37 2. We analyze an algorithm proposed by our laboratory supervisors, namely the *weighted*  
38 *Behavioral Cloning*, understanding how its formulation changes under this different type  
39 of environment. We detail all the subtleties of this approach, hoping to shed lights on its  
40 potential pitfalls.
- 41 3. We show how Behavioral Cloning and Filtered Behavioral cloning represent limit cases of  
42 Weighted-BC, as its hyperparameter  $\beta$  tends to infinity or to zero, respectively.
- 43 4. We experiment these methods on different quality datasets, to see if any real benefit we can  
44 get from weighted-BC.

## 45 2 Methods

46 In this section, we formalize the environments used in our work, then present the RL algorithms used  
47 to fine-tune our LLMs. As pre-trained language models, we used Pythia-14m (Biderman et al. 2023)  
48 for the word replacement task and OPT-350M (S. Zhang et al. 2022) for Wordle.

### 49 2.1 Environments

50 We experiment with a simple single-step environment, namely "word replacement", in order to  
51 perform quick tests of our method, but the main focus is the Wordle environment. That is why we  
52 formalize Wordle in detail and we define the reward-weighted behavioral cloning considering the  
53 underlying Markov decision process. Both environments provide explicit and easily computable  
54 reward and we have access to a simulator to run online evaluation of fine-tuned LLMs.

#### 55 2.1.1 Word Replacement

56 Word replacement is a single-turn toy environment where the player’s task consists of replacing some  
57 words in an input sentence, following a replacement map learned during training. Specifically, given a  
58 vocabulary  $\mathcal{V}$  and a replacement map  $d : \mathcal{K} \rightarrow \mathcal{T}$  where  $\mathcal{K}, \mathcal{T} \subseteq \mathcal{V}$ , the model task is simply taking an  
59 input sentence and outputting the same sentence where all occurrences of a word  $k \in \mathcal{K}$  are replaced  
60 with the corresponding target word  $t = d(k)$ . A reward between between 0 and 8 was assigned to a  
61 response proportionally to the number of words correctly replaced, and set to -8 if the model was  
62 replacing a word that shouldn’t have been replaced, in order to discourage illegal substitutions.

#### 63 2.1.2 Wordle

64 In this section, we formalize the Markov Decision Process (MDP) of the Wordle game. We denote  
65 as  $\Sigma$  the lower-case English letters alphabet, (i.e.,  $\Sigma = \{a, b, \dots, z\}$ ) and we choose a vocabulary  
66  $V \subset \Sigma^5$ , where, for a set  $A$ ,  $A^n = A \times A \times \dots \times A$  indicates the Cartesian product of  $A$  with himself  
67  $n$ -times. Note that  $V$  contains only English five-letter words, and in the original game  $|V| = 2309$   
68 (while for our simplified experiments we choose a smaller dictionary,  $|V| = 100$ ). In every game,  
69 a random "secret" word is chosen uniformly form the vocabulary, i.e.  $w \sim \mathcal{U}(V)$  and represents  
70 the word to guess. We also define the feedback set  $\mathcal{F} = \{\langle - \rangle, \langle y \rangle, \langle g \rangle\}^5$ , where a  $\langle - \rangle$  represents a  
71 letter that is not present in the secret word  $w$ ,  $\langle y \rangle$  represents a letter that is present in  $w$  but in another  
72 position,  $\langle g \rangle$  indicates that the letter is correctly positioned as in  $w$ . So the feedback is deterministic  
73 when  $w$  is fixed. In general,  $w$  is not known so the feedback will be a random function.

