

# Online Learning and Equilibrium Computation with Ranking Feedback

Mingyang Liu, Yongshan Chen, Zhiyuan Fan, Gabriele Farina, Asuman Ozdaglar, Kaiqing Zhang



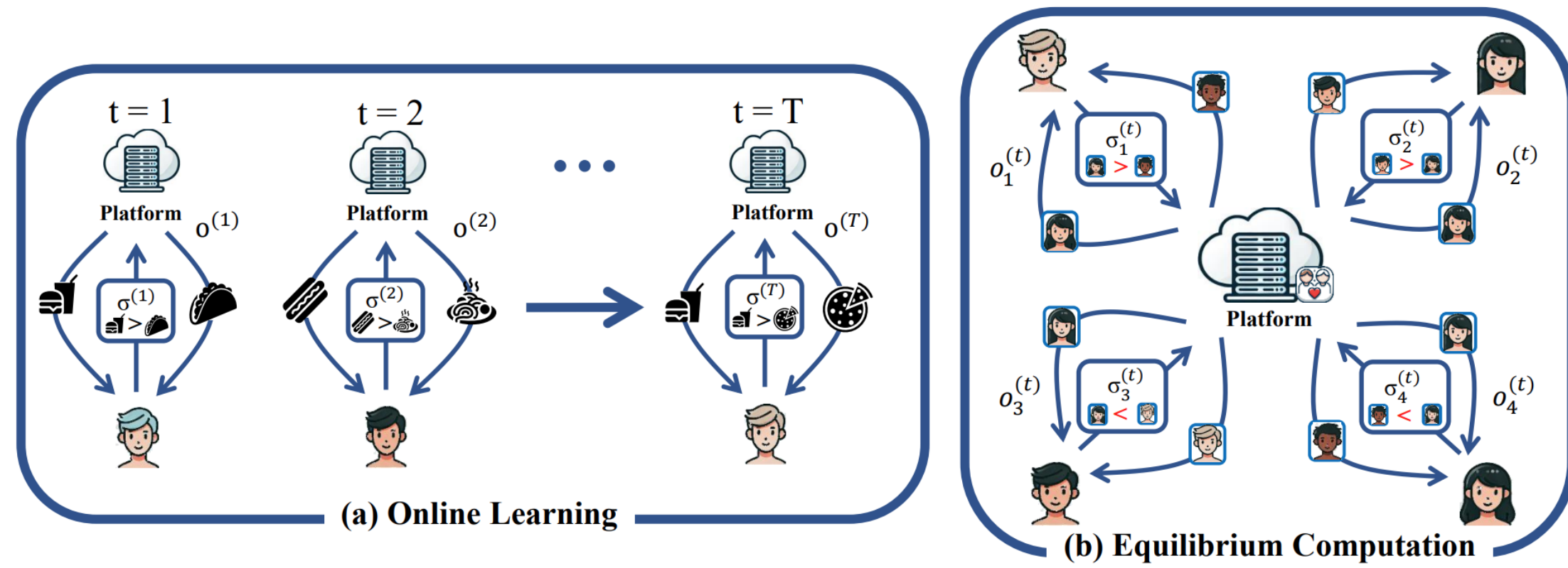
## Examples

### Numeric v.s. Ranking Feedback

- Easier for humans to rank than score,
- Numeric utilities may be private or unavailable.

### Applications

- LLM routing with user preference rankings,
- Recommendation / matching system.



