

Learning in Games with Best-Response Oracles

Tomer Koren (Google Research)

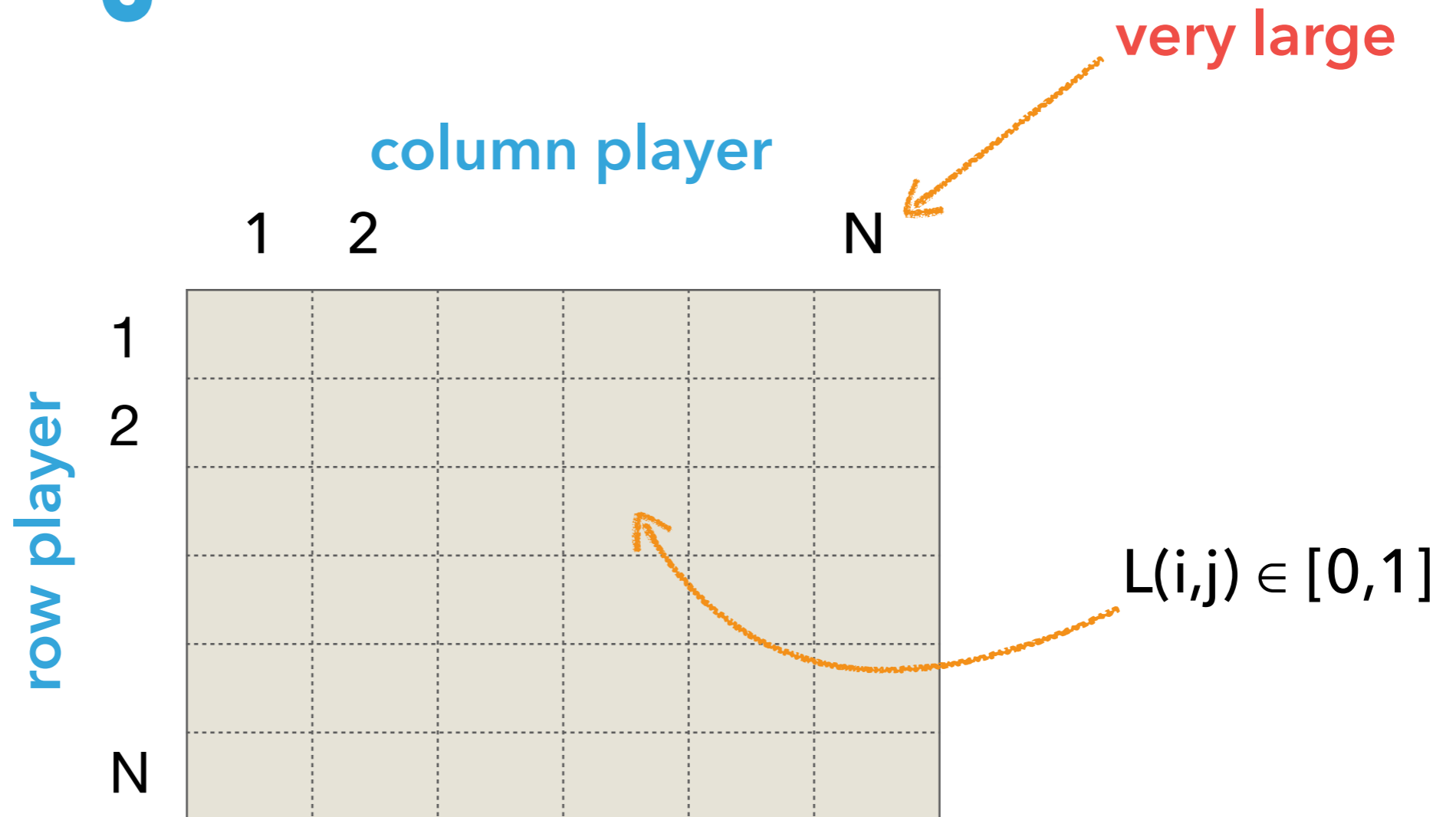
joint work with

Elad Hazan (Princeton)



PRINCETON
UNIVERSITY

Zero-sum games



TIME TO
COMPUTE /
APPROX?

von Neumann equilibrium:

$$\lambda = \min_{p \in \Delta} \max_{q \in \Delta} \mathbb{E}_{\substack{i \sim p \\ j \sim q}} [L(i,j)]$$

EFFICIENT STRATEGIES
FOR PLAYERS ?

Fictitious Play [Brown '49]

players plays their **best-response** to empirical dist. of opponent's past plays

- ▶ FP converge to minimax [Robinson '51]
- ▶ but, might need $\Omega(2^N)$ iterations [Brandt+ '10]
- ▶ convergence rate is $\Omega(T^{-1/N})$ [Daskalakis-Pan '14] (refutes Karlin's conjecture [Karlin '59])