74 The Wordle game MDP is a tuple  $(\mathcal{S}, \mathcal{A}, r, P, \mu, \gamma)$ , where

- 75 1. A state  $s = (g_1, f_w(g_1), g_2, f_w(g_2), \dots, g_6, f_w(g_6)) \in (\Sigma^5 \times \mathcal{F})^6 =: \mathcal{S}$  represents a valid  
76 board.
  - 77 •  $g_i \in V \cup \{\perp\}, \forall i \in \{1, 2, \dots, 6\}$ , i.e. each guess on the board is a five-letters English  
78 word or is empty. The latter case is represented by the special symbol  $\perp$ .
  - 79 •  $f_w(g_i) \in \mathcal{F} \cup \{\perp\}$ , i.e. each feedback follows a guess and depends on the secret  
80 word, which is not observed. Note than  $f_w(g) = \perp \iff g = \perp, \forall w \in V$

81 Note that not every board states is valid. Informally, after an empty guess  $\perp$  there could not  
 82 be any non-empty guess on the board and the feedback for each guess has to adhere with the  
 83 secret word  $w$ . With that said, we denote the set of all the valid boards given the secret word  
 84  $w$  as  $\mathcal{D}(w)$ .

85 2. An *action*  $a = g \in V$  is a valid guess, so the action space  $\mathcal{A} = V$  is discrete and any guess  
 86 can be performed at any state.

87 3. The *reward*  $r_w(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  depends on the secret word  $w$  and is defined as  $r_w(s, a) =$   
 88  $h(f_w(a))$ . In our case,  $h = 0 \cdot \#(\langle g \rangle) - 3 \cdot \#(\langle y \rangle) - 5 \cdot \#(\langle - \rangle)$  counts the occurrences  
 89 of each symbol in the feedback and assign them a score (0 for  $\langle g \rangle$ , -3 for  $\langle y \rangle$ , -5 for  $\langle - \rangle$ ).  
 90 This way, the range of the reward is  $[-25, 0]$ , with the lowest reward achievable when the  
 91 feedback is  $(\langle - \rangle, \langle - \rangle, \langle - \rangle, \langle - \rangle, \langle - \rangle)$  and the highest corresponding to the correct guess  
 92  $(\langle g \rangle, \langle g \rangle, \langle g \rangle, \langle g \rangle, \langle g \rangle)$ .

93 4. The transition probabilities  $P(s'|s, a)$  are characterized as follows:

$$P(s'|s, a) = \sum_{w \in V} p(w, s'|s, a) = \sum_{w \in V} p(w)p(s'|s, a, w) = \frac{1}{|V|} \sum_{w \in V} p(s'|s, a, w)$$

94 Note that  $p(s'|s, a, w)$  is actually deterministic because an action uniquely determines the  
 95 next guess present on the board and the feedback is a deterministic function of  $w$ . Moreover,  
 96 if we define  $s$  and  $s|(a, w)$  to be of the form

$$s = (g_1, f_w(g_1), \dots, g_i, f_w(g_i), \perp, \perp, \dots)$$

$$s|(a, w) := (g_1, f_w(g_1), \dots, g_i, f_w(g_i), a, f_w(a), \dots)$$

97 when the action  $a$  and the secret word  $w$  are fixed, then for all  $s \in \mathcal{S}$ :

$$P(s'|s, a, w) = \begin{cases} 1 & \text{if } s' = s|(a, w), \\ 0 & \text{otherwise} \end{cases}$$

98 Actually, the transition probability  $P(s'|s, a)$  is a mixture of deterministic discrete distribu-  
 99 tions where all components are weighted uniformly and the stochasticity of the transition  
 100 is directly linked to the stochasticity of the feedback, which depends on the latent variable  
 101  $w$ . Also note that the transition probabilities are actually independent of the state  $s$  (that is,  
 102 all the previous guesses and feedback received) and really depend only on the action and  
 103 the secret word. With that said, we have to include all this information because the policy  
 104 learned to solve the game *has to* use it to make new guesses.

105 5. The *initial distribution*  $\mu$  is deterministic since the starting state for every game is  
 106  $(\perp, \perp, \dots)$ .