## Comparison to Dueling Bandit

$$R^T = \max_{\hat{a} \in \mathcal{A}} \sum_{t=1}^T u^t(\hat{a}) - \frac{1}{2} \left( u^t(a^t) + u^t(b^t) \right). \quad (1.1)$$

### Adversarial Dueling Bandit [Saha et al., 2021]

Receive  $U^{(1)}, \dots, U^{(T)} \in [0, 1]^{\mathcal{A} \times \mathcal{A}}$

$$\Pr(a^{(t)} > b^{(t)} | o^{(t)} = \{a^{(t)}, b^{(t)}\}) = U^{(t)}(a^{(t)}, b^{(t)})$$

$$u^t(a) = \sum_{a' \neq a} U^t(a, a')$$

### Ours

Receive  $u^{(1)}, \dots, u^{(T)} \in [-1, 1]^{\mathcal{A}}$

$$\Pr(a^{(t)} > b^{(t)} | o^{(t)} = \{a, b\}) = \frac{\exp(r^t(a))}{\exp(r^t(a)) + \exp(r^t(b))}$$

$$r^t = \begin{cases} u^t & \text{Instantaneous Utility} \\ \frac{1}{t} \sum_{s=1}^t u^s & \text{Time-average Utility (Full-Info)} \\ \frac{\sum_{s=1}^t u^s(a) \sum_{a' \in o^{(s)}} \mathbb{1}(a=a')}{\sum_{s=1}^t \sum_{a' \in o^{(s)}} \mathbb{1}(a=a')} & \text{Time-average Utility (Bandit)} \end{cases}$$

## Online Learning with Ranking Feedback

### Procedure

- Update the strategy  $\pi^t \in \Delta^{\mathcal{A}}$
- Propose a multiset (possibly include repeated elements) of actions  $o^{(t)} \subseteq \mathcal{A}$
- Receive ranking  $\sigma^{(t)}$  over options in  $o^{(t)}$

Plackett-Luce (PL) model [Luce, 1959, Plackett, 1975]

$$\Pr(\sigma^{(t)} | o^{(t)}) = \prod_{k_1=1}^K \frac{\exp\left(\frac{1}{\tau} r^t(\sigma^{(t)}(k_1))\right)}{\sum_{k_2=1}^K \exp\left(\frac{1}{\tau} r^t(\sigma^{(t)}(k_2))\right)}. \quad (1.2)$$

$r^t$  is a function of  $\{u^{(1)}, \dots, u^{(t)}\}$  to be specified later.

## Ranking with Instantaneous Utility

$$r^t = u^t.$$

### Full-Information Feedback

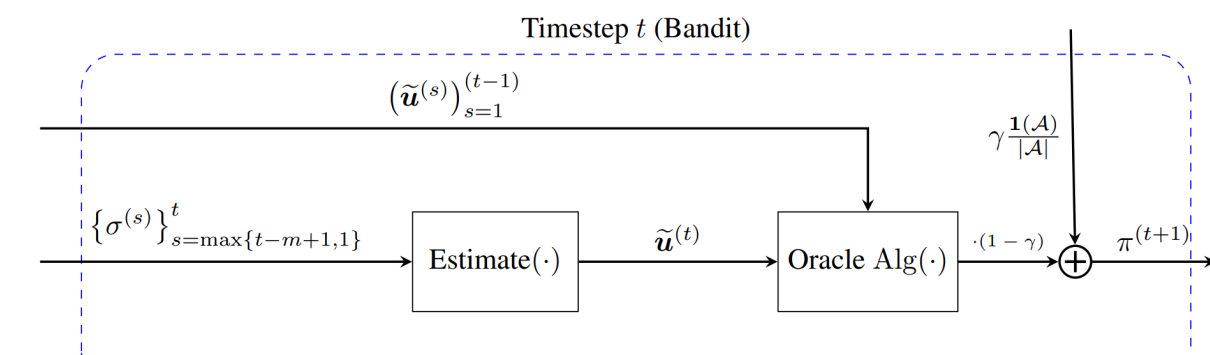
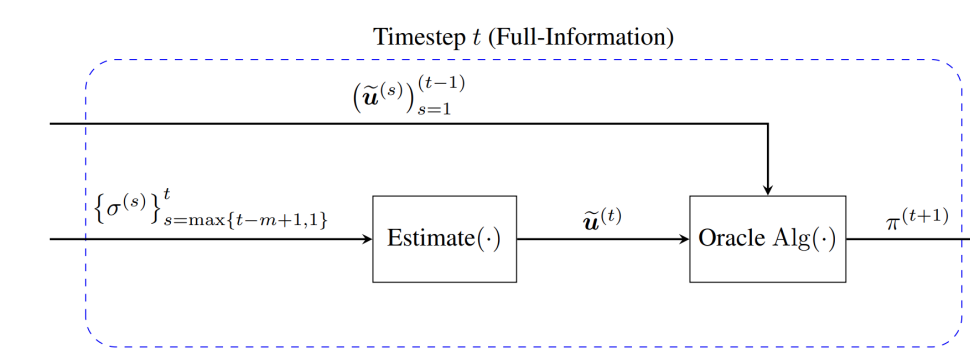
$$o^{(t)} = \mathcal{A}$$