$$y_3 = \operatorname{argmin}_y \frac{1}{2}L(x_1, y) + \frac{1}{2}L(x_2, y)$$

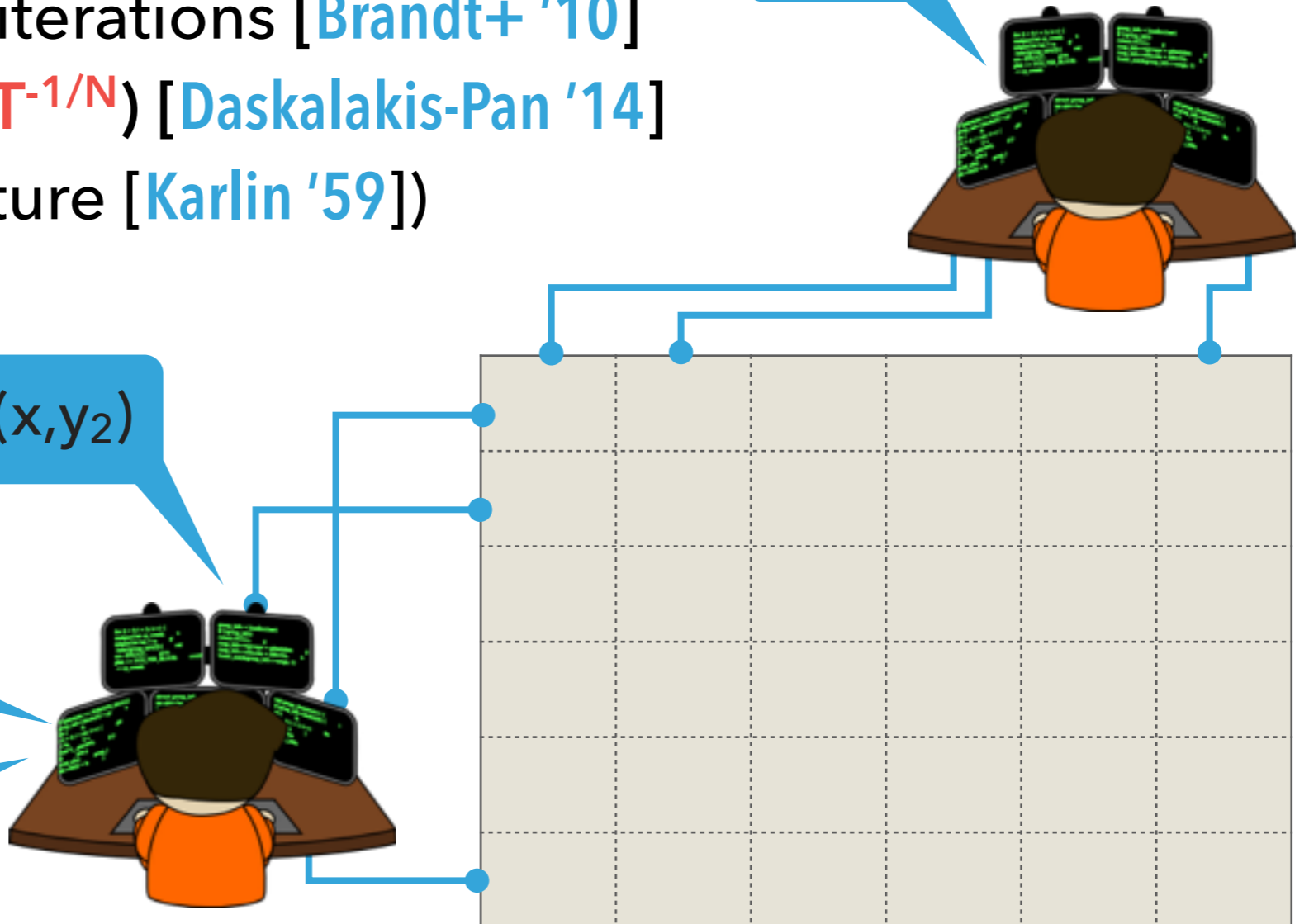
$$y_2 = \operatorname{argmin}_y L(x_1, y)$$

y_1

$$x_3 = \operatorname{argmin}_x \frac{1}{2}L(x, y_1) + \frac{1}{2}L(x, y_2)$$

$$x_2 = \operatorname{argmin}_x L(x, y_1)$$

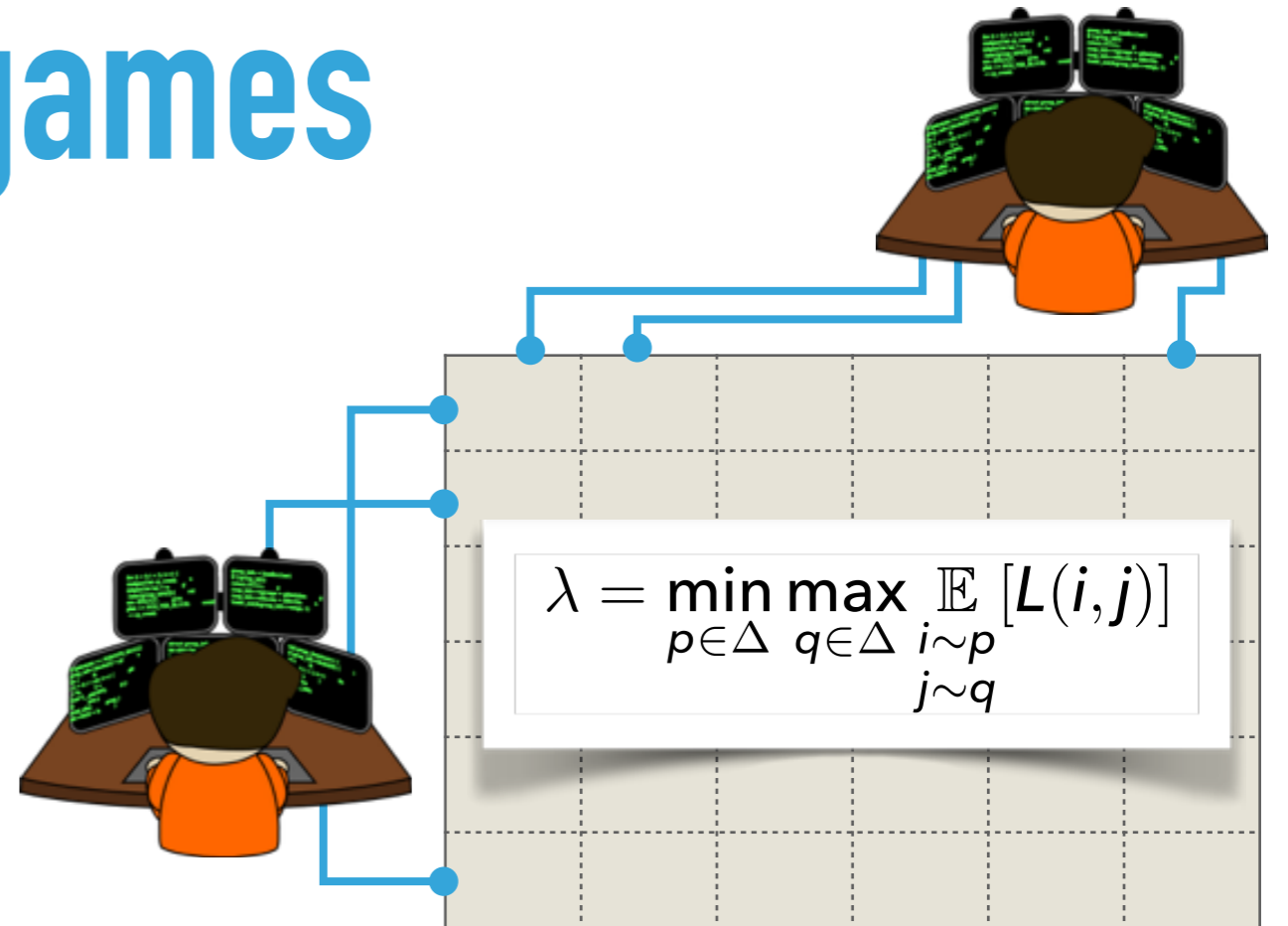
x_1



Solving zero-sum games

- ▶ Poly-time since the 70's...
(equivalent to LP)
- ▶ state of the art:
 $\tilde{O}(N)$ time algorithm, **tight**
[Grigoriadis-Khachiyan '95]
- ▶ $\tilde{O}(N)$ time via regret minimization
[Freund-Schapire '99]
- ▶ $\tilde{O}(N)$ time for generalized
minimax problems [Clarkson+ '10]