107 6. Since we are dealing with a finite horizon environment, the discount factor  $\gamma = 1$  (note: we  
 108 will change this setting in section 2.4)

109 In our work, the policy  $\pi(\cdot|s)$  is parameterized by a language model and we encode a state  $s$  in a  
 110 prompt-string. Given a state  $s$ :

$$s = ((a, r, i, s, e), (\langle - \rangle, \langle y \rangle, \langle - \rangle, \langle y \rangle, \langle y \rangle),$$

$$(r, o, u, t, e), (\langle g \rangle, \langle - \rangle, \langle y \rangle, \langle - \rangle, \langle y \rangle),$$

$$(r, u, l, e, s), (\langle g \rangle, \langle y \rangle, \langle - \rangle, \langle y \rangle, \langle g \rangle),$$

$$\perp, \perp, \perp, \perp, \perp, \perp)$$

111 Its string encoding is reported in figure 1.

## 112 2.2 Behavioral Cloning

113 A commonly used offline RL approach for language model training is *behavioral cloning* (BC),  
 114 which consists of performing supervised fine-tuning on a dataset  $\mathcal{D}$  composed of  $(s, a)$  pairs where  
 115  $a \sim \pi_E(s)$  and  $\pi_E$  is an expert policy which we seek to clone. The maximum likelihood objective:

$$\pi_{\text{MLE}} = \operatorname{argmax}_{\pi} \sum_{(s,a) \in \mathcal{D}} \log \pi(a|s)$$

116 A common variation of BC is the *filtered behavioral cloning*, which first filters the data following  
 117 some criteria, in order to learn only from the most successful samples, rejecting the others. In our  
 118 case, a natural criterion for filtering is the reward  $r$ .

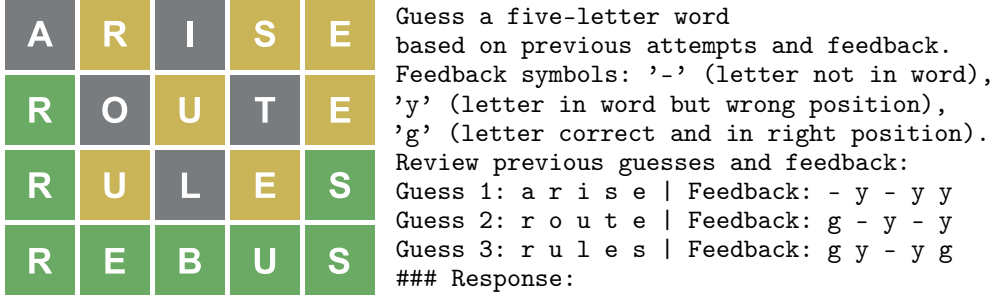


Figure 1: An example of a prompt for a Wordle game guess. Note that letters are separated by white spaces to make sure that they are all tokenized separately. The expected correct guess is "r e b u s".

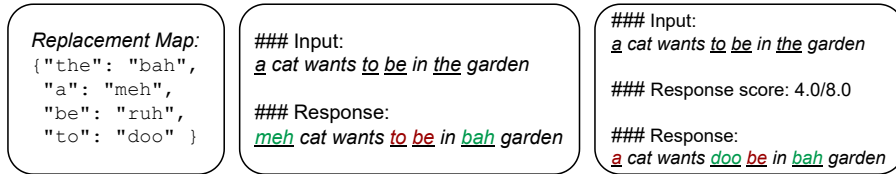


Figure 2: Example of data sample for word replacement task, generated using the replacement map on the left and probability of replacement  $p = 0.5$ , therefore providing a moderately noisy signal to the model. In the center and on the right respectively, two samples without and with reward conditioning.

119 **2.3 Performance conditioning Behavioral Cloning**

120 A surprisingly effective method to leverage the generalization and in-context learning abilities of  
121 LLMs is *performance conditioned behavioral cloning* (Shypula et al. 2023). Inspired by prompting  
122 strategies (T. Zhang et al. 2023) and offline RL (L. Chen et al. 2021), this method adapts BC for  
123 LLMs by directly introducing a reward signal into the training process by explicitly incorporating  
124 this information in the prompt. In practice, this means modifying each  $(s, a, r)$  triple to a new one  
125  $(\tilde{s}, a, r)$ , where  $\tilde{s}$  encodes both information about the current state and the reward  $r = r(s, a)$ . This  
126 allows the model to implicitly learn a mapping between responses and their correctness, which can  
127 be used during inference to improve results by asking the model to provide a response with the best  
128 reward.

