# Learning Stronger Notions of Equilibrium

Brian Zhang

# Recap: CCEs in Normal-Form Games

$X_i$ = set of **pure** strategies of player $i$

Correlated strategy profile:

*Note: not* $\Delta(X_1) \times \cdots \times \Delta(X_n)$

$$\bar{\mu}^T := \frac{1}{T} \sum_{t=1}^{T} (\mu_1^t \otimes \mu_2^t \otimes \cdots \mu_n^t) \in \Delta(X_1 \times \cdots \times X_n)$$

*the product distribution in* $\Delta(X_1) \times \cdots \times \Delta(X_n)$
*whose marginal on* $X_i$ *is* $\mu_i^t \in \Delta(X_i)$

Regret guarantee: for all players $i$:

$$\max_{x_i^*} \frac{1}{T} \sum_{t=1}^{T} \left[ u_i(x_i^*, x_{-i}^t) - u_i(x_i^t, x_{-i}^t) \right] \leq O_n\left(\frac{1}{\sqrt{T}}\right)$$

$$= \max_{x_i^*} \mathbb{E}_{x \sim \bar{\mu}^T} \left[ u_i(x_i^*, x_{-i}) - u_i(x_i, x_{-i}) \right]$$

$\bar{\mu}^T$ is an $\epsilon$-*"coarse-correlated equilibrium" (CCE)* where $\epsilon = O_n\left(1/\sqrt{T}\right)$

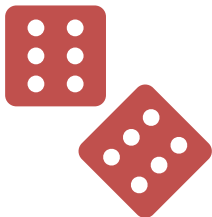Works for extensive-form games too: use CFR!

# Coarse-Correlated Equilibria

**Def:** $\mu \in \Delta(X_1 \times \cdots \times X_n)$ is a coarse-correlated equilibrium (CCE) if

$$\mathbb{E}_{x \sim \mu} \left[ u_i(x_i^*, x_{-i}) - u_i(x_i, x_{-i}) \right] \leq 0$$

for all players $i$ and all strategies $x_i^* \in X_i$

"Correlation device"
"Mediator"

I will sample $x \sim \mu$. You can either **commit to playing the strategy I sample**, or **play a strategy of your choice**

**CCE:**

I will **commit to playing your sampled strategy**, whatever it is.

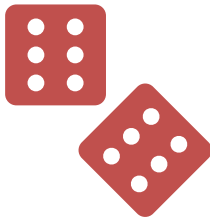Player $i$

# Coarse-Correlated Equilibria

**Def:** $\mu \in \Delta(X_1 \times \cdots \times X_n)$ is a coarse-correlated equilibrium (CCE) if

$$\mathop{\mathbb{E}}_{x \sim \mu} \left[ u_i(x_i^*, x_{-i}) - u_i(x_i, x_{-i}) \right] \leq 0$$

for all players $i$ and all strategies $x_i^* \in X_i$

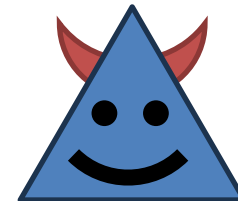"Correlation device"
"Mediator"

I will sample $x \sim \mu$. You can either **commit to playing the strategy I sample**, or **play a strategy of your choice**

**Not CCE:**

I think $x_i^*$ is a unilaterally profitable deviation, and I'll play that instead

Player $i$

Fairly **weak notion**: Player must commit **before seeing the sampled strategy**
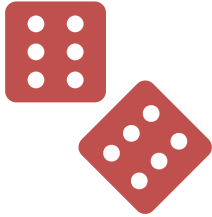e.g., CCEs can include **dominated strategies** (HW1)

# Correlated Equilibria

**Def:** $\mu \in \Delta(X_1 \times \cdots \times X_n)$ is a correlated equilibrium (CE) if

$$\mathop{\mathbb{E}}_{x \sim \mu} \left[ u_i(\phi_i(x_i), x_{-i}) - u_i(x_i, x_{-i}) \right] \leq 0$$

for all players $i$ and all functions $\phi_i : X_i \to X_i$
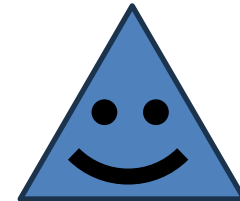
"Correlation device"
"Mediator"

I will sample $x \sim \mu$, and **tell you** $x_i$. Then you can choose what action you want to play.

$x_i$

Player $i$

**CE:**

Okay, I will play $x_i$

# Correlated Equilibria

**Def:** $\mu \in \Delta(X_1 \times \cdots \times X_n)$ is a correlated equilibrium (CE) if

$$\mathop{\mathbb{E}}_{x \sim \mu} \left[ u_i(\phi_i(x_i), x_{-i}) - u_i(x_i, x_{-i}) \right] \leq 0$$

for all players $i$ and all functions $\phi_i : X_i \to X_i$

"Correlation device"
"Mediator"

I will sample $x \sim \mu$, and **tell you** $x_i$. Then you can choose what action you want to play.

