

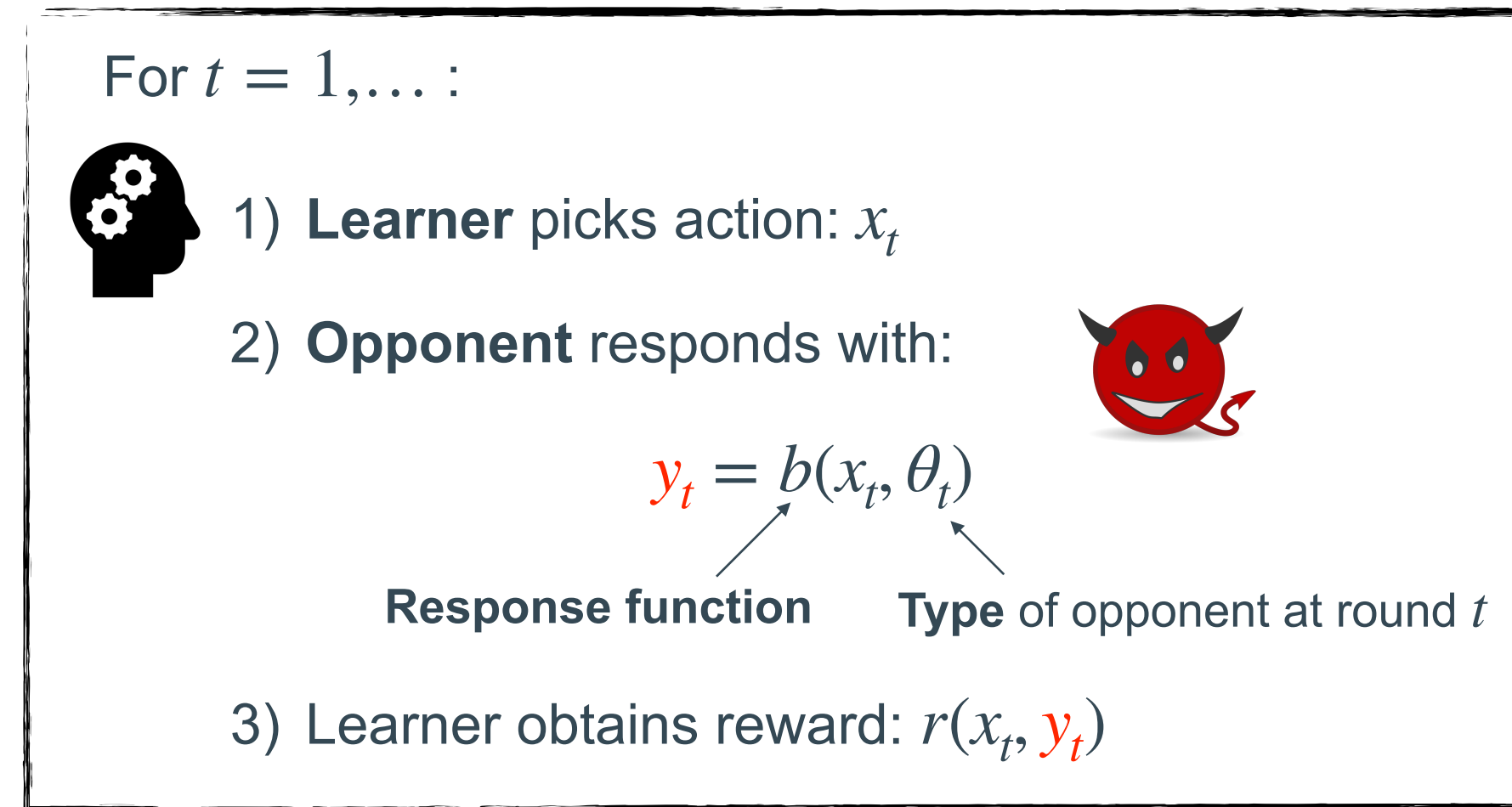
Learning to Play Sequential Games versus Unknown Opponents

Pier Giuseppe Sessa
Ilija Bogunovic
Maryam Kamgarpour
Andreas Krause

ETH zürich



Repeated Sequential Game Setup



- Response function $b(\cdot, \cdot)$ is **unknown** to the learner. Learner only observes $y_t + \epsilon_t$, with ϵ_t zero-mean sub-Gaussian.
- Opponent's types can be adversarially selected.