$$R^{(T), \text{external}} := \max_{\hat{\pi} \in \Delta^{\mathcal{A}}} \sum_{t=1}^T \langle u^{(t)}, \hat{\pi} - \pi^{(t)} \rangle$$

### Bandit Feedback

$$o^{(t)} \sim \pi^{(t)}, \quad |o^{(t)}| = K$$

$$R^T := \max_{\hat{\pi} \in \Delta^{\mathcal{A}}} \sum_{t=1}^T \left( \langle u^{(t)}, \hat{\pi} \rangle - \frac{1}{K} \sum_{a \in o^{(t)}} u^t(a) \right)$$



## Ranking with Time-average Utility

### Full-Information Feedback

$$r^t = \frac{1}{t} \sum_{s=1}^t u^s$$

$$o^{(t)} = \mathcal{A}$$

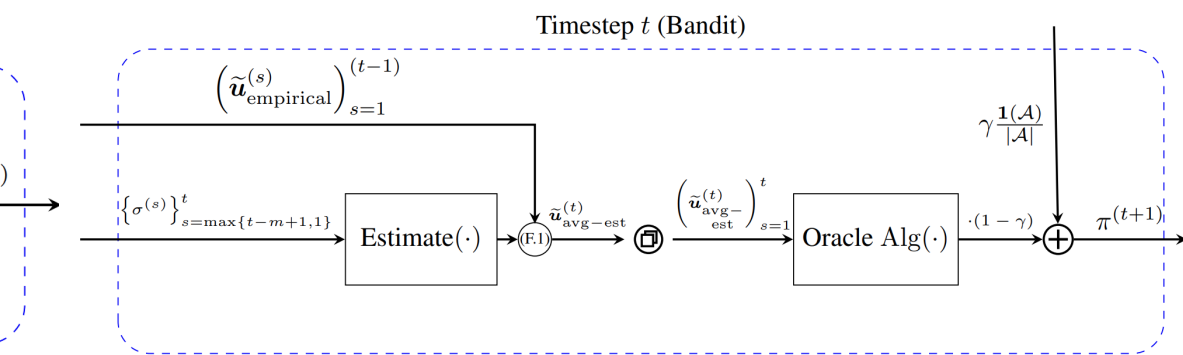
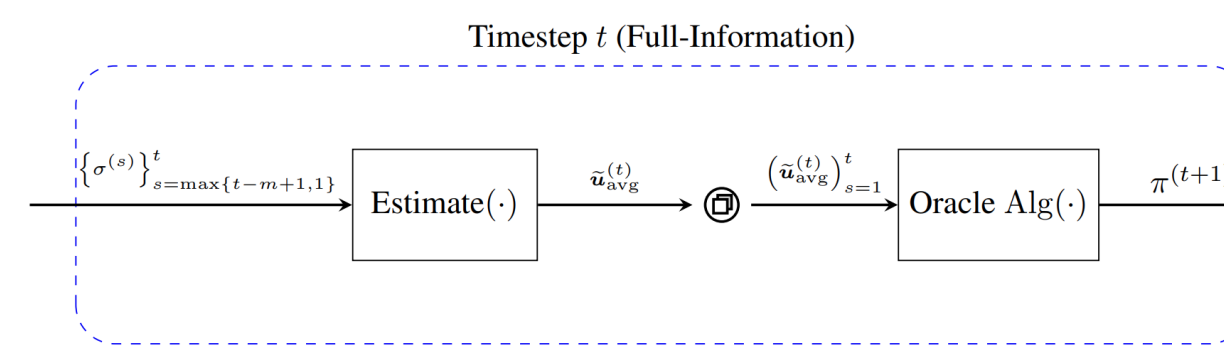
$$R^{(T), \text{external}} := \max_{\hat{\pi} \in \Delta^{\mathcal{A}}} \sum_{t=1}^T \langle u^{(t)}, \hat{\pi} - \pi^{(t)} \rangle$$