$x_i$

Player $i$

**Not CE:**

Given your recommendation $x_i$, I think $x'_i := \phi_i(x_i)$ is a better action, so I'll play that instead.

# Correlated Equilibria
# in Normal-Form Games

**Chicken**

| | Stop | Go |
|---|---|---|
| **Stop** | | |
| **Go** | | |

# Correlated Equilibria
# in Normal-Form Games

**Chicken**

|  | Stop | Go |
|---|---|---|
| **Stop** |  |  |
| **Go** | -5, -5 |  |

# Correlated Equilibria
# in Normal-Form Games

**Chicken**

|  | Stop | Go |
|---|---|---|
| **Stop** |  | 0, 1 |
| **Go** |  | -5, -5 |

# Correlated Equilibria
# in Normal-Form Games

**Chicken**

| | Stop | Go |
|---|---|---|
| **Stop** | | 0, 1 |
| **Go** | 1, 0 | -5, -5 |

# Correlated Equilibria in Normal-Form Games

**Chicken**

|  | Stop | Go |
|---|---|---|
| **Stop** | 0, 0 | 0, 1 |
| **Go** | 1, 0 | -5, -5 |

# Correlated Equilibria in Normal-Form Games

**Chicken**

|  | Stop | Go |
|---|---|---|
| **Stop** | 0, 0 <br> 0 | 0, 1 <br> p |
| **Go** | 1, 0 <br> 1-p | -5, -5 <br> 0 |

# Correlated Equilibria
# in Normal-Form Games

**Chicken**

|  | Stop | Go |
|---|---|---|
| **Stop** | 0, 0 <br> 0 | 0, 1 <br> p |
| **Go** | 1, 0 <br> 1-p | -5, -5 <br> 0 |

$$\mu = \frac{1}{2}(\text{Stop}, \text{Go}) + \frac{1}{2}(\text{Go}, \text{Stop})$$

is a CE (and a CCE)

# CCEs can be learned using any no-regret algorithm.

## Question: Can CEs?

# Normal-Form Strategy Maps

A map $\phi : X \rightarrow X$, where $X := \{\boldsymbol{e}_1, \ldots, \boldsymbol{e}_n\} \subset \mathbb{R}^n$, is given by a matrix $\boldsymbol{M} \in \mathbb{R}^{n \times n}$ whose $i$th column specifies $\phi(\boldsymbol{e}_i) \in X$.

e.g.,

$$\boldsymbol{M} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\phi(\boldsymbol{x}) = \boldsymbol{M}\boldsymbol{x}$$

# Normal-Form Strategy Maps

A randomized map $\phi : X \to \mathrm{conv}(X)$, where $X :=$ $\{\boldsymbol{e}_1, \dots, \boldsymbol{e}_n\} \subset \mathbb{R}^n$, is given by a matrix $\boldsymbol{M} \in \mathbb{R}^{n \times n}$ whose $i$th column specifies $\phi(\boldsymbol{e}_i) \in \mathrm{conv}(X)$.

e.g.,

$$\boldsymbol{M} = \begin{bmatrix} 0.7 & 1 & 0.2 \\ 0.3 & 0 & 0.6 \\ 0 & 0 & 0.2 \end{bmatrix}$$

$$\phi(\boldsymbol{x}) = \boldsymbol{M}\boldsymbol{x}$$

# No-(External-)Regret Learning

Pure strategy set $X := \{\boldsymbol{e}_1, \dots, \boldsymbol{e}_n\} \subset \mathbb{R}^n$

On each iteration:

- player outputs **mixed strategy** $\boldsymbol{x}^t \in \text{conv}(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in [-1,1]^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize **regret** after $T$ timesteps

$$R_X(T) := \max_{\boldsymbol{x}^* \in X} \sum_{t=1}^{T} \langle \boldsymbol{u}^t, \boldsymbol{x}^* - \boldsymbol{x}^t \rangle$$

# No-Swap-Regret Learning

Pure strategy set $X := \{e_1, \ldots, e_n\} \subset \mathbb{R}^n$

On each iteration:

- player outputs **mixed strategy** $x^t \in \text{conv}(X)$
- environment outputs (possibly adversarial) **utility vector** $u^t \in [-1,1]^n$
- player observes $u^t$ and gets reward $\langle u^t, x^t \rangle \in [-1,1]$

Goal: minimize **swap regret** after $T$ timesteps

$$R_X^{\text{Swap}}(T) := \max_{M \in S_n} \sum_{t=1}^{T} \langle u^t, Mx^t - x^t \rangle$$

$S_n$ = set of $n \times n$ stochastic matrices

**Proposition:**
If all players in a game achieve swap regret $\epsilon T$, then the average strategy profile $\bar{\mu}$ is an $\epsilon$-correlated equilibrium.