SUBLINEAR IN
INPUT SIZE



- ▶ More recent results:
poly(N) / T convergence rates
[Daskalakis+ '11, Rakhlin-Sridharan '13]

this talk:
focus on **N**

Learning in zero-sum games

[Freund-Schapire '99]

Players use **online learning** algos
(e.g., Multiplicative Weights)

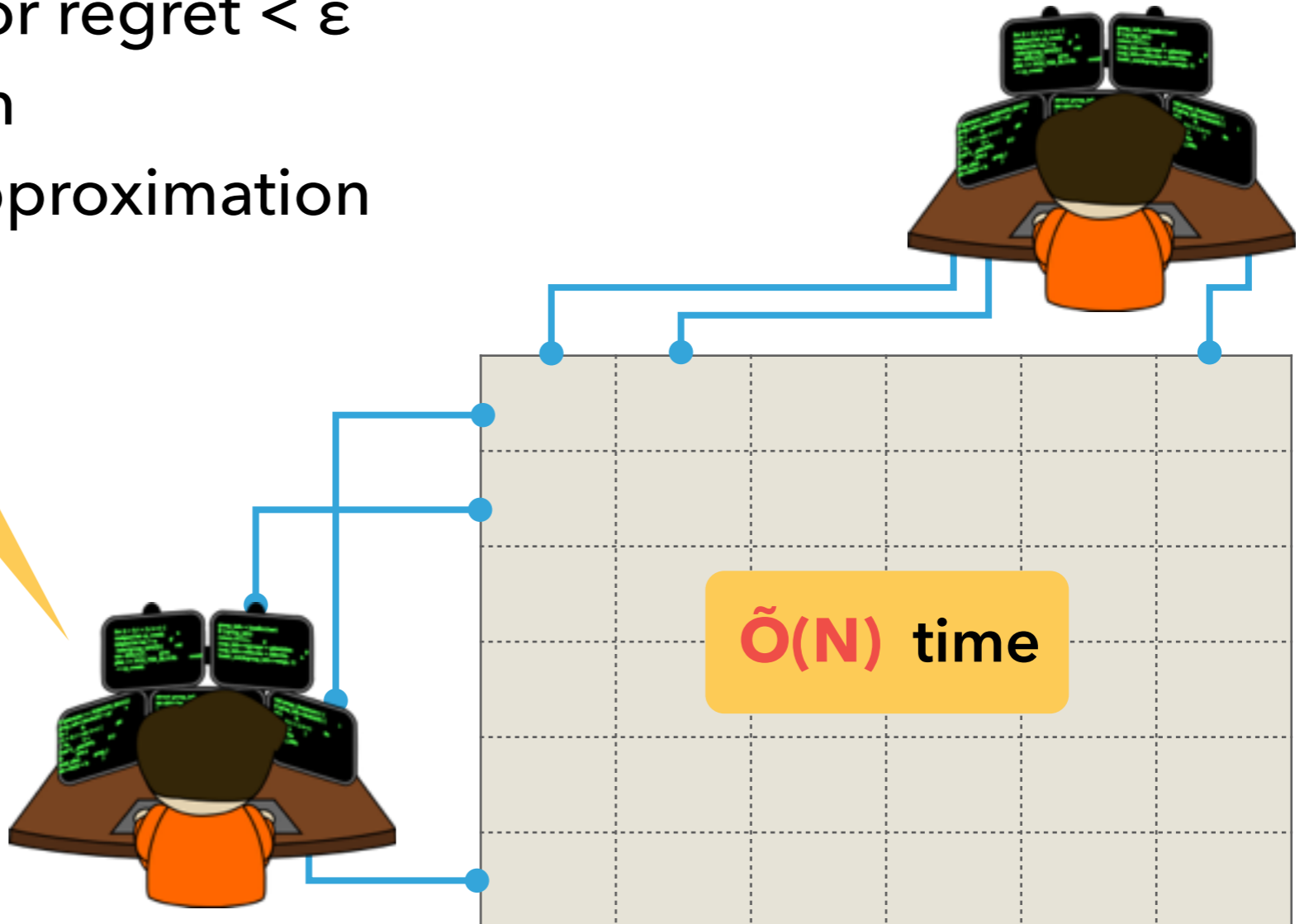
- ▶ $\log(N) / \epsilon^2$ iterations for regret $< \epsilon$
- ▶ $O(N)$ time per iteration
- $O(N / \epsilon^2)$ time for ϵ -approximation

REGRET: $\sqrt{\frac{\log N}{T}}$

$O(N)$ time

REGRET: $\sqrt{\frac{\log N}{T}}$

$O(N)$ time



Can we do better?