### Bandit Feedback

$$r^t = \frac{\sum_{s=1}^t u^s(a) \sum_{a' \in o^{(s)}} \mathbb{1}(a=a')}{\sum_{s=1}^t \sum_{a' \in o^{(s)}} \mathbb{1}(a=a')}$$

$$o^{(t)} \sim \pi^{(t)}, \quad |o^{(t)}| = K$$

$$R^T := \max_{\hat{\pi} \in \Delta^{\mathcal{A}}} \sum_{t=1}^T \left( \langle u^{(t)}, \hat{\pi} \rangle - \frac{1}{K} \sum_{a \in o^{(t)}} u^t(a) \right)$$



## Utility Estimation

**Algorithm 1** Utility Estimation with Action Permutations: Estimate  $\left( \left\{ \sigma^{(s)} \right\}_{s=1}^{m'} \right)$

- 1: **Input:** A set consisting of  $m'$  permutations of actions:  $\{\sigma^{(s)}\}_{s=1}^{m'}$  with  $|\sigma^{(s)}| = K$  for all  $s \in [m']$ , and temperature  $\tau > 0$ .
- 2: **for**  $j = 1, 2, \dots, |\mathcal{A}| - 1$  **do**
- 3:   **for**  $s = 1, \dots, m'$  **do**
- 4:     Calculate  $n_{j,1}^{(s)}, n_{j,2}^{(s)}$  defined as
 
$$n_{j,1}^{(s)} := \sum_{i,k \in [K]} \mathbb{1}(\sigma^{(s)}(i) = a^j, \sigma^{(s)}(k) = a^{|\mathcal{A}|} \text{ and } i < k),$$

$$n_{j,2}^{(s)} := \sum_{i,k \in [K]} \mathbb{1}(\sigma^{(s)}(i) = a^j, \sigma^{(s)}(k) = a^{|\mathcal{A}|} \text{ and } i > k).$$
- 5:   **end for**
- 6:   Let  $\mathcal{T}_j := \{s \in [m'] : n_{j,1}^{(s)} + n_{j,2}^{(s)} > 0\}$
- 7:   Let  $\text{sig}^{-1}(x): (0, 1) \rightarrow \mathbb{R} := \log \frac{x}{1-x}$  be the inverse function of  $\text{sig}(\cdot)$ . The utility of action  $a^j$  is then estimated as
 
$$\tilde{u}(a^j) = \begin{cases} \text{Proj}_{[-1,1]} \left( \tau \text{sig}^{-1} \left( \frac{1}{|\mathcal{T}_j|} \cdot \sum_{s \in \mathcal{T}_j} \left( \frac{n_{j,1}^{(s)}}{n_{j,1}^{(s)} + n_{j,2}^{(s)}} \right) \right) \right) & |\mathcal{T}_j| > 0 \\ 0 & |\mathcal{T}_j| = 0. \end{cases}$$
- 8:   **end for**
- 9: **Return**  $\tilde{\mathbf{u}} = (\tilde{u}(a^1), \tilde{u}(a^2), \dots, \tilde{u}(a^{|\mathcal{A}|}), 0)$