# The GGM Framework

**Idea:** Use
- a regret minimizer $\mathcal{R}_\Phi$ on $S_n$ (stochastic matrices) with regret $R_\Phi(T)$, and
- fixed points

we'll discuss how to do this in a minute

**Algorithm:** For each iteration $t = 1, \dots, T$:
1. Obtain matrix $\boldsymbol{M}^t$ from $\mathcal{R}_\Phi$
2. Compute $\boldsymbol{x}^t \in \text{conv}(X)$ such that $\boldsymbol{M}^t \boldsymbol{x}^t = \boldsymbol{x}^t$
3. Play $\boldsymbol{x}^t$, observe utility $\boldsymbol{u}^t$
4. Feed to $\mathcal{R}_\Phi$ the utility $\boldsymbol{M} \mapsto \langle \boldsymbol{u}^t, \boldsymbol{M}\boldsymbol{x}^t \rangle$

**Regret analysis:** $\quad R_X^{\text{Swap}}(T) = \max_{\boldsymbol{M} \in S_n} \sum_{t=1}^{T} \langle \boldsymbol{u}^t, \boldsymbol{M}\boldsymbol{x}^t - \boldsymbol{x}^t \rangle$

# The GGM Framework

Blum, Mansour (*JMLR* 2007); Gordon, Greenwald, Marks (*ICML* 2008)

**Idea:** Use
- a regret minimizer $\mathcal{R}_\Phi$ on $S_n$ (stochastic matrices) with regret $R_\Phi(T)$, and
- fixed points

we'll discuss how to
do this in a minute

**Algorithm:** For each iteration $t = 1, \dots, T$:

1. Obtain matrix $\boldsymbol{M}^t$ from $\mathcal{R}_\Phi$
2. Compute $\boldsymbol{x}^t \in \mathrm{conv}(X)$ such that $\boldsymbol{M}^t \boldsymbol{x}^t = \boldsymbol{x}^t$
3. Play $\boldsymbol{x}^t$, observe utility $\boldsymbol{u}^t$
4. Feed to $\mathcal{R}_\Phi$ the utility $\boldsymbol{M} \mapsto \langle \boldsymbol{u}^t, \boldsymbol{M}\boldsymbol{x}^t \rangle$

**Regret analysis:** $\quad R_X^{\mathrm{Swap}}(T) = \max_{\boldsymbol{M} \in S_n} \sum_{t=1}^{T} \langle \boldsymbol{u}^t, \boldsymbol{M}\boldsymbol{x}^t - \boldsymbol{M}^t \boldsymbol{x}^t \rangle = R_\Phi(T)$

# Regret Minimization Over $n \times n$ Stochastic Matrices



{ Sequence-form strategies in this tree-form decision problem }

$\cong$

{ 4×4 stochastic matrices }

**Use CFR!**

$$R_X^{\text{Swap}}(T) = R_\Phi(T) \in \mathcal{O}\left(n\sqrt{T \log n}\right)$$

*with MWU at every decision point*

Tighter analysis is possible: Blum-Mansour shows $\sqrt{Tn \log n}$

**Theorem** [Blum & Mansour *JMLR* 2007]
There exists an algorithm for learning CE in normal-form games with convergence rate $\sqrt{(n \log n)/T}$ .

# More Generally: $\Phi$-Equilibria

**Def:** Given a tuple of subsets $\Phi = \{\Phi_i\}_{i \in [n]}$ where $\Phi_i \subseteq X_i^{X_i}$, correlated distribution $\mu \in \Delta(X_1 \times \cdots \times X_n)$ is a $\Phi$-equilibrium if

$$\mathbb{E}_{x \sim \mu} \left[ u_i(\phi_i(x_i), x_{-i}) - u_i(x_i, x_{-i}) \right] \le 0$$

for all players $i$ and all functions $\phi_i \in \Phi_i$

**Special cases:**

- CCE (constant functions): $\Phi_i = \{\phi_{x_i^*} : x^* \in X_i\}$ where $\phi_{x_i^*}(x_i) = x_i^*$ for all $x_i$

- CE (all functions): $\Phi_i = X_i^{X_i}$

# No-(External-)Regret Learning in Extensive-Form Games

Pure strategy set $X \subseteq \{0,1\}^n$

On each iteration:

- player outputs **tree-form strategy** $\boldsymbol{x}^t \in \mathrm{conv}(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in \mathbb{R}^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize **regret** after $T$ timesteps

$$R_X(T) := \max_{\boldsymbol{x}^* \in X} \sum_{t=1}^{T} \langle \boldsymbol{u}^t, \boldsymbol{x}^* - \boldsymbol{x}^t \rangle$$

