

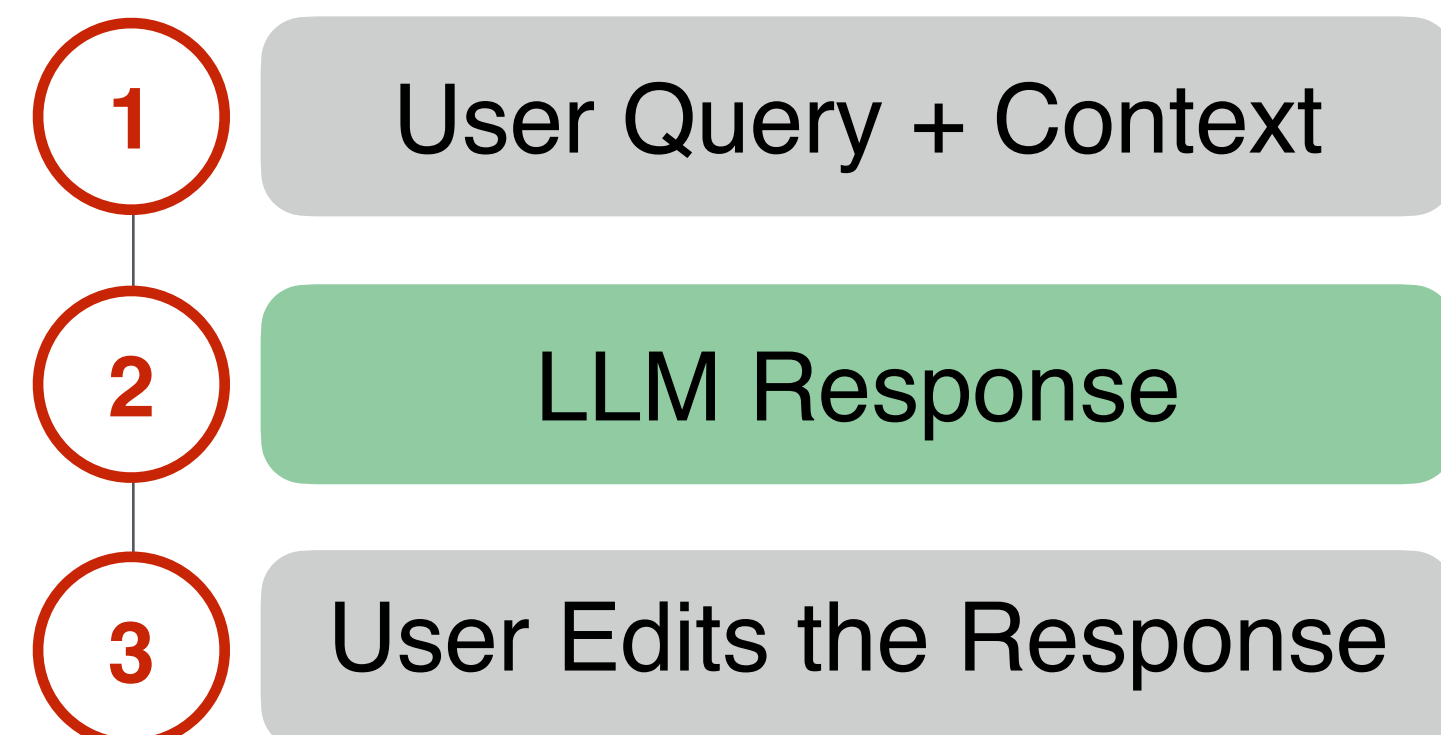
Principled Fine-tuning of LLMs from User-Edits: A Medley of Preference, Supervision, and Reward

Dipendra Misra*, Aldo Pacchiano*, Ta-Chung Chi, and Ge Gao



Post-Training LLMs with User-Edits Deployment Data

- Post-training LLMs on an in-distribution labeled data can lead to improvement in the performance.
- Challenge:** where to get the training data for post-training?
 - Get annotations from third-party. **Cons:** expensive, not in-distribution, need to ensure diversity.
 - What if we could use deployment data? **Pros:** in-distribution, abundant.
- We focus on **user-edits** — a common type of deployment data:
 - Common in coding and writing AI assistants.
 - Provides naturally occurring user feedback for improving LLMs!