Learner's **regret** as a performance indicator:

$$R(T) := \max_{x \in \mathcal{X}} \sum_{t=1}^T r(x, b(x, \theta_t)) - \sum_{t=1}^T r(x_t, y_t).$$

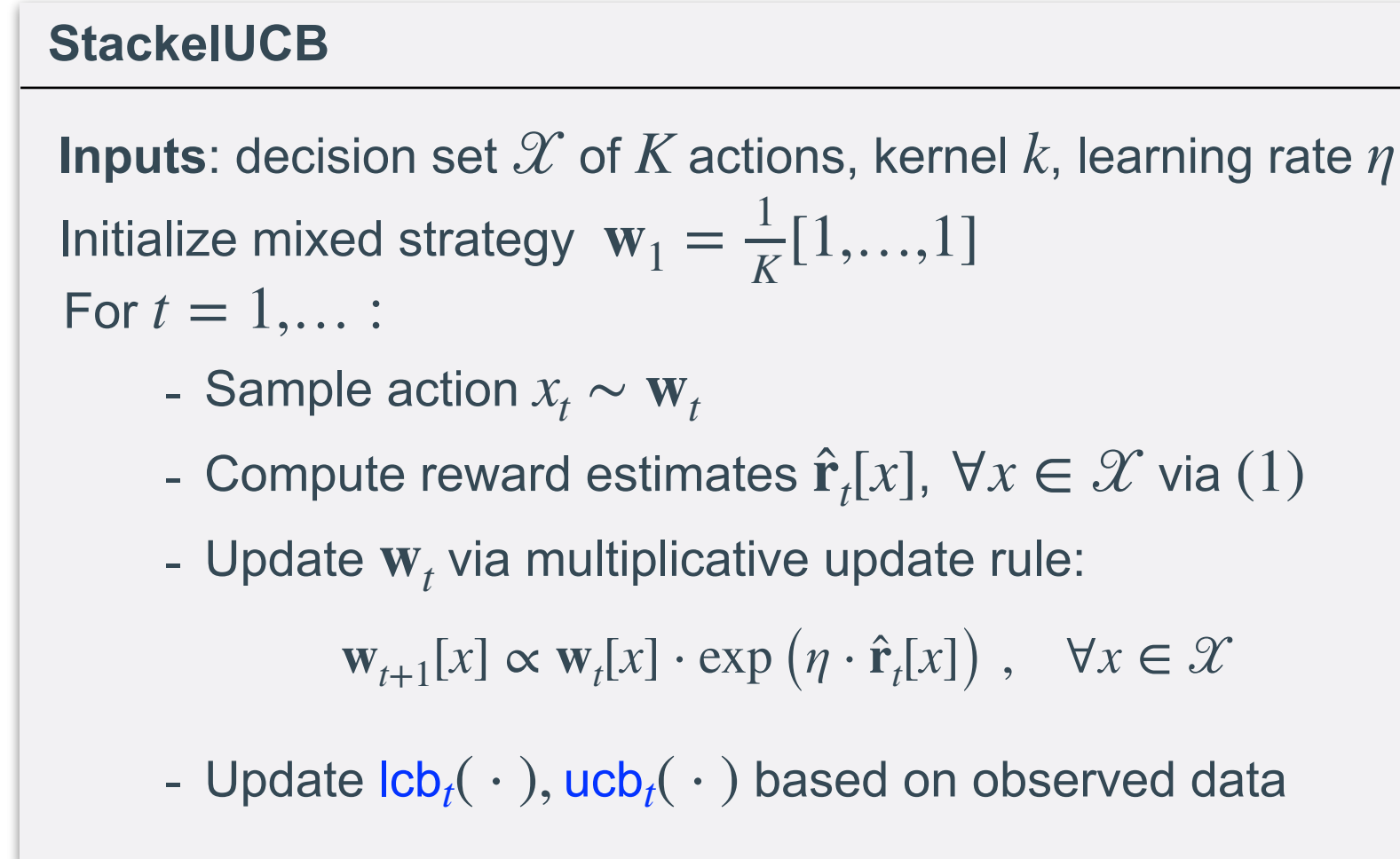
High-Level Ideas

- Iteratively learn the opponent response function via **kernel-ridge regression**
- Use an online learning strategy with **optimistic reward estimates**:

$$\hat{\mathbf{r}}_t[x] = \max_y r(x, y) \quad (1)$$

s.t. $y \in [\text{lcb}_t(x, \theta_t), \text{ucb}_t(x, \theta_t)]$.

The **StackelUCB** Algorithm



Thm (informal): Assume $\|b(\cdot)\|_k$ is bounded and r is L -Lipschitz. Then, w.h.p.,

$$R(T) \leq \mathcal{O}\left(\sqrt{T \log K}\right) + \mathcal{O}\left(L\gamma_T \sqrt{T}\right).$$

Kernel-dependent **max info. gain** [2] γ_T :

Maximal uncertainty reduction about $b(\cdot)$ after T noisy observations. E.g., $\gamma_T = \mathcal{O}((\log T)^{d+1})$ for SE kernels with domain dimension d .

Single Opponent Type (i.e., when $\theta_t = \bar{\theta} \forall t$)

Corollary: Consider the strategy: $x_t = \arg \max_{x \in \mathcal{X}} \hat{\mathbf{r}}_t[x]$

Then, w.h.p.,

$$R(T) \leq \mathcal{O}\left(L\gamma_T \sqrt{T}\right).$$

Stackelberg Games with unknown followers' utilities

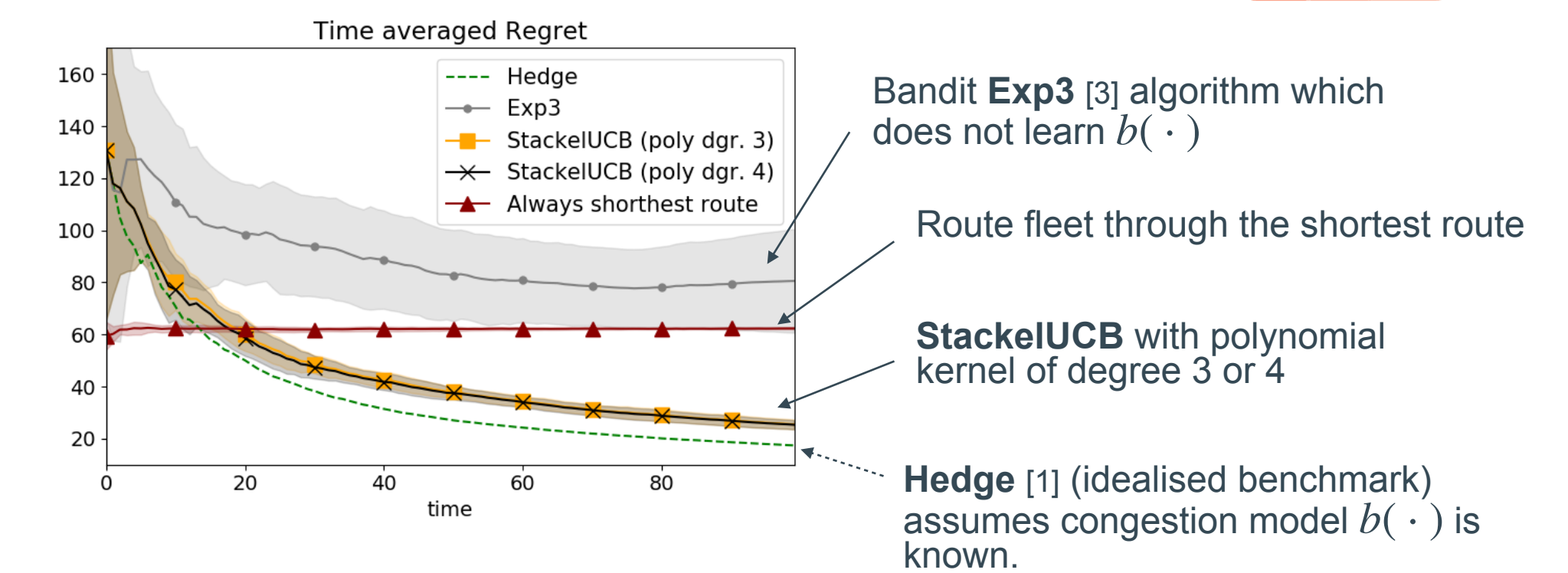
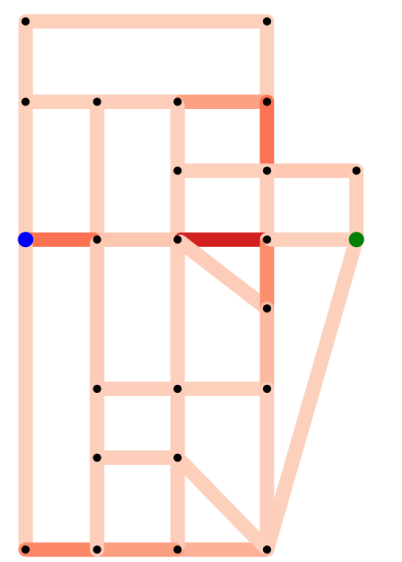
$\mathcal{X} = n_t$ -dimensional simplex and $b(\cdot, \theta_t) = \arg \max_y U_f(\cdot, \theta_t)$.