# No-(External-)Regret Learning in Extensive-Form Games

Pure strategy set $X \subseteq \{0,1\}^n$

On each iteration:

- player outputs **mixed strategy** $\mu^t \in \Delta(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in \mathbb{R}^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\underset{\boldsymbol{x}^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize **regret** after $T$ timesteps

$$R_X(T) := \max_{\boldsymbol{x}^* \in X} \sum_{t=1}^{T} \underset{\boldsymbol{x}^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \boldsymbol{x}^* - \boldsymbol{x}^t \rangle$$

# No-Φ-Regret Learning

Pure strategy set $X \subseteq \{0,1\}^n$, set of deviations $\Phi \subseteq X^X$

On each iteration:

- player outputs **mixed strategy** $\mu^t \in \Delta(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in \mathbb{R}^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\underset{x^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize Φ**-regret** after $T$ timesteps

$$R_X^{\Phi}(T) := \max_{\phi \in \Phi} \sum_{t=1}^{T} \underset{x^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \phi(\boldsymbol{x}^t) - \boldsymbol{x}^t \rangle$$

| Φ | Notion of Regret | Corresponding Notion of Equilibrium |
|---|---|---|
| $\Phi_{\text{Ext}} = \{\text{constant functions}\}$ | External | Coarse-Correlated |
| $\Phi_{\text{Swap}} = X^X$ (all functions) | Swap | Correlated |

# No-Φ-Regret Learning

Pure strategy set $X \subseteq \{0,1\}^n$, set of deviations $\Phi \subseteq X^X$

On each iteration:

- player outputs **mixed strategy** $\mu^t \in \Delta(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in \mathbb{R}^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\underset{x^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize Φ-**regret** after $T$ timesteps

$$R_X^\Phi(T) := \max_{\phi \in \Phi} \sum_{t=1}^{T} \underset{x^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \phi(\boldsymbol{x}^t) - \boldsymbol{x}^t \rangle$$

**Proposition**
If all players in a game run Φ-regret minimizers that achieve Φ-regret $\epsilon T$, then the average strategy profile $\bar{\mu}$ is an $\epsilon$-approximate Φ-equilibrium.

# Swap Regret in Extensive-Form Games

**Q:** Can **swap regret** be efficiently minimized in *extensive-form* games?

**Theorem**
[Corollary of Blum-Mansour]

There exists a swap regret minimizer for tree-form strategy sets whose swap regret is $\epsilon T$ after $\mathcal{O}(n \cdot 2^n/\epsilon^2)$ **iterations.**

Bad per-iteration complexity and convergence rate

**Theorem**
[*Special case of* Peng & Rubinstein *STOC'24*; Dagan, Daskalakis, Fishelson, Golowich *STOC'24*]

There exists a swap regret minimizer for tree-form strategy sets* whose swap regret is $\epsilon T$ after $n^{\widetilde{\mathcal{O}}(1/\epsilon)}$ **iterations.**

*or, indeed, any set $X \subset \mathbb{R}^n$ for which *external* regret is minimizable

$\Rightarrow$ For **constant** $\epsilon$, an $\epsilon$-CE can be computed in **polynomial time!**

**Theorem**
[Daskalakis, Farina, Golowich, Sandholm, Zhang *arXiv'24*]

There is a constant $c > 0$ such that achieving swap regret $\epsilon T$ in tree-form strategy sets requires $\exp(\Omega(\min\{n, 1/\epsilon\}^c))$ **iterations.**

**Open question:** Can $\epsilon$-CE be computed in time $\text{poly}(n, 1/\epsilon)$ or even $\text{poly}(n, \log(1/\epsilon))$?

(using something other than adversarial no-swap-regret learning)

# Digression: Nonlinear strategy maps

Pure strategy set $X \subseteq \{0,1\}^n$ , set of deviations $\Phi \subseteq X^X$

External regret minimizer on $X$ outputs points in $\text{conv}(X)$

**Q:** For $\boldsymbol{x}^* \in \text{conv}(X)$ and $\phi : X \to X$, what does $\phi(\boldsymbol{x}^*)$ mean?

**A1:** When $X = \{\boldsymbol{e}_1, \dots, \boldsymbol{e}_n\}$ is a normal-form strategy set, $\text{conv}(X) = \Delta(X)$ and $\phi(\boldsymbol{x}) = \boldsymbol{Mx}$ for some $\boldsymbol{M}$, so we can set $\phi(\boldsymbol{x}^*) = \sum_i \boldsymbol{x}_i^* \phi(\boldsymbol{e}_i) = \boldsymbol{Mx}^*$.

**A2:** Take **any** distribution $\mu \in \Delta(X)$ with $\boldsymbol{x}^* = \mathop{\mathbb{E}}_{\boldsymbol{x} \sim \mu} \boldsymbol{x}$, and define