129 **2.4 Reward-Weighted Behavioral Cloning**

130 Our main contributions revolve around the *Reward-Weighted Behavioral Cloning (Weighted-BC)*.  
131 In this section, we introduce its formulation and analyze the approximation involved, identifying  
132 where potential pitfalls may reside. Moreover, we show how this approach generalizes both BC and  
133 Filtered-BC.

134 **Reward optimization as a limit** Given the Wordle environment introduced in section 2.1.2 and for  
135  $\gamma \in (0, 1)$ , the value function of a policy  $\pi$  is defined as follows:

$$V^\pi(s) := \mathbb{E} \left[ \sum_{t=1}^6 \gamma^{t-1} r_w(s_t, a_t) | s_1 = s, \pi \right]$$

136 As we take the limit as  $\gamma \rightarrow 0$ , the horizon degenerates and only the immediate reward is considered:

$$\lim_{\gamma \rightarrow 0} V^\pi(s) = \mathbb{E}_{w \sim \mathcal{U}(W), s \sim \mathcal{D}(w), a \sim \pi(\cdot|s)} [r_w(s, a)]$$

137 Note that for a single-step environments and  $\gamma = 1$  this limit matches  $V^\pi(s)$ . Conversely, in a general  
138 multi-step environment, if we approximate the value function with this limit then the quality of our  
139 approximation is going to be worse as the number of steps grows. Our working assumption: in  
140 few-steps environment (such as Wordle, where the number of steps  $T = 6$ ), we assume that such  
141 approximation is "good enough", in the sense that by optimizing the immediate reward we could  
142 good results as if we optimized the quality/value function.

143 **Starting point: DPO loss** The loss commonly used for RLHF is the following (cfr DPO paper  
 144 Rafailov et al. 2023, we adapted the expected value on Wordle data):

$$\max_{\theta} \mathbb{E}_{w \sim \mathcal{U}(V), s \sim \mathcal{D}(w), a \sim \pi_{\theta}(\cdot|s)} [r_w(s, a)] - \beta D_{KL}[\pi_{\theta}(\cdot|s) \| p(\cdot|s)]$$

145 where  $p$  is a reference policy (a pre-trained LLM or an handcrafted policy), defined over states  $s$ , and  
 146  $\beta \geq 0$  is a regularization parameter that modulates how far we may stray from  $p(\cdot|s)$ .

147 The expected value is taken with respect to the space of the observable data, which represents all the  
 148 possible wordle games. A sample from this dataset space is drawn first by sampling the secret word  
 149 uniformly from the vocabulary  $w \sim \mathcal{U}(V)$  and then sampling a "random" state  $s \sim \mathcal{D}(w)$ , where  
 150  $\mathcal{D}(w)$  is the set of legal board states which are consistent with the secret word  $w$ . We can assume  
 151 that a sample is draw uniformly random from this set, but this assumption is actually violated in the  
 152 datasets that we are going to use, where full trajectories are drawn (see section 3.2).

153 As shown in Rafailov et al. 2023 (and here adapted due to the introduction of the expected value on  
 154  $w$ ), the analytical solution to this objective is:

$$\pi^*(a|s, w) = \frac{1}{Z(s, w)} p(a|s) \exp\left(\frac{1}{\beta} r_w(s, a)\right)$$

155 where  $Z(s, w)$  is the (intractable) partition function that ensures that  $\sum_{a \in \mathcal{A}} \pi^*(a|s, w) = 1$ .