## Overview

$$R^T := \sum_{t=2}^T \|u^{(t)} - u^{(t-1)}\| \leq \mathcal{O}(T^\tau). \quad (2.1)$$

| Lower Bound  | Full-Information   | Bandit   |
|--|--|--|
| InstUtil Rank                                      | $\Omega(T)$ for $\tau \leq \mathcal{O}(1)$                                     |  |
| AvgUtil Rank                                       | $\tilde{\Omega}(T)$ for $\tau \leq \mathcal{O}\left(\frac{1}{T \log T}\right)$ | $\Omega(T)$ for $\tau \leq \mathcal{O}\left(\frac{1}{\log T}\right)$ |
| Upper Bound  |  |  |
| $(\tau = \mathcal{O}(1), \text{Sublinear Regret})$ | Full-Information   | Bandit   |
| InstUtil Rank                                      | $q < 1$  |  |
| AvgUtil Rank                                       | $\checkmark$   | $q < \frac{1}{3}$  |

Table: Summary of our contributions, including the *negative results* (top) and the *positive results* (bottom).  $\checkmark$  indicates that no additional assumptions are needed. Here,  $\tau > 0$  denotes the temperature parameter of the ranking model in the PL model.

## Computing Equilibrium in Games

**Assumption 6.1.** The (full-information) online learning algorithm Alg satisfies the following condition: for any  $T > 0$ ,  $t \in [T]$ , and sequences of utilities  $(u^{(s)})_{s=1}^t, (u^{(s)})_{s=1}^t \in (\mathbb{R}^{\mathcal{A}})^t$ , we have

$$\left\| \text{Alg} \left( (u^{(s)})_{s=1}^t \right) - \text{Alg} \left( (u^{(s)})_{s=1}^t \right) \right\| \leq L \left\| \sum_{s=1}^t u^{(s)} - \sum_{s=1}^t u^{(s)} \right\|,$$

where  $L = \Theta(T^{-c})$  for some constant  $c \in (0, 1)$ .

**Assumption 7.1** (Sublinear variation of strategies). The (full-information) online learning algorithm Alg needs to satisfy the following condition: for any  $T > 0$ ,  $t \in [T-1]$ , and sequence of utility vectors  $(u^{(s)})_{s=1}^{t+1} \in ([-1, 1]^{\mathcal{A}})^{t+1}$ , we have  $\left\| \text{Alg} \left( (u^{(s)})_{s=1}^t \right) - \text{Alg} \left( (u^{(s)})_{s=1}^{t+1} \right) \right\| \leq \eta$ , where  $\eta = \Theta(T^{-w})$  for some constant  $w \in (0, 1)$ .

**Theorem 7.2.** Consider InstUtil Rank with constant  $\tau > 0$  and Algorithm 2. For any  $\delta \in (0, 1)$ ,  $T > 0$ , and any full-information no-regret learning algorithm with numeric utility feedback, Alg, that satisfies Assumption 7.1, by choosing  $M, m, \gamma$  according to Theorem 5.2, we have that with probability at least  $1 - \delta$ , the algorithm finds an  $\epsilon$ -CCE, with

$$\epsilon \leq \max_{i \in [N]} \left\{ \frac{1}{T} R_i^{(T), \text{external}} \left( \text{Alg}, (\tilde{u}_i^{(t)})_{t=1}^T \right) \right\} + \mathcal{O} \left( \eta^{\frac{1}{3}} \left( \log \left( \frac{T}{\delta} \right) \right)^{\frac{1}{3}} \right) \quad (\text{Full-Information})$$

$$\epsilon \leq \max_{i \in [N]} \left\{ \frac{1}{T} R_i^{(T), \text{external}} \left( \text{Alg}, (\tilde{u}_i^{(t)})_{t=1}^T \right) \right\} + \mathcal{O} \left( \eta^{\frac{1}{5}} \log \left( \frac{T}{\delta} \right) \right). \quad (\text{Bandit})$$

