**Algorithm 1** NLPO- Natural Language Optimization

1: **Input:** Dataset $\mathcal{D} = \{(\mathbf{x}^i, \mathbf{y}^i)\}_{i=1}^N$ of size $N$
2: **Input:** initial policy parameter $\pi_{\theta_0}$
3: **Input:** initial LM $\pi_0$
4: **Input:** initial value function parameters $V_{\phi_0}$
5: **Input:** initialize parameterized masked policy $\pi_{\psi_0}(\cdot|\cdot, \pi_{\theta_0})$ with parameterized top-p policy $\pi_{\theta_0}$
6: **Input:** policy update frequency $\mu$
7: **repeat**
8:    Sample mini-batch $\mathcal{D}_m = \{(\mathbf{x}^m, \mathbf{y}^m)\}_{m=1}^M$ from $\mathcal{D}$
9:    Collect trajectories $\mathcal{T}_{\tau_i}$ by running policy $\pi_{\psi_n}$ in for batch $\mathcal{D}_m$ in env
10:    Compute Preference and KL penalty rewards $\hat{R}_t$
11:    Compute the advantage estimate $\hat{A}_t$
12:    Update the policy by maximizing the PPO-Clip objective:
13:

$$\pi_{\theta_{m+1}} = \text{argmax}_\theta \frac{1}{\mathcal{D}_m T} \sum_{\tau \in \mathcal{D}_m} \sum_{\tau=0}^T \min(r_t(\theta) A^{\pi_{\theta_m}}, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) A^{\pi_{\theta_m}})$$

14:
15:    where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_m}(a_t|s_t)}$.
16:
17:    Update the value function:
18:

$$V_{\phi_{m+1}} = \text{argmin}_\phi \frac{1}{\mathcal{D}_m T} \sum_{\tau \in \mathcal{D}_m} \sum_{t=0}^T \left(V_\phi(s_t) - \hat{R}_t\right)^2$$

19:    Update the parameterized masked poicy every $\mu$ iterations:
20:

$$pi_{\psi_{n+1}}(\cdot|\cdot, \pi_{\theta_{m+1}})$$

21: **until** convergence and **return** $\pi_\theta$