156 *Nota Bene:*  $\pi^*$  is conditioned on both  $s$  and  $w$ , but the policy that we are interested to clone should  
 157 not be conditioned on  $w$ . A natural way to obtain it is marginalizing out  $w$  (overloading  $\pi^*$  symbol):

$$\pi^*(a|s) = \sum_{w \in V} p(w|s) \cdot \pi^*(a|s, w)$$

158 Another *crucial approximation* that we do next is assuming  $p(w|s) = p(w)$ , i.e.  $w$  is independent of  
 159  $s$ , so that:

$$\mathbb{E}_{w \sim \mathcal{U}(V)} [\pi^*(a|s, w)] = \sum_{w \in V} p(w) \pi^*(a|s, w) = \pi^*(a|s)$$

160 **Weighted-BC Objective** *Proposed idea*<sup>1</sup>: *optimize a cross-entropy objective whose solution is*  
 161  *$\pi^*(\cdot|s)$  (under the approximation as before)*

$$\begin{aligned} \mathcal{L}(\theta) &= \mathbb{E}_{w \sim \mathcal{U}(V), s \sim \mathcal{D}(w)} \left[ \sum_{a \in \mathcal{A}} \pi^*(a|s, w) \log \pi_{\theta}(a|s) \right] \\ &= \mathbb{E}_{w \sim \mathcal{U}(V), s \sim \mathcal{D}(w)} \left[ \frac{\sum_{a \in \mathcal{A}} p(a|s) \exp\left(\frac{1}{\beta} r_w(s, a)\right) \log \pi_{\theta}(a|s)}{\sum_{a \in \mathcal{A}} p(a|s) \exp\left(\frac{1}{\beta} r_w(s, a)\right)} \right] \\ &= \mathbb{E}_{w \sim \mathcal{U}(V), s \sim \mathcal{D}(w), a \sim p(\cdot|s)} \left[ \frac{1}{Z(s, w)} \exp\left(\frac{1}{\beta} r_w(s, a)\right) \log \pi_{\theta}(a|s) \right] \end{aligned}$$

162 Since the explicit value of  $p(\cdot|s)$  is not directly accessible, for each state  $s$  we sample  $n$  times from  
 163 this distribution and denote the obtained set of samples as  $A_s$ , with  $|A_s| = n$ . Consequently, we  
 164 set  $\hat{p}(\cdot|s) := \mathcal{U}(A_s)$ , so we can approximate both the expected value over  $p(\cdot|s)$  and  $Z(s, w)$  using  
 165  $\hat{p}(\cdot|s)$  instead. Our sampled loss is denoted as  $\hat{\mathcal{L}}$ :

$$\mathcal{L}(\theta) \approx \hat{\mathcal{L}}(\theta) := \mathbb{E}_{w \sim \mathcal{U}(V), s \sim \mathcal{D}(w), a \sim \hat{p}(\cdot|s)} \left[ \frac{\exp\left(\frac{1}{\beta} r_w(s, a)\right)}{\sum_{a \in A_s} \exp\left(\frac{1}{\beta} r_w(s, a)\right)} \log \pi_{\theta}(a|s) \right]$$

166 Some considerations:

167 • The objective is a likelihood where actions are weighted using their reward. The population  
 168 population weights are, for all  $w \in V$ ,  $s \in \mathcal{D}(w)$  and  $a \in \mathcal{A}$ :

$$\rho(a|s, w) = \frac{\exp\left(\frac{1}{\beta} r_w(s, a)\right)}{\sum_{a \in \mathcal{A}} \exp\left(\frac{1}{\beta} r_w(s, a)\right)}$$

<sup>1</sup>Explain exactly who proposed it and how

169 • These weights are approximated using sets of samples  $A_s$ :

$$\hat{\rho}(a|s, w) = \frac{\exp\left(\frac{1}{\beta}r_w(s, a)\right)}{\sum_{a \in A_s} \exp\left(\frac{1}{\beta}r_w(s, a)\right)}$$

170 • The sampled loss  $\hat{\mathcal{L}}(\theta)$  is a biased estimator of  $\mathcal{L}(\theta)$  because the expected value does not  
171 distribute over products and divisions (unless independence holds, which is not the case).

172 **Weighted-BC generalizes BC and Filtered-BC** Let us analyze the limit behavior of the weighted-  
173 BC. Let  $n = |A_s|$  for every  $s \in \mathcal{D}$  (where  $\mathcal{D}$  is a generic dataset) and let  $b$  be the batch size used for  
174 our updates.

