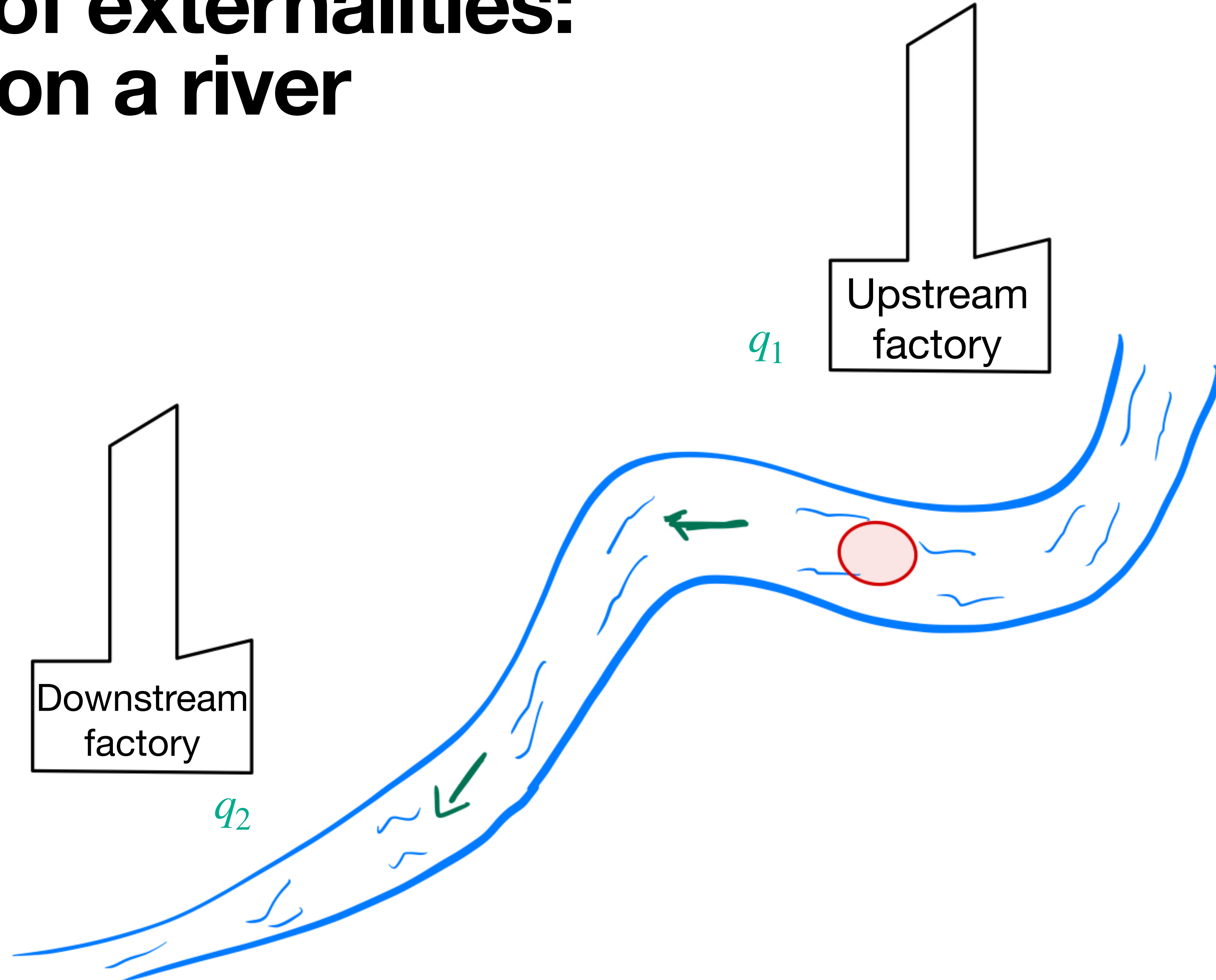


Mitigating Externalities while Learning

A. Scheid, A. Capitaine, E. Boursier, E. Moulines, M. Jordan, A. Durmus

NeurIPS 2024

Example of externalities: factories on a river



Example: factories on a river

Upstream utility: $\pi_1(q_1)$ increases with q_1

Downstream utility: $\pi_2(q_1, q_2)$ increases with q_2 and decreases with q_1



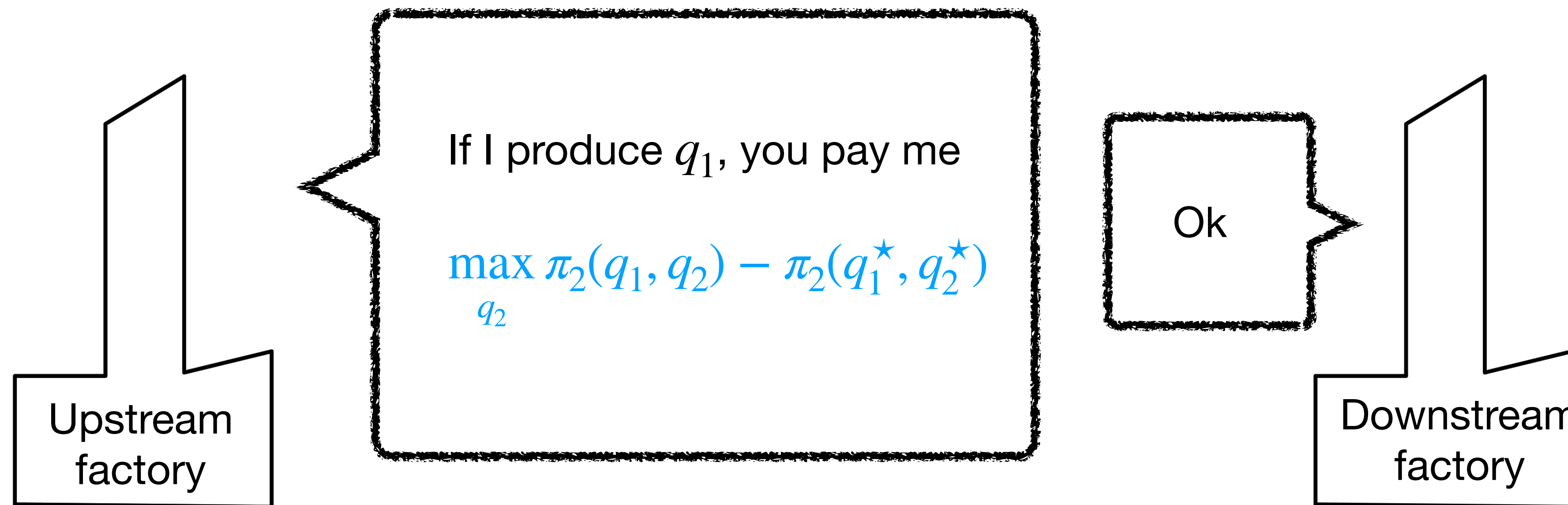
