

### FedRLHF: A Convergence-Guaranteed Federated Framework for Privacy-Preserving and Personalized RLHF

Flint Xiaofeng Fan, Cheston Tan, Yew-Soon Ong, Roger Wattenhofer, and Wei Tsang Ooi

AAMAS 2025

Flint Xiaofeng Fan Postdoc, ETH Zurich Scientist, A\*STAR CFAR A\*STAR International Fellow



• RL from Human Feedback

-> to align with human preferencc

- Key Applications
  - Robotics
  - Recommendation
  - Language Models

-> e.g, ChatGPT



• RL from Human Feedback

-> to align with human preferencc

- Key Applications
  - Robotics
  - Recommendation
  - Language Models

-> e.g, ChatGPT

### **Policy Model**





• RL from Human Feedback

-> to align with human preferencc

- Key Applications
  - Robotics
  - Recommendation
  - Language Models

-> e.g, ChatGPT





#### You're giving feedback on a new version of ChatGPT.

Which response do you prefer? Responses may take a moment to load.





• RL from Human Feedback

-> to align with human preferencc

- Key Applications
  - Robotics
  - Recommendation
  - Language Models

-> e.g, ChatGPT



## **Challenge: Centralization**

• Privacy

7

- All user data and human preference are collected (HIPAA, GDPR, PDPA, etc.)
- Personalization





## **Challenge: Centralization**

- Privacy
  - All user data and human preference are collected (HIPAA, GDPR, PDPA, etc.)
- Personalizatieference model for all users





## **Challenge: Centralization**

• Privacy

=> All user data and human preference are collected (HIPAA, GDPR, etc.)

- Personalization
  - => One preference model for all users



Α

### **Solution: FedRLHF**

### RLHF



Centralized data/feedback collection



Single global preference model

FedRLHF

- Local training on local data
- Local feedback models
- Only model weights shared
- Personalized policies





For each hospital

- Trains its RL model locally using patient data
- Environmental rewards: Objective health metrics (e.g., blood sugar levels, recovery time)





For each hospital

- Trains its RL model locally using patient data
- Environmental rewards: Objective health metrics (e.g., blood sugar levels, recovery time)
- Doctor feedback provides qualitative signals (e.g., treatment appropriateness)

A central party (central hospital, government, etc.)

- Aggregates model updates from all hospitals
- Captures shared knowledge



(s<sub>t</sub>,a<sub>t</sub>) Too aggressive? Just right? Too conservative?

- Local personalization
- Hospitals receive the global model and
- Refine further using local data and feedback
- Adapts to local patient demographics and doctor preferences
- Balance between global collaboration and highly personalized care



### **FedRLHF – Problem Formulation**

System setup: K clients, each client  $k \in 1, 2, ..., K$ 

Local MDP

$$M_k = (\mathcal{S}, \mathcal{A}, P_k, R_k^0, 
ho_0(s), \gamma)$$

- Local feedback
- Shaped reward
- Local objective





### **FedRLHF – Problem Formulation**

System setup: K clients, each client  $k \in 1, 2, ..., K$ 

 $H_k(s,a)$ 

- Local MDP
- Local feedback
- Shaped reward
- Local objective

$$M_k = (\mathcal{S}, \mathcal{A}, P_k, R_k^0, 
ho_0(s), \gamma)$$

$$R_k(s,a) = R_k^0(s,a) + \lambda H_k(s,a),$$

- Server-clients communication FedAvg
  - Server broadcasts global model
  - Clients send parameters

**Global objective** 

$$J( heta) = rac{1}{K}\sum_{k=1}^K J_k( heta)$$



### **FedRLHF - Convergence**

Assumption 6 (Bounded Human Feedback). For all  $s \in S$ ,  $a \in \mathcal{A}$ , and  $k \in [K]$ , the human feedback is bounded:

 $|H_k(s,a)| \leq H_{\max}.$ 

REMARK. Assumption 6 limits the variance introduced by human feedback in the learning process. In our experiments with the Movie-Lens task, we implement this by bounding feedback values and options (Section 6.1.2), similar to practical systems like ChatGPT that curate feedback for consistency.