175 1. When  $\beta \rightarrow 0$ , the weights  $\rho(a|s, w) \rightarrow 1$  on the action with maximum reward (assuming it  
176 is unique). It corresponds to Filtered-BC, where the batch size is  $b/n$ .

177 2. When  $\beta \rightarrow \infty$ ,  $\rho(a|s, w) = 1/n, \forall a \in A_s$ . It corresponds to BC, since the constant  $1/n$   
178 does not affect the solution of the optimization problem.

### 179 3 Experiments

180 In this section, we evaluate the mentioned methods on both environments. In our experiments, the  
181 main research questions are:

182 Q1. Does weighted-BC empirically match our theoretical expectation? How the hyper-parameter  
183  $\beta$  affects the performance? Can weighted-BC outperform the baselines?

184 Q2. Does the performance-conditioning help in solving the tasks?

185 Q3. What changes if we have a dirty dataset, which comes from a policy with sub-optimal  
186 behavior?

#### 187 3.1 Word Replacement

188 Given the simplicity of this task, the experiments were conducted on a less powerful language model  
189 and using a quite noisy signal. In particular, the performance of different models discussed in section  
190 2 were compared using a dataset generated with a probability of replacement of the key-words of  
191  $p = 0.25$ . The results, in figure 3, show that despite having a relatively noisy signal, filtered BC  
192 proved to always achieve better rewards, both in the performance conditioned and unconditioned  
193 cases. A possible explanation behind this unexpected behavior stands in the way the reward is defined,  
194 in particular because in all cases where there is no word to be replaced in the input, the non-filtered  
195 methods will see more samples being identical while not being useful for the model to learn the task,  
196 despite their high reward. 3.

#### 197 3.2 Wordle

198 We test our wordle environment using two datasets, the first one generated using an handcrafted  
199 expert policy  $\pi_E$  and the second generated using a noisy policy, which consists of a mixture between  
200  $\pi_E$  and a random policy.

201 The datasets are generated in the following way:

202 1. A secret word  $w$  is sampled from  $\mathcal{U}(V)$ . Then, a Wordle game with that secret word is  
203 played by drawing samples from  $p(\cdot|s)$ . At each turn of a game,  $N$  samples are generated  
204 and they are gathered in the set  $A_s$ , along with their rewards. Since multiple actions are  
205 generated each turn, a random one is then kept for next turn and the state  $s$  is updated.

206 2. This process is repeated  $M$  times, so that at the end we will have collected trajectories that  
207 correspond to  $M$  Wordle games.

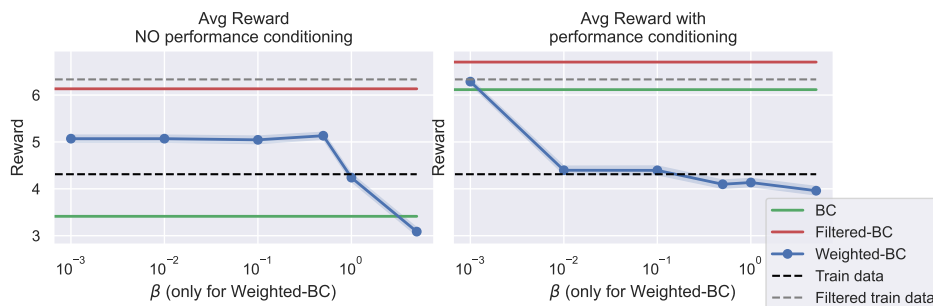


Figure 3: Average reward (with 0.95 confidence interval) for word replacement single-turn task. Weighted BC was tested over different  $\beta$ s and both BC and Weighted BC were trained with suboptimal supervisions (dashed line), while Filtered BC’s reference dataset had more quality samples. Weighted BC’s performance are between BC and Filtered BC, as theoretically expected, approaching Filtered BC and standard respectively with small  $\beta$  and high  $\beta$ .

