# Initialization

**Initialization**

Supervised
Dataset $\mathcal{D}_0$

Supervised
Training

GPT-2
Weights

# Learning from User Behavior

**Learning from User Behavior**

Rounds $r = 1, 2, 3, \dots$

User Interactions



Dataset $\mathcal{D}_r$
Construction

Natural Language
Generation Model
$\bar{x} \sim P(\cdot \mid \cdot, \cdot ; \theta_r)$

Datasets
$\mathcal{D}_0, \dots, \mathcal{D}_r$

Contextual Bandit Training