Games are often **exponentially large**

- ▶ $X = \{ \text{all } (s,t)\text{-paths in a given graph} \}$
 $Y = \{ \text{costs on edges} \}$
- ▶ $X = \{ \text{all permutations over } [n] \}$
 $Y = \{ \text{value assignments to items} \}$
- ▶ $X = \{ \text{subsets of } [n] \}$
 $Y = \{ \text{submodular evaluation functions} \}$



But best-response / optimization is **poly-time** = **poly(log N)**

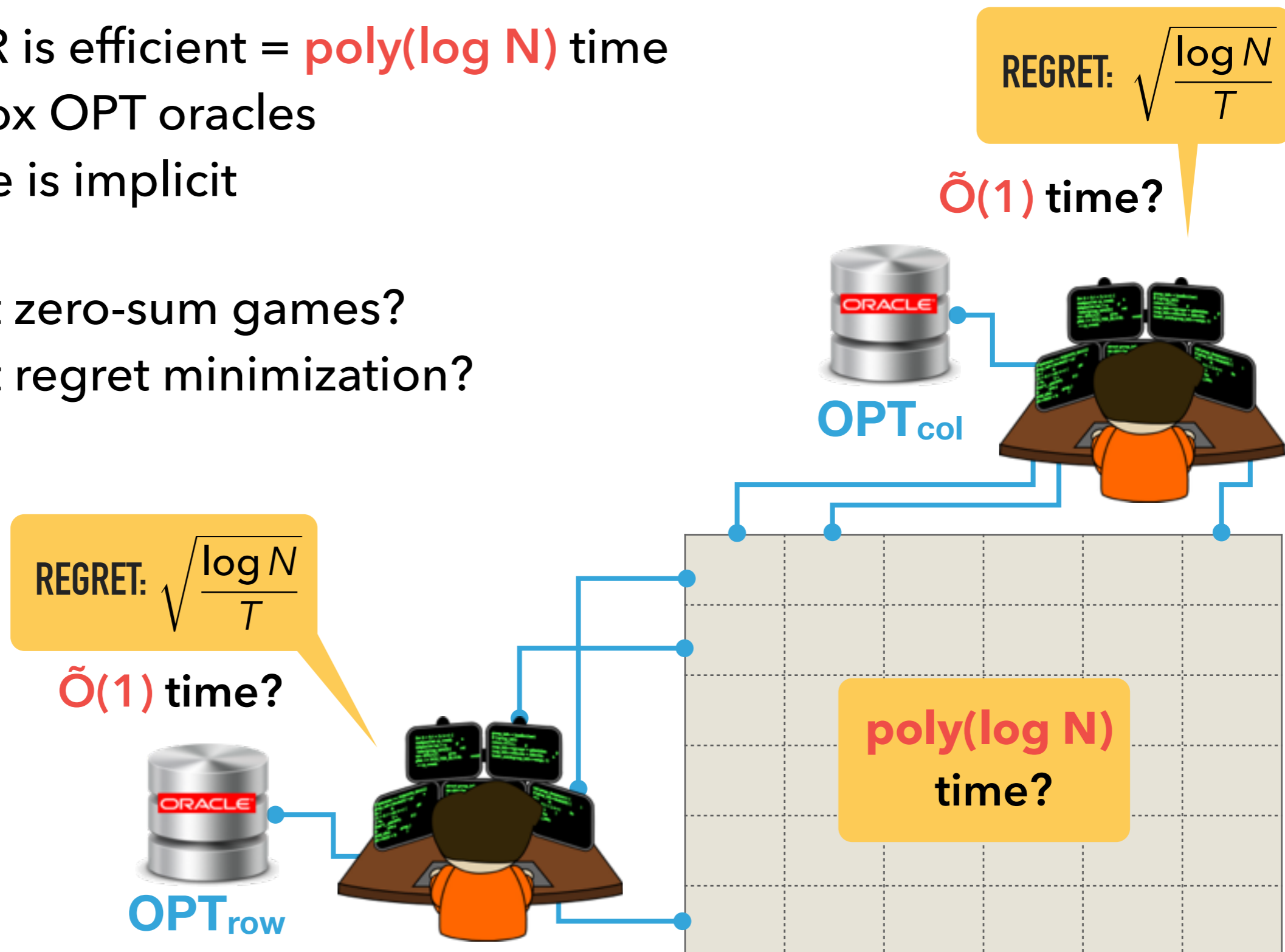
Best response oracles

assume BR is efficient = **poly(log N)** time

- ▶ black-box OPT oracles
- ▶ structure is implicit

→ efficient zero-sum games?

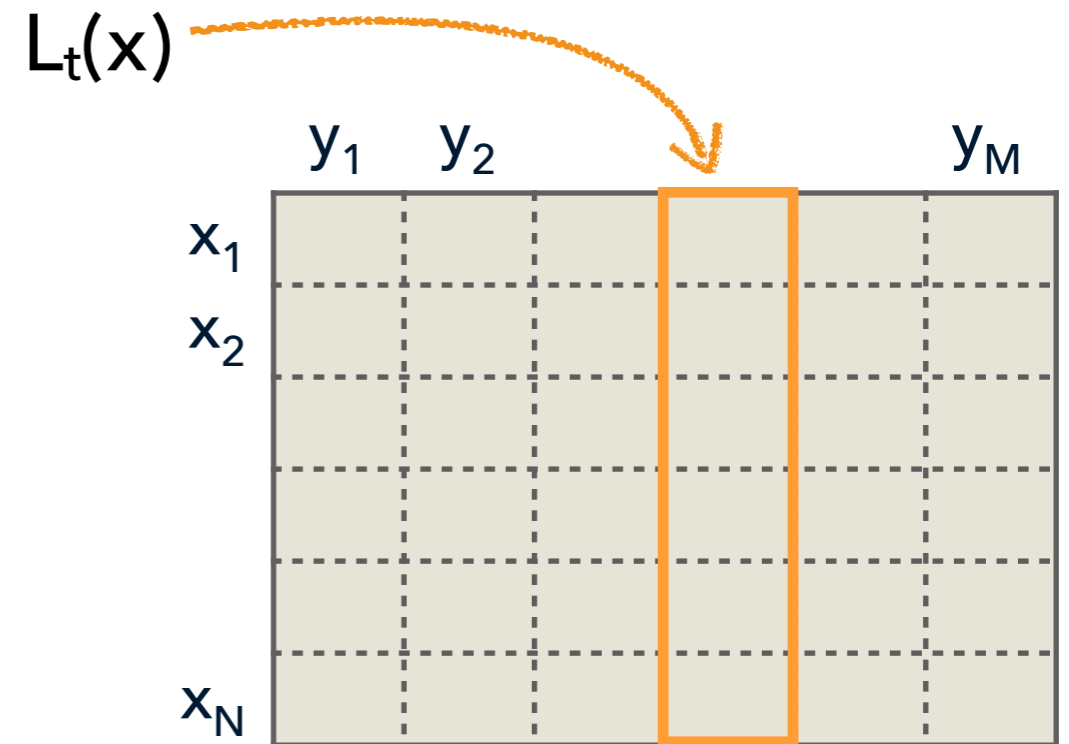
→ efficient regret minimization?