**Theorem 7.3.** Consider AvgUtil Rank with constant  $\tau > 0$  and Algorithm 3. For any  $\delta \in (0, 1)$ ,  $T > 0$ , and any full-information no-regret learning algorithm with numeric utility feedback, Alg, that satisfies Assumption 6.1, by choosing  $M, m, \gamma$  according to Theorem 6.2, we have that with probability at least  $1 - \delta$ , the algorithm finds an  $\epsilon$ -CCE under full-information feedback, with

$$\epsilon \leq \max_{i \in [N]} \left\{ \frac{1}{T} R_i^{(T), \text{external}} \left( \text{Alg}, (u_i^{(t)})_{t=1}^T \right) \right\} + \mathcal{O} \left( L T^{\frac{1}{5}} \log \left( \frac{T}{\delta} \right) \right). \quad (\text{Full-Information})$$

When  $M, m$ , and  $\gamma$  are chosen as in Theorem 6.3, and both Assumption 6.1 and Assumption 7.1 hold, we have that with probability at least  $1 - \delta$ , the algorithm finds an  $\epsilon$ -CCE under bandit feedback, with

$$\epsilon \leq \max_{i \in [N]} \left\{ \frac{1}{T} R_i^{(T), \text{external}} \left( \text{Alg}, (u_i^{(t)})_{t=1}^T \right) \right\} + \tilde{\mathcal{O}} \left( \left( \log \left( \frac{1}{\delta} \right) \right)^2 \left( L^{\frac{1}{5}} \eta^{\frac{1}{5}} + L^{\frac{1}{5}} \right) T^{\frac{1}{5}} \right). \quad (\text{Bandit})$$

## LLM Routing with Ranking Feedback

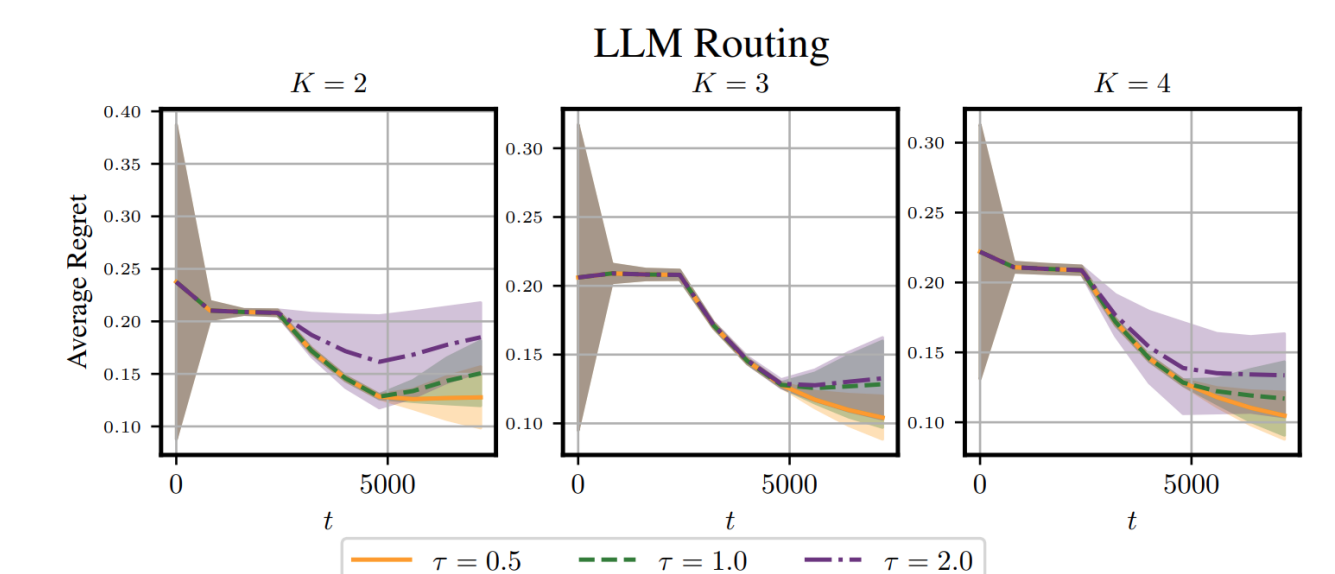


Figure:  $R^T$  with AvgUtil under bandit feedback for different temperatures  $\tau$  and numbers of proposed actions  $K$  in the online learning setting.