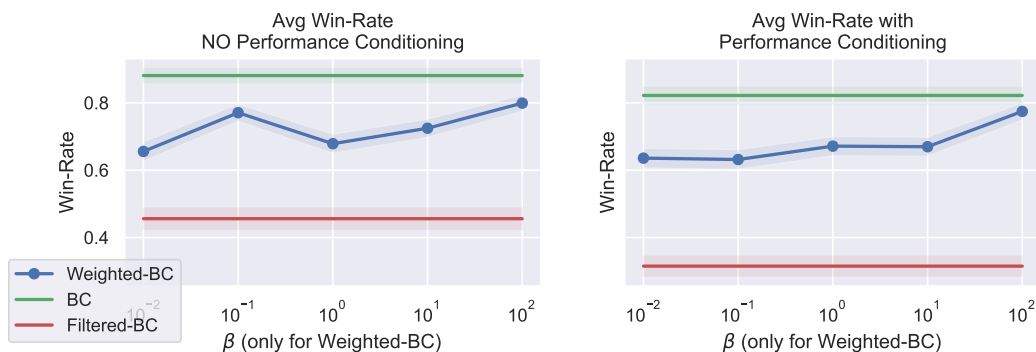


Figure 4: Win-rate (over 1000 games) of models trained on the **expert policy dataset**. **5000 games were used for training**, where at every step 5 independent guesses are sampled from the expert policy. Then, **the guess associated with the highest reward is kept for Filtered-BC while a random guess is used to continue the game**. The batch size is 128 for BC and Weighted-BC, 25 for Filtered-BC, and the learning rate is consistently  $5e-05$ . The same number of gradient updates was performed across models, ensuring the convergence of the training.

### 208 3.2.1 Expert policy dataset

209 The expert policy is handcrafted and at each step filters all the words in the vocabulary based on the  
 210 known constraints and randomly samples among the filtered words. The win-rate of this policy over  
 211 5000 games is 99.84%.

212 The results are showed in fig. 4, and the answers to our research questions:

213 A1. The empirical results match the expectation. Discarding data is harmful because all guesses  
 214 are generated by the expert policy  $\pi_E$  and, despite some fluctuations, by modulating  $\beta$  it  
 215 appears that we approach both limits (i.e., Filtered-BC for small  $\beta$ , BC for high  $\beta$ ), closing  
 216 the gap between the two methods.

217 A2. Performance-conditioning leads to worse win-rates in both baselines, while it does not  
 218 change much for weighted-BC. At this stage, it is not really clear if these results are  
 219 statistically significant, so that runs with multiple seeds would help to assess whether there  
 220 is a real difference between the two plots of fig. 4.

### 221 3.2.2 Mixture of Expert and Random Policy dataset

222 The "mixture policy" is a mixture between the expert policy and a random policy, with the mixture  
 223 weight  $p = 0.4$ . This means that, with 40% probability, a sample is drawn from a random policy (so

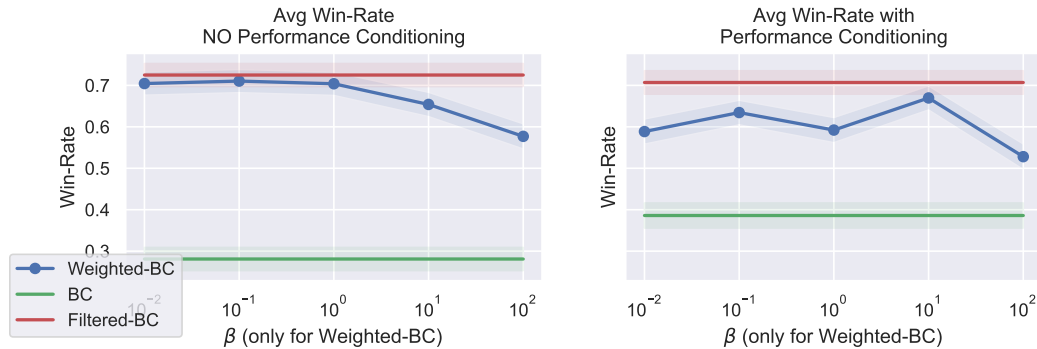


Figure 5: Win-rate of models trained on the **mixture policy dataset**. Same training settings as in fig. 4