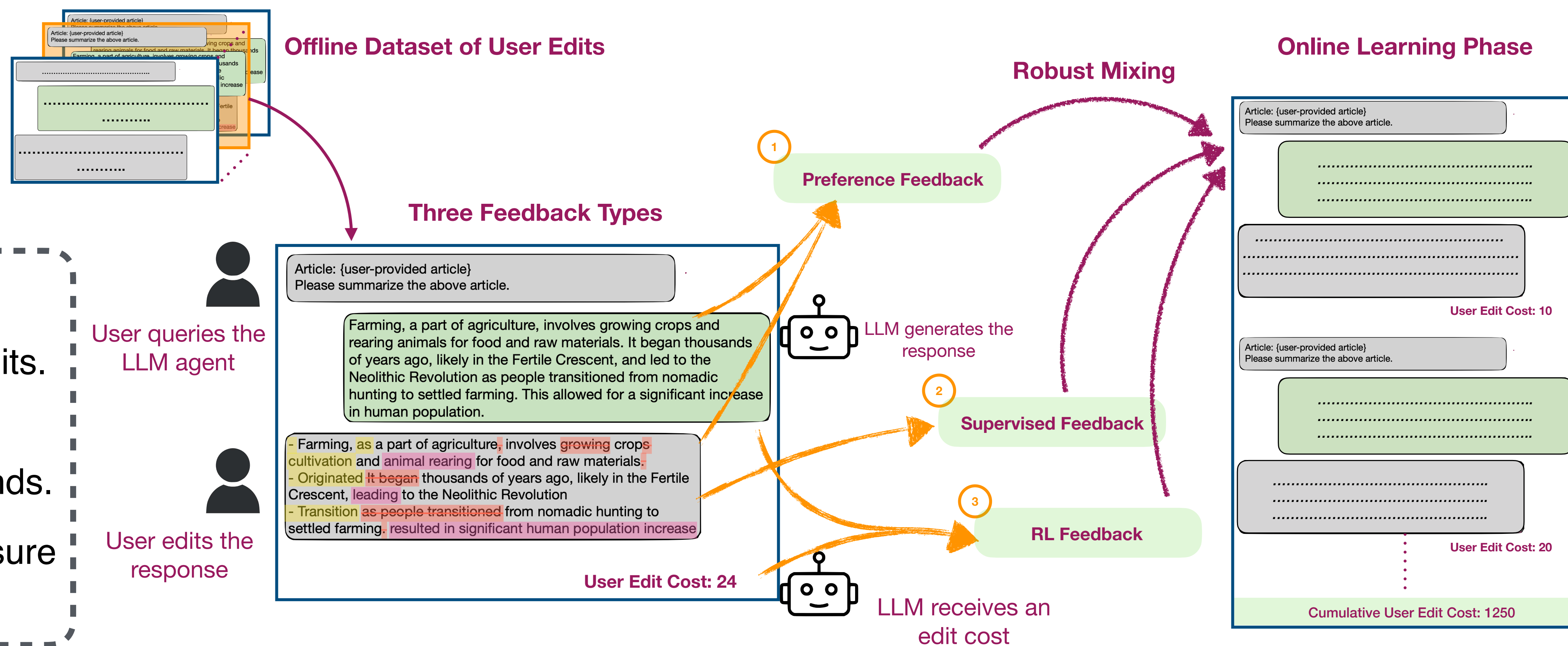


Principled Learning from User-Edits

- Initiate principled algorithm development in learning from user-edits setting.
- Rich algorithm design landscape: three different learning feedback type available!

Learning Setup

- Offline Learning:** Learn from a dataset of generations and user-edits.
- Online Learning:** Evaluate with a user over a small number of T rounds.
- Edit Distance Cost:** Used to measure performance (lower \rightarrow better).



Theoretical Results & Main Algorithm

Theoretical Results: Derive bounds on sub-optimality (SubOpt) of algorithms that use the three individual feedback types:

- 1 SFT (labeled) 2 DPO (preference data) 3 Offline RL (cost)

Informal Theorem: Sub-optimality of DPO policy $\hat{\pi}_{\text{DPO}}$ on user-edits preference data satisfies:

$$\text{SubOpt}(\hat{\pi}_{\text{DPO}}) \leq \epsilon, \quad \text{if } n \geq \Omega\left(\frac{C_{\text{PREF}} \log \frac{|\Pi|}{\delta}}{\beta^2 \sigma'(-V_{\max})^2 \epsilon^2}\right)$$

- DPO is more affected by coverage while SFT by quality of feedback.

Main Algorithm: Combines different user feedback types in two ways:

- Early Ensemble: train a policy to optimize the loss of different types jointly.
- Late Ensemble: run a bandit algorithm in the online phase on the set of policies trained with various feedback types.

Results

- Evaluate on two domains from Gao et al. 2024, using simulated LLM user.
- A weak user converges more slowly π^* to than a strong user.

Method	Summarization		Email Writing		Max SubOpt
	Strong User	Weak User	Strong User	Weak User	
Base	0.9455 \pm 0.01	0.9445 \pm 0.02	0.5108 \pm 0.03	0.4923 \pm 0.01	0.7364 \pm 0.10
SFT	0.5377 \pm 0.02	0.9304 \pm 0.19	0.4159 \pm 0.05	0.4539 \pm 0.03	0.5772 \pm 0.19
DPO	1.0790 \pm 0.06	0.8267 \pm 0.06	0.3365 \pm 0.00	0.3368 \pm 0.01	0.8698 \pm 0.11
EarlyEnsemble	0.2092 \pm 0.09	0.3586 \pm 0.01	0.3438 \pm 0.06	0.4864 \pm 0.01	0.1612 \pm 0.01
LateEnsemble	0.2768 \pm 0.13	0.4403 \pm 0.03	0.4202 \pm 0.11	0.3739 \pm 0.04	0.1586 \pm 0.04