# Theorem 4.1 (Convergence of FedRLHF). The output of Algorithm FedRLHF satisfies:

$$\mathbb{E}[J(\theta^{*}) - J(\theta_{\text{avg}})] \leq \underbrace{\frac{L}{\mu T}(J(\theta^{*}) - J(\theta_{0}))}_{\text{$\downarrow$}} + \underbrace{\frac{1}{2\mu K}(G^{2} + \sigma^{2})}_{\text{$\downarrow$}} + \underbrace{\frac{L}{\mu}\lambda H_{\text{max}}}_{\text{$\downarrow$}}.$$

$$\underbrace{\bigcup_{\substack{O(1/T):\\ \text{linear convergence rate}\\ \text{$w.r.t. T}}}_{\text{$w.r.t. K}} \underbrace{\bigcup_{\substack{O(1/K):\\ \text{$scalability}\\ \text{$w.r.t. K}}}}_{\text{$w.r.t. K}} \underbrace{\bigcup_{\substack{O(1):\\ \text{$bounded impact}\\ \text{$from feedback}}}}_{\text{$from feedback}}$$



### **FedRLHF - Sample Complexity**

Theorem 4.2 (Sample Complexity of FedRLHF). To achieve an expected optimality gap of  $\mathbb{E}[J(\theta^*) - J(\theta_{avg})] \leq \epsilon$ , the total number of samples required across all clients is:

 $N = O\left(\frac{L(G^2 + \sigma^2)}{\mu^2 \epsilon^2}\right)$ 

Task-dependent properties

resulting in per-client sample complexity:

$$N_c = rac{N}{K} = O\left(rac{L(G^2 + \sigma^2)}{K\mu^2\epsilon^2}
ight)$$

combined effect of gradient bound and variance

decreases proportionally with K



## **FedRLHF - Personalization vs Performance**

A fundamental trade-off in collaborative learning:

VS

### Global Learning

- Shared knowledge
- Collaborative gain
- Model consistency

## Local Adaptation

- Personal preferences
- Client-specific behavior
- Divergent policies

Definition 5.2 (Personalization Score).

$$P_k( heta) = \mathbb{E}_{s \sim 
ho}[D_{ ext{KL}}(\pi_k(\cdot|s, heta) \parallel \pi(\cdot|s, heta))]$$

Definition 5.3 (Global Performance Metric).

$$J_g( heta) = rac{1}{K}\sum_{k=1}^K J_k^0(\pi)$$

Theorem 5.1 (Personalization-Performance Trade-off). Using the FedRLHF algorithm, the global performance metric (Definition 5.3) satisfies:

$$J_g( heta) \geq rac{1}{K}\sum_{k=1}^K J_k^0(\pi_k) - C \cdot \left(rac{1}{K}\sum_{k=1}^K \sqrt{P_k( heta)}
ight)$$

Theorem 5.2 (Impact of Human Feedback):a) The average personalization scores increases atb) The global performance decreases at $O(\lambda)$ c) The system sample complexity increases at



### **FedRLHF - Personalization vs Performance**

Takeaway:

1. Trade-off Highlight: Improving personalization comes at the cost of global performance.

2. 
$$R_k(s,a)=R_k^0(s,a)+\lambda H_k(s,a),$$

- λ † Local Adaptation † Global Consistency
- $\lambda \downarrow$  Local Adaptation  $\downarrow$  Global Consistency

### **Experiment: Sentiment-Controlled Text Generation**

### Dataset & Task:

- 50,000 labelled reviews from IMDb dataset (positive/negative sentiment)
- Partitioned among K=5 clients, each with ~10k reviews
- Objective: Fine-tune a GPT-like model to generate text that matches a *desired sentiment style*. Different clients want different positivity levels.

### Setup:

- Transformer RL (TRL)<sup>2</sup> library from hugging face
- Flower API<sup>3</sup> for simulating realistic distributed network, governed by a server





<sup>2</sup>https://github.com/huggingface/trl/tree/main <sup>3</sup>https://flower.ai/docs/framework/index.html Please refer to our paper for more details and results.



### **Implementation & Local Feedback Mechanism**

**Model** pretrained from HuggingFace

- GPT-2 (125M parameters + PPO overhead for RL) Intrinsic rewards  $R_k^0 \in [0,1]$
- Based on log-likelihood of coherent text (fluency)

Feedback/Sentiment  $R_{ ext{sentiment}} \in [0,1]$ 

• Outputs from pretrained DistBERT sentiment classifier Combined reward  $R_k = \lambda_k \cdot R_{ ext{sentiment}} + (1 - \lambda_k) \cdot R_k^0$ 

### **Results – FedRLHF vs Centralized RLHF**

- FedRLHF catches up in average reward, eventually surpassing centralized RLHF after ~1500–2000 samples.
- FedRLHF starts with higher loss in early rounds, due to variance among clients, but quickly drops due to the exploration of more clients.
- FedRLHF achieves comparable, if not better, performance than centralized RLHF, while preserving privacy



### **Results – Personalization & λ Tuning**



### Conclusion

### 1. FedRLHF decentralizes RLHF in FL

- Personalized preference/reward model
- 2. Theoretical Foundations
  - Convergence guarantees
  - Sample complexity analysis
  - Personalization-performance tradeoff
- 3. Looking ahead
  - Relaxed assumptions
  - Further privacy enhancements
  - Robust aggregation
  - Principled implementation of feedback



### fxf@u.nus.edu



## **THANK YOU**