Upstream optimal strategy: $q_1^* = \arg \max_q \pi_1(q)$

Downstream optimal strategy: $q_2^* = \arg \max_q \pi_2(q_1^*, q)$

Social inefficiency: $\pi_1(q_1^*) + \pi_2(q_1^*, q_2^*) < \max_{q_1, q_2} \{ \pi(q_1) + \pi(q_1, q_2) \}$

Recovering social optimum

Idea: allowing proprietary rights to restore social efficiency



→ recover **social optimum**

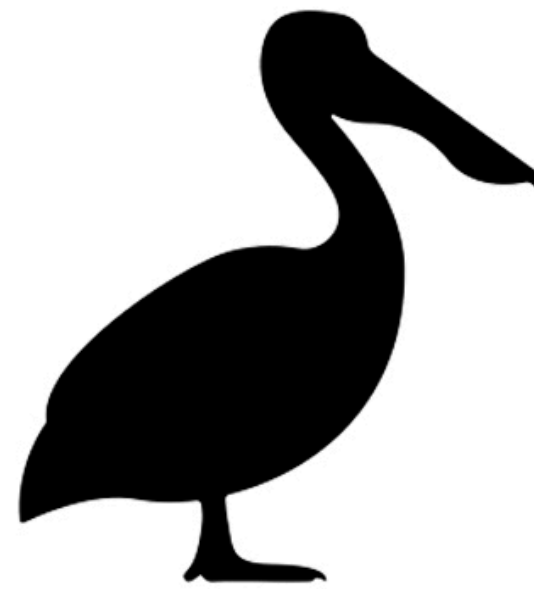
$$\max_{q_1, q_2} \pi_1(q_1) + \pi_2(q_1, q_2)$$