Online Learning

Iteratively, for $t=1,2,\dots,T$:

- (1) player: $x_t \in X$ ("expert")
- (2) adversary: $y_t \in Y$
- (3) player's loss = $L(x_t, y_t) = L_t(x_t)$



- ▶ Goal: minimize **regret**:
(average, expected)

$$\frac{1}{T} \sum_{t=1}^T L_t(x_t) - \frac{1}{T} \min_{x^*} \sum_{t=1}^T L_t(x^*) \rightarrow 0$$

- ▶ **Value access** to matrix:

$$\text{VAL}(x, y) = L(x, y)$$

$\tilde{O}(1)$ time

- ▶ **Best Response oracle**:

$$\text{OPT}(S \subseteq Y) = \operatorname{argmin}_x \sum_{y \in S} L(x, y)$$

Learn **efficiently**? Regret $< \frac{1}{4}$ in **poly(log N)** time?

Learning-theoretic motivation

- ▶ Fundamental question in learning theory:
generic & efficient reduction of online learning to optimization?
(analogous to fundamental theorem of statistical learning)
- ▶ Many **specialized** online algorithms for optimizable settings:
submodular opt., network routing, online PCA, contextual bandits,
online ranking,...
- ▶ Practical – numerous previous attempts:
Online convex optimization [[Zinkevich '03](#), [Hazan+ '06](#)],
Follow the Perturbed Leader (FPL) [[Hannan '57](#), [Kalai-Vempala '06](#)],
Dropout perturbation [[vanErven-Kotlowski-Warmuth '14](#)],
Contextual bandits [[Agarwal+ '14](#)], ...
☞ typically **poly(log N)** computation, but need **explicit structure**

Results

In OPT oracle model:

- ▶ **Thm 1.** Any algo that approximates $N \times N$ zero-sum games to within $\varepsilon = 1/4$ runs in total time $\tilde{O}(\sqrt{N})$
▶ $\tilde{O}(N)$ time needed to minimize regret
- ▶ **Thm 2.** There exists (new) online learning algo that attains regret $< \varepsilon$ in total time $\tilde{O}(\sqrt{N} / \varepsilon^2)$, tight
▶ vs. $\Theta(N/\varepsilon^2)$ time w/o OPT oracle
- ▶ **Corr.** There exists (new) algo that approximates $N \times N$ zero-sum games in total time $\tilde{O}(\sqrt{N})$, tight
▶ vs. $\Theta(N)$ time without oracles

4TH ROOT OF
INPUT SIZE

Aldous' random walk function

[Aldous '83, Aaronson '06]

Dist. over functions that have a single local min.

(1) start RW from uniform vertex of d -dim cube

(2) $f(i)$ = time to hit vertex i

$$V = \{0,1\}^d$$

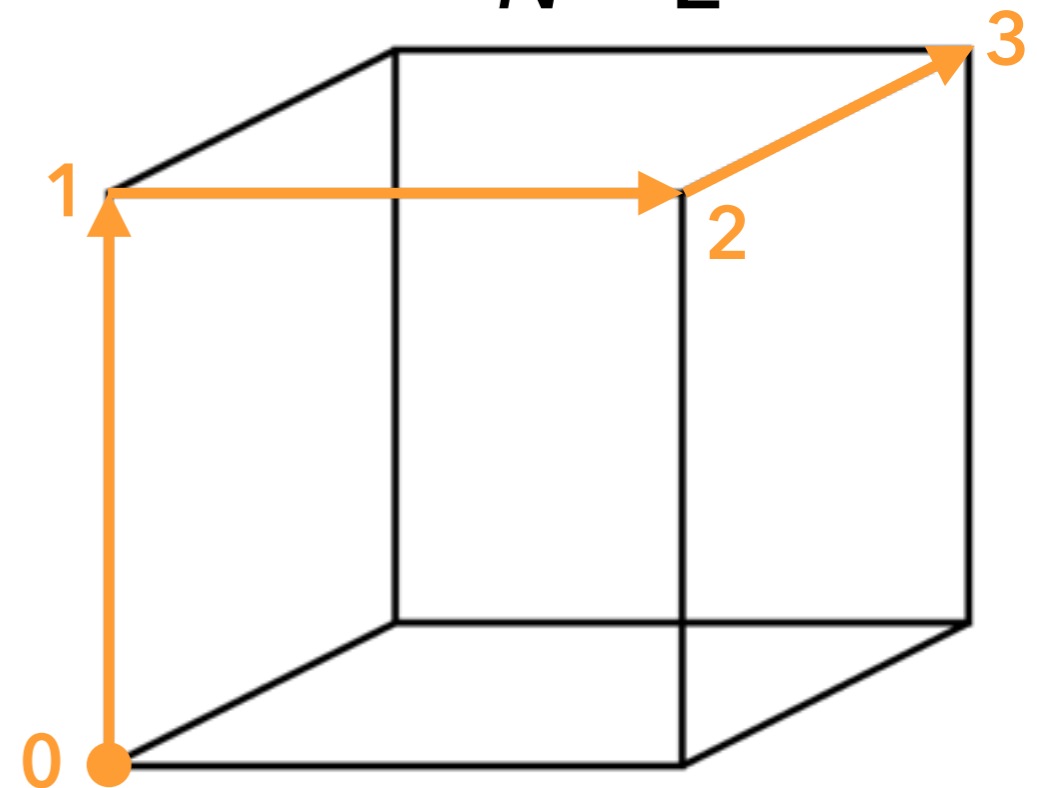
$$N = 2^d$$

▶ Any such f has a single local min? ✓

▶ #queries to find minimum?

$$O(\sqrt{N}) = O(2^{d/2})$$

LOCAL
SEARCH...