224 from a uniform distribution over the whole dictionary) and with 60% probability it is drawn from  $\pi_E$ .  
 225 The win-rate of the mixture policy over 5000 games is 92.4%. This good performance is explained by  
 226 the fact that sampling random words helps exploration of new constraints, so subsequent non-random  
 227 guesses will still leverage the constraints discovered. Results on this dataset are reported in fig. 5

228 A1. Also in this case, empirical results match the expectation. Opposite to the previous case,  
 229 Filtered-BC dominates by only retaining the best examples (1/5 of them) and we can clearly  
 230 observe weighted BC approaching filtered-BC for small  $\beta$ .

231 A2. Similar considerations than before can be drawn. The only difference is that we observe a  
 232 boost in BC win-rate, but we are far from being able to reach any conclusion on performance  
 233 conditioning.

234 A3. Generally, in the expert policy dataset we can reach higher win-rates (> 80%) and, according  
 235 to our intuition, leveraging all the available high-quality data, as well as filtering only the  
 236 good-data in a dirty dataset, is the strategy that achieves the highest win-rate.

## 237 4 Discussion

238 As claimed in DPO paper, the hyper-parameter  $\beta$  regulates the "strength" of the KL regularizer term.  
 239 This intuition is consistent with our experiments on weighted-BC. In particular, high  $\beta$  means strong  
 240 adherence to the reference policy. This will lead to better results if the data was generated using an  
 241 expert policy. Conversely, for small  $\beta \approx 0$ , we optimize the reward without adhering to the reference  
 242 policy, thing that is desirable when the data is dirty or noisy.

243 In the high quality data scenario, as expected, the best results are obtained by the methods that are  
 244 leveraging more data, i.e. BC and Weighted BC with high  $\beta$ , while for poor quality data the best  
 245 results are achieved by methods which leverage samples with higher rewards, either by filtering or by  
 246 using a very accentuated weighting.

247 It is worth noting that if we do not have prior information about the cleanliness of our data then  
 248 weighted-BC seems to be a proper choice to achieve reasonable performance in different scenarios.

249 Future works should firstly address a more exhaustive experimentation to confirm the validity of  
 250 our results and possibly compare with more complex algorithms capable of optimizing over the  
 251 cumulative reward or an approximation of it, going over the limitation of optimizing on the immediate  
 252 reward.

## 253 References

- 254 Biderman, Stella et al. (2023). "Pythia: A suite for analyzing large language models across training  
 255 and scaling". In: *International Conference on Machine Learning*. PMLR, pp. 2397–2430.  
 256 Chen, Lili et al. (2021). "Decision transformer: Reinforcement learning via sequence modeling". In:  
 257 *Advances in neural information processing systems* 34, pp. 15084–15097.



- 258 Ouyang, Long et al. (2022). “Training language models to follow instructions with human feedback”.  
259 In: *Advances in neural information processing systems* 35, pp. 27730–27744.
- 260 Rafailov, Rafael et al. (2023). *Direct Preference Optimization: Your Language Model is Secretly a*  
261 *Reward Model*. arXiv: 2305.18290 [cs.LG].
- 262 Shypula, Alexander et al. (2023). “Learning performance-improving code edits”. In: *arXiv preprint*  
263 *arXiv:2302.07867*.
- 264 Wardle, Josh (2021). *Wordle*. <https://www.nytimes.com/games/wordle/index.html>.
- 265 Zhang, Susan et al. (2022). “Opt: Open pre-trained transformer language models”. In: *arXiv preprint*  
266 *arXiv:2205.01068*.
- 267 Zhang, Tianjun et al. (2023). “The wisdom of hindsight makes language models better instruction  
268 followers”. In: *International Conference on Machine Learning*. PMLR, pp. 41414–41428.