Corollary: When **StackelUCB** is run over a finite discretisation of \mathcal{X} , w.h.p.

$$R(T) \leq \mathcal{O}\left(\sqrt{T n_t \log(L\sqrt{n_t T})}\right) + \mathcal{O}\left(L\gamma_T \sqrt{T}\right).$$

Routing Vehicles in Congested Traffic Network

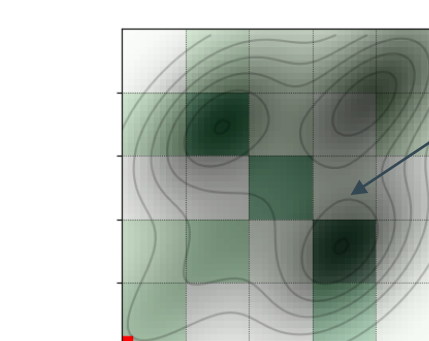
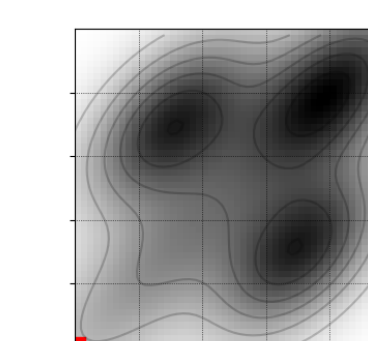
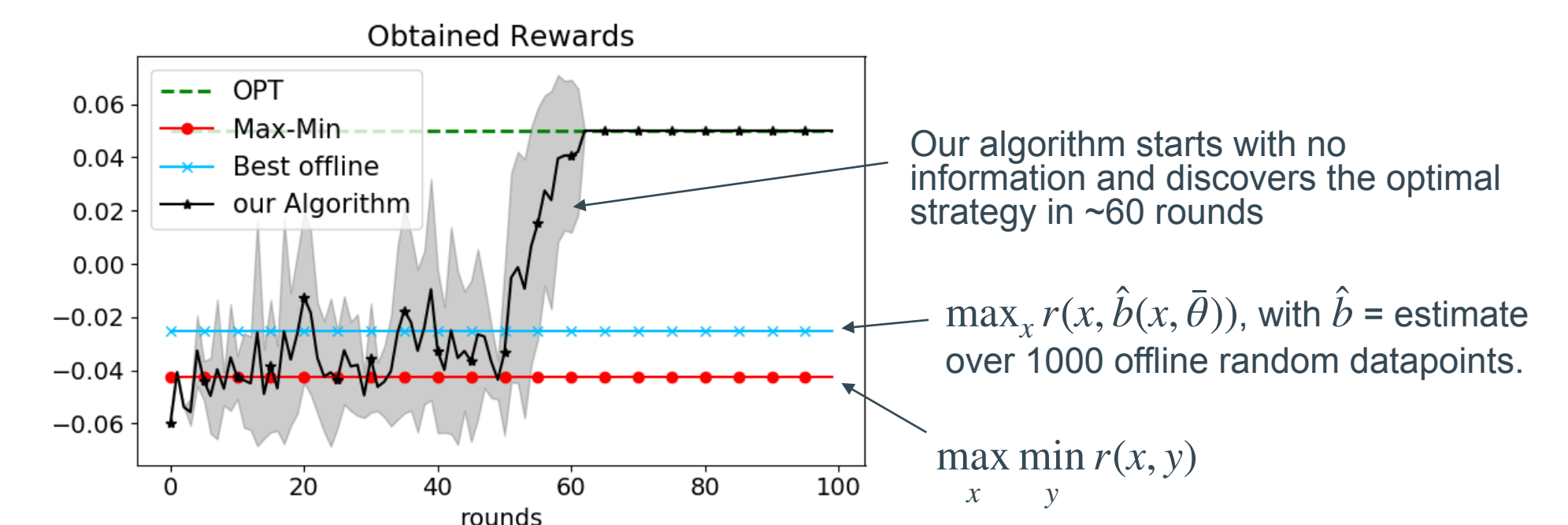
- Network operator chooses routes x_t for fleet of vehicles
 - Based on x_t , users in the network pick their routes and cause congestion: $y_t = b(x_t, \theta_t)$
- ↑
users' demand profile
(i.e., users origins and destinations)



- StackelUCB** leads to low regret and reduces the network's congestion

Wildlife Protection against Poaching Activities

- Park rangers choose a patrol strategy x_t (i.e., $x_t[i]$ = probability of patrolling area i in the park)
 - Poachers choose poaching location $y_t = b(x_t, \bar{\theta})$
- ↑
park animal density
poachers' preference model



Park animal density Rangers patrol strategy.

High prob. assigned to areas with high animal density and closer to poachers' base location (unknown to the algorithm)



References

- [1] Y. Freund and R. E. Schapire. "A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting". In *J. Comput. Syst. Sci.*, 1997.
- [2] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger. "Gaussian process optimization in the bandit setting: No regret and experimental design". In *ICML*, 2010.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. "The Nonstochastic Multiarmed Bandit Problem". In *SIAM J. Comput.*, 2003.