Thm [Aldous, Aaronson]: this is tight

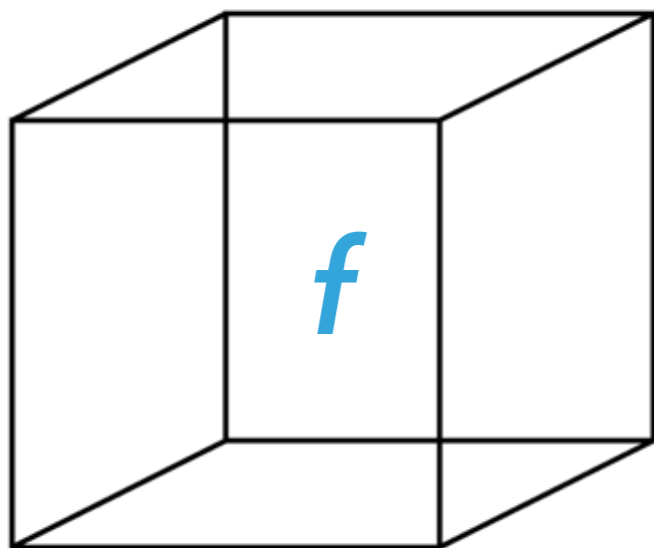
Any algo that tells whether argmin vertex is odd/even w.p. $> 2/3$

would need $\tilde{\Omega}(\sqrt{N})$ queries to f

Local search → Learning in games

Reduction: oracle access to f → VAL + OPT oracles for game

argmin f is **odd?** **even?**



$N = 2^d$

| | | | | | |
|---|---|---|---|---|-----------|
| 1 | 1 | 1 | 1 | 1 | 1 |
| 0 | 1 | 1 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 1 | 1 |
| 0 | 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1/4? 3/4? |

$N = 2^d$

$$L(i, j) = \begin{cases} 0 & \text{if } f(i) < f(j) \\ 1 & \text{otherwise} \end{cases}$$

game
equilibrium

Local search → Learning in games

Reduction: oracle access to $f \rightarrow$ VAL + OPT oracles for game

VAL oracle?

two queries to compute $L(i,j)$

OPT oracle?

$$\text{OPT}(S) = \arg \min_i \left\{ \sum_{j \in S} L(i,j) \right\}$$

▶ find i s.t.: $f(i) < \min(f(S))$

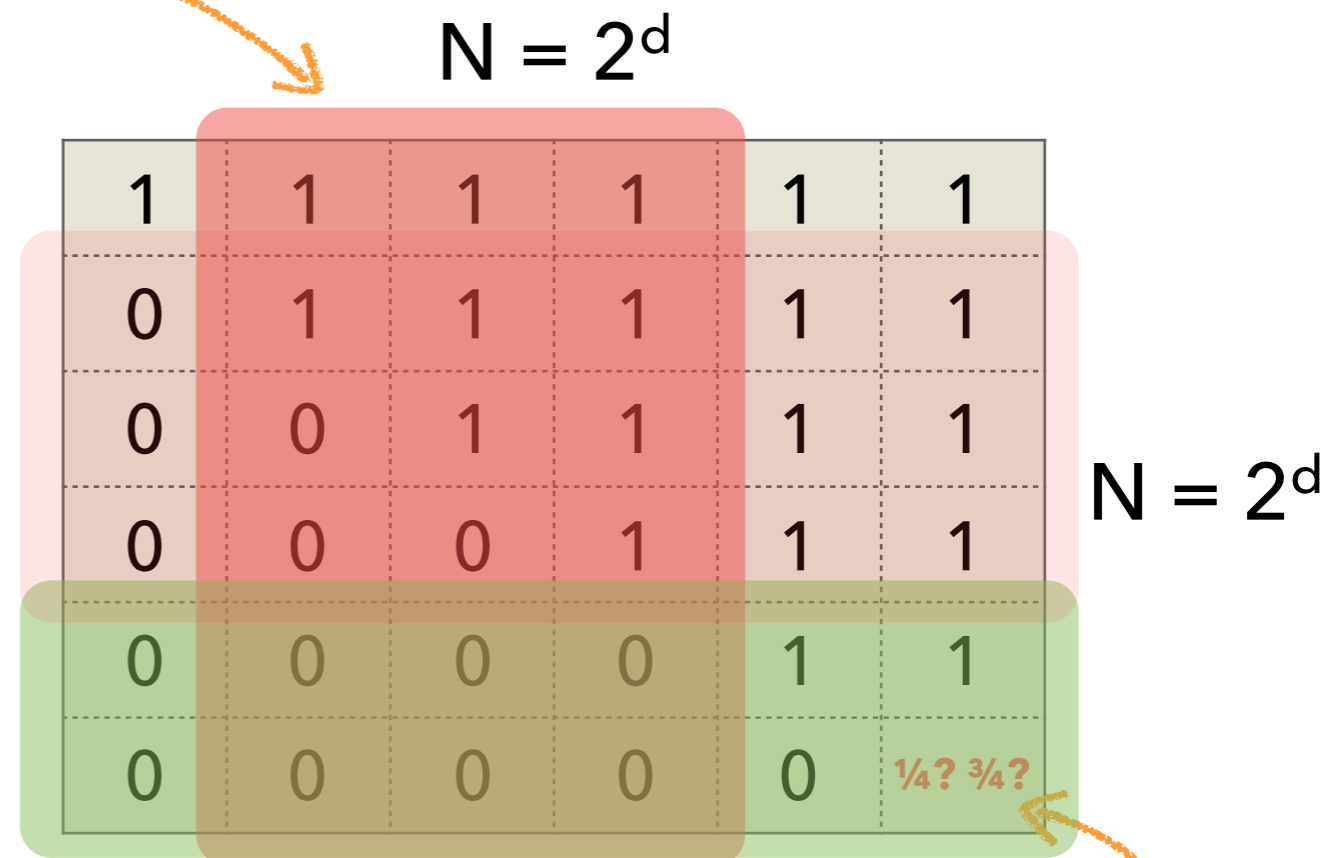
▶ local search:

$\leq \log(N)$ queries

per each $j \in S$

$$L(i,j) = \begin{cases} 0 & \text{if } f(i) < f(j) \\ 1 & \text{otherwise} \end{cases}$$

game equilibrium



Results

In OPT oracle model:

- ▶ **Thm 1.** Any algo that approximates $N \times N$ zero-sum games to within $\varepsilon = 1/4$ runs in total time $\tilde{O}(\sqrt{N})$
▶ $\tilde{O}(N)$ time needed to minimize regret
- ▶ **Thm 2.** There exists (new) online learning algo that attains regret $< \varepsilon$ in total time $\tilde{O}(\sqrt{N} / \varepsilon^2)$, tight
▶ vs. $\Theta(N/\varepsilon^2)$ time w/o OPT oracle
- ▶ **Corr.** There exists (new) algo that approximates $N \times N$ zero-sum games in total time $\tilde{O}(\sqrt{N})$, tight
▶ vs. $\Theta(N)$ time without oracles