Issue: the players do not know their utilities!

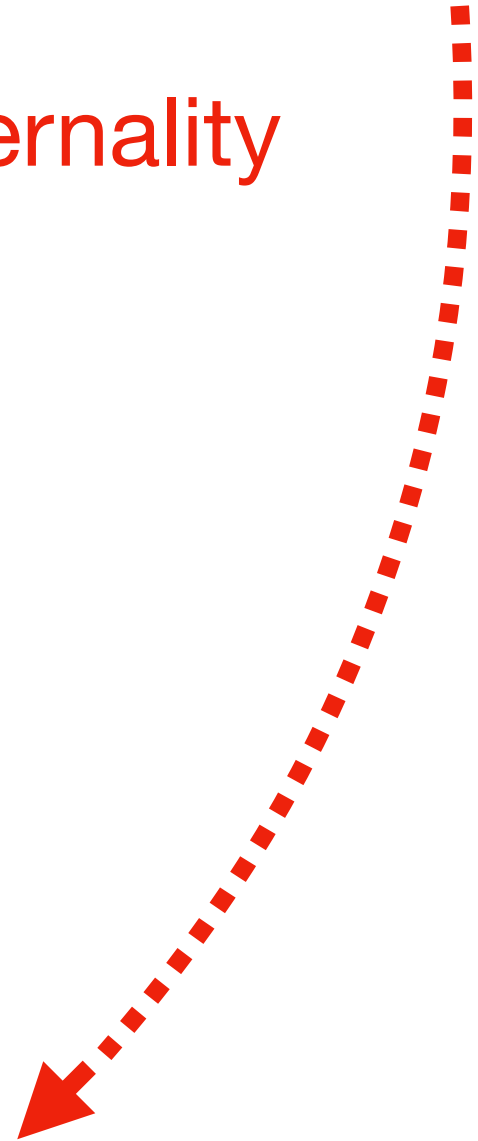
What if we learn utility functions over time?

We can recover social optimum through transfers

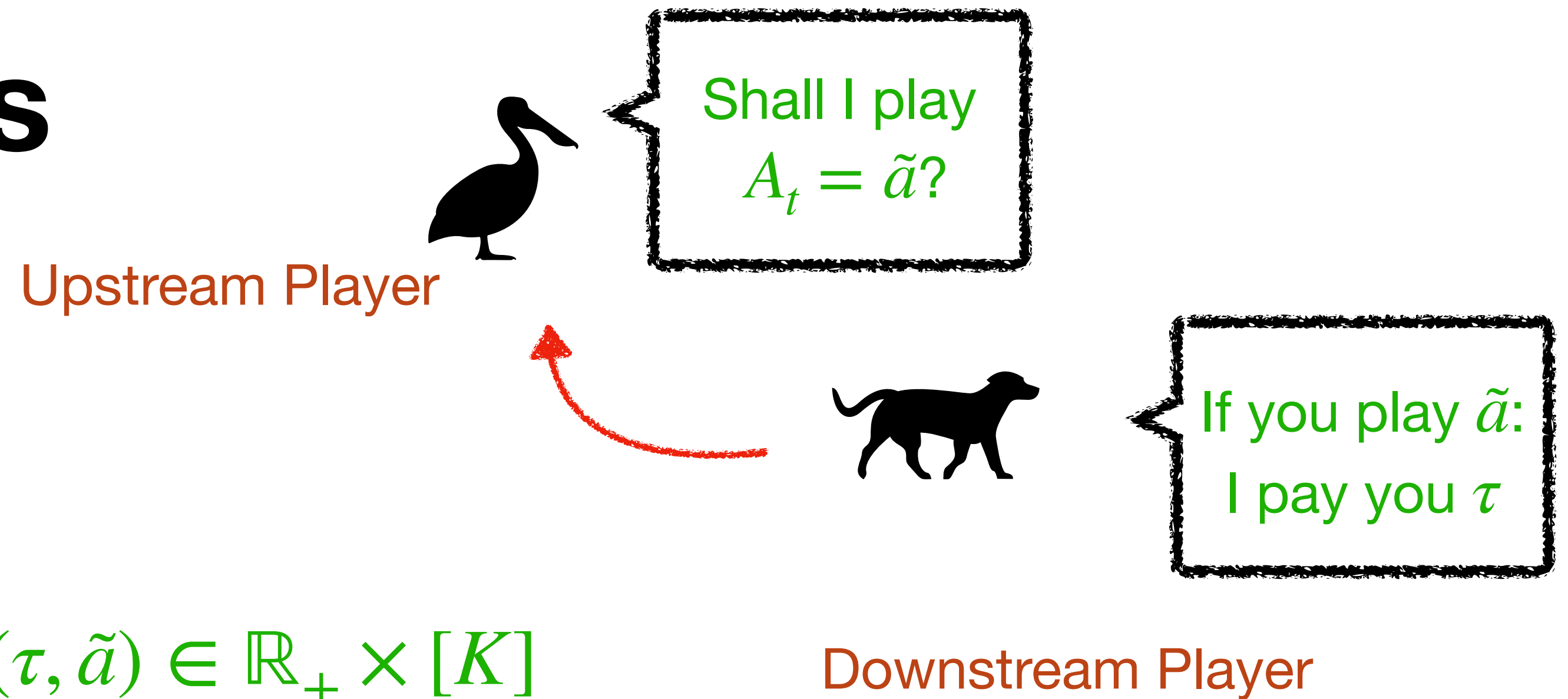
*For instance:
Resource sharing*



Externality

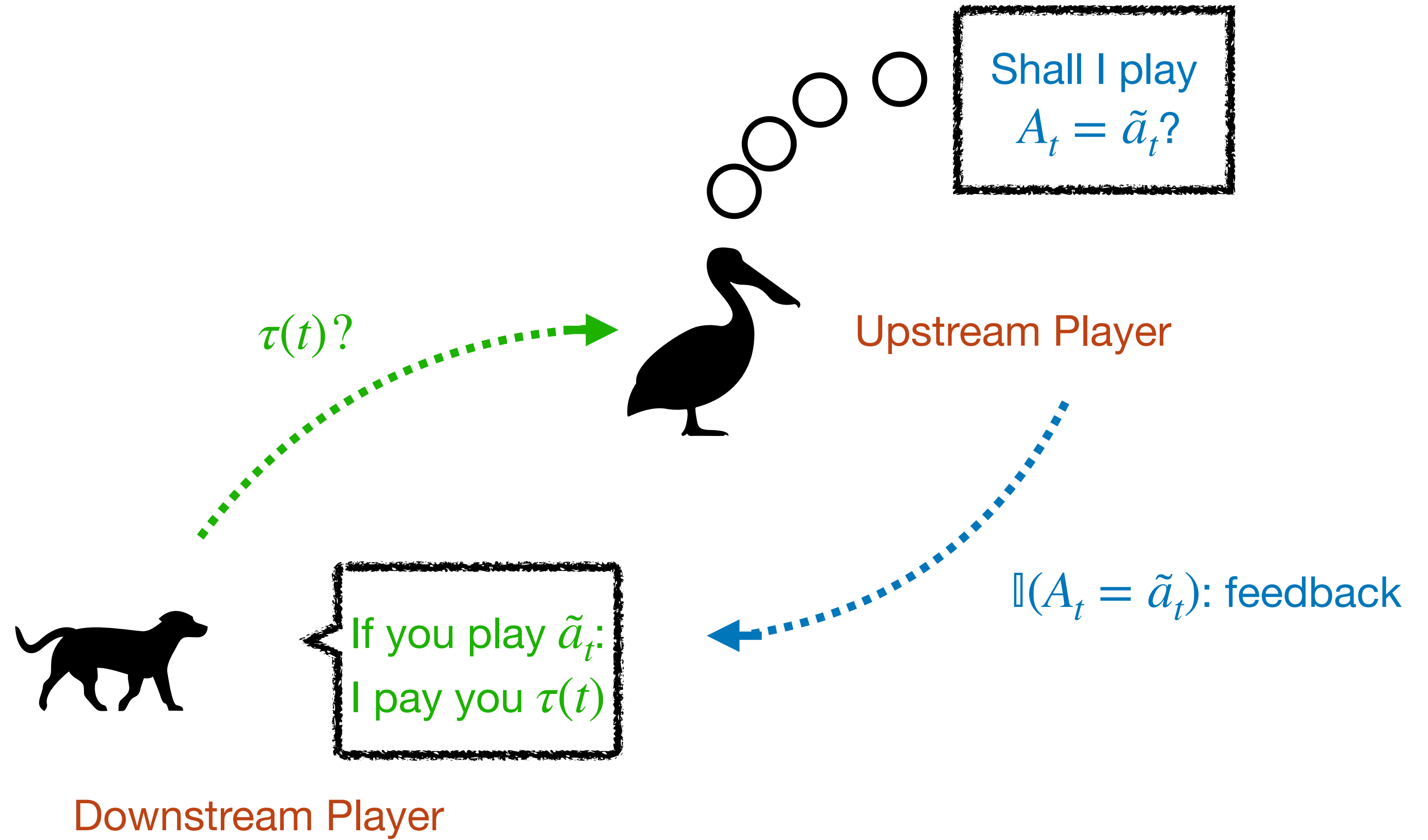


Model with transfers



At each time t :

1. Downstream player proposes payment $(\tau, \tilde{a}) \in \mathbb{R}_+ \times [K]$
2. Upstream player:
 - observes $(\tau(t), \tilde{a}_t)$
 - plays $A_t \in [K]$
 - receives and observes: $Z_{A_t}(t) + \mathbf{1}_{A_t=\tilde{a}} \tau$
3. Downstream player:
 - plays $B_t \in [K]$
 - gets reward $X_{A_t, B_t}(t) - \mathbf{1}_{A_t=\tilde{a}} \tau$
 - observes separately $X_{A_t, B_t}(t)$ and $\mathbf{1}_{A_t=\tilde{a}}$



Designing no regret strategies

Upstream player is a no-regret learner

Downstream policy

Idea: run batched binary searches to find τ_a^\star : *the minimal incentive to have $A_t = a$*

- Propose payment (a, τ) for T^α successive time steps
- Observe T^\neq the number of times upstream did **not** pull a

Using the **no-regret assumption**, w.h.p.

- If $T^\neq > CT^{\alpha\kappa+\beta}$, then $\tau_a^\star > \tau - \frac{1}{T^\beta}$

For any β s.t. $\alpha\kappa + \beta < \alpha$

- If $T^\neq < T^\alpha - CT^{\alpha\kappa+\beta}$, then $\tau_a^\star < \tau + \frac{1}{T^\beta}$

Downstream policy

Theorem

For a good choice of parameters:

$$\mathcal{R}^d(T) = o(T)$$

- most of the regret is due to *waiting for the upstream player learning*
- the faster does the upstream player learns, the better for the downstream one
- upstream and downstream players can typically use UCB as a subroutine

Conclusion

Summary:

- Study a repeated two player games with an upstream/downstream relation
- We propose a downstream algorithm that works for general upstream policies

Direct extensions:

- Instance dependent bounds
- Anytime policy
- Extension to linear contextual case

Perspectives

Higher level questions:

- More general interactions between multiple learning agents
- Propose *black box independent* strategies
- Potential long term strategic manipulations

Thank you!