Intuition

Idea: reduce effective #experts from N to \sqrt{N}

Interpolate two extreme cases:

1. **There are few leaders:**

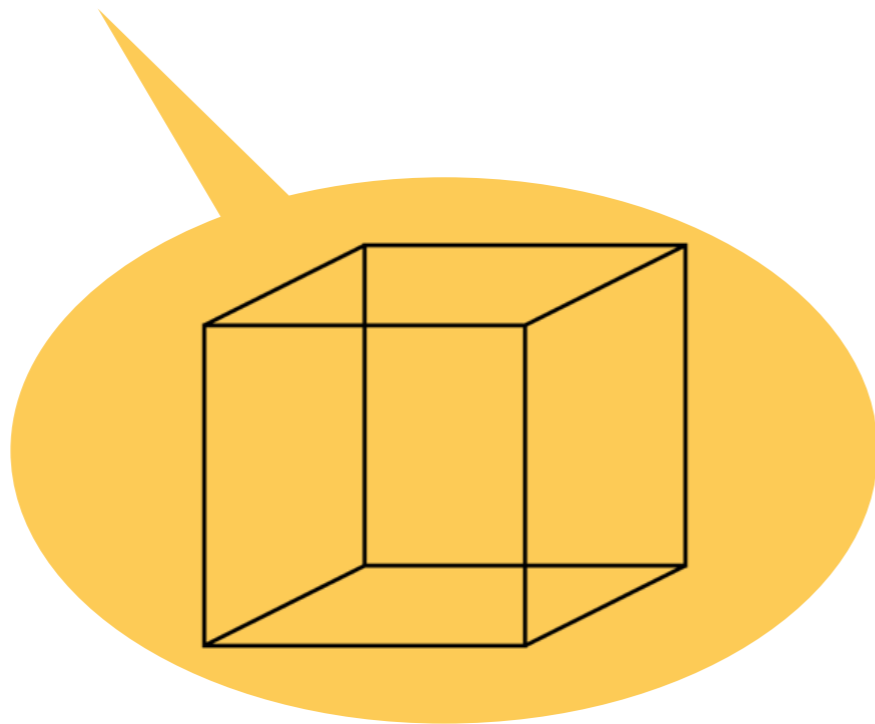
→ **OPT oracle** is useful

2. **Leader keeps changing:**

→ **sampling** $\sim\sqrt{N}$ experts will get us into \sqrt{N} "finalists"

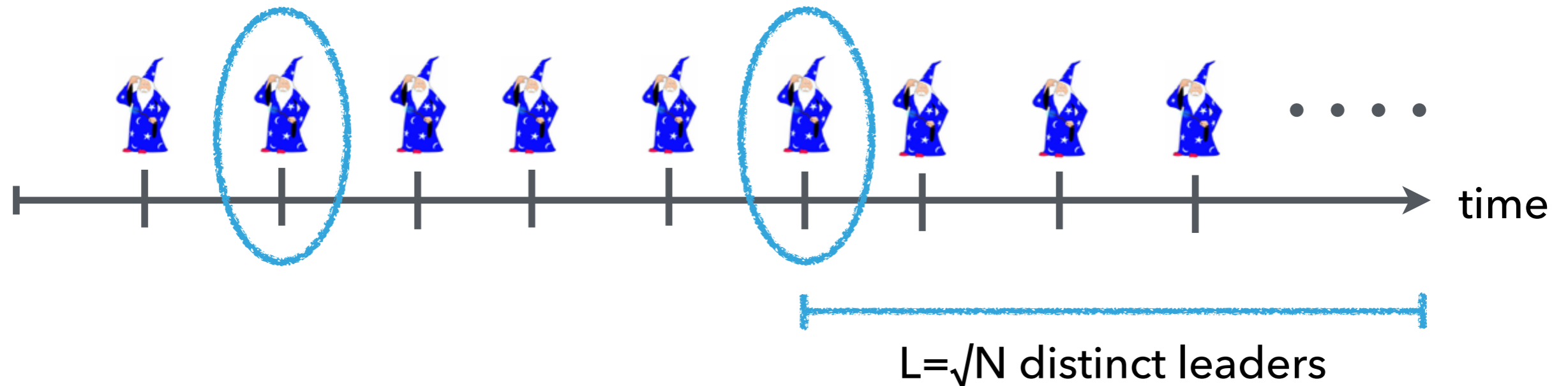


$$x_t^* = \text{OPT}(y_1, \dots, y_{t-1})$$



Stream of leaders

- ▶ Sort leaders by "death time" = last time ever to appear as leader



- ▶ **Sampling \sqrt{N} experts:**

1. w.h.p. gets us to last \sqrt{N} leaders

2. EXP3 regret vs \sqrt{N} sampled experts $\leq \sqrt{\frac{\sqrt{N}}{T}}$

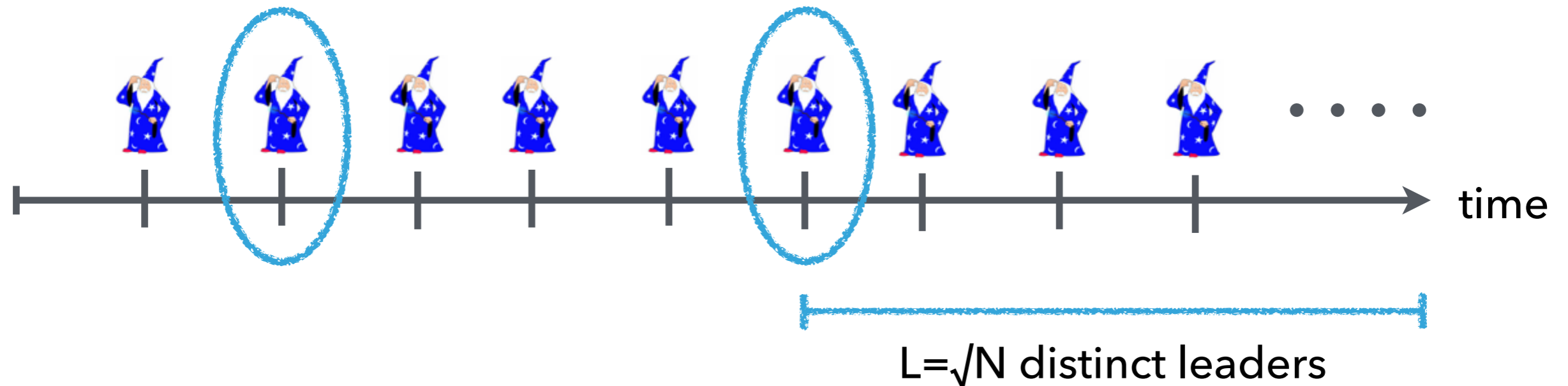
- ▶ "Only" need to get low-regret in last time interval:

$$\sqrt{\frac{L}{T}}$$

LEADERS

Stream of leaders

- ▶ Sort leaders by “death time” = last time ever to appear as leader



Combine two algorithms:

- (1) Bandit algo over random sample of \sqrt{N} experts
- (2) “Leaders” algorithm

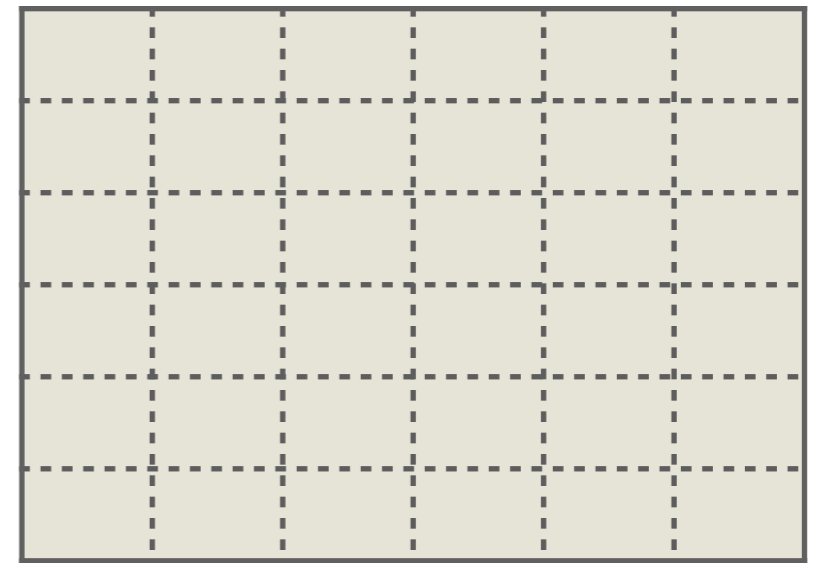
- ▶ **Leaders:** for any sequence with at most L distinct leaders:
needs $\tilde{O}(1)$ time per round

$$\frac{1}{T} \sum_{t=1}^T L_t(x_t) - \frac{1}{T} \sum_{t=1}^T L_t(x^*) \lesssim \sqrt{\frac{L}{T}}$$

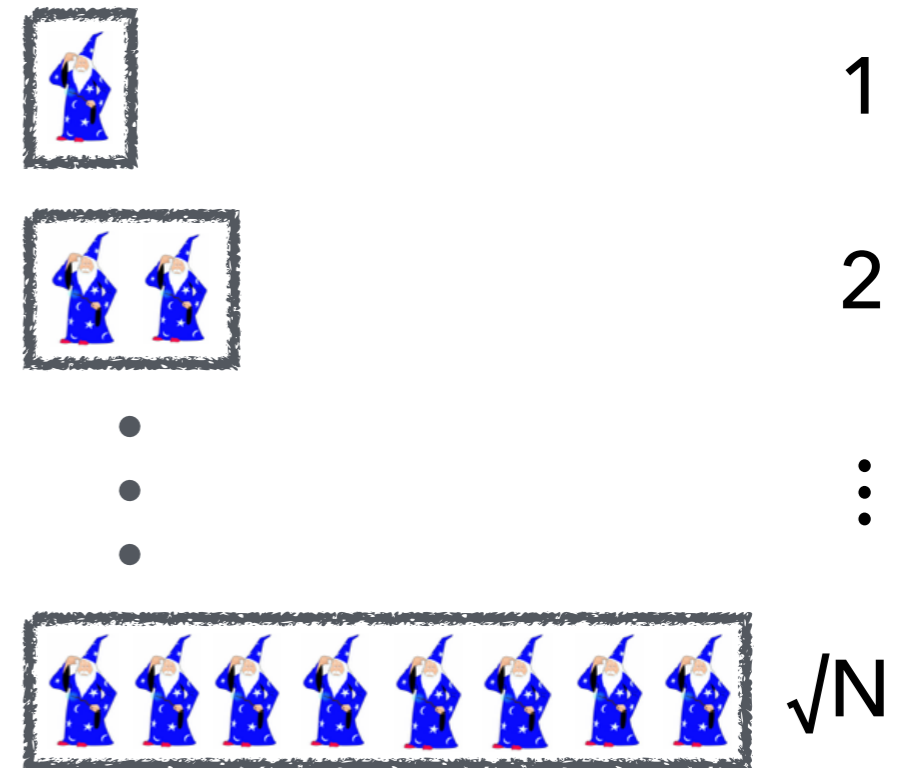
Final algorithm

bandit on
 $\sqrt{N} + \log(N)$
 "meta-arms"

$|S| = \sqrt{N}$
 random



$\log(N)$
 "sliding
 windows"



Thm. for any sequence y_1, \dots, y_T ,

w.p. $\geq 1 - \delta$:

$$\frac{1}{T} \sum_{t=1}^T L_t(x_t) - \frac{1}{T} \sum_{t=1}^T L_t(x^*) \lesssim \sqrt{\frac{\sqrt{N}}{T} \log \frac{1}{\delta}}$$

Bottom line

- ▶ efficient OPT \Rightarrow efficient online learning
- ▶ but it helps, **quadratically**
- ▶ intriguing connections to **local search**

Many questions:

- ▶ stronger positive results? what assumptions?
(e.g., [[Daskalakis-Syrgkanis '16](#)])
- ▶ what about oracle complexity? (lower b.)
- ▶ approximate optimization? (upper b.)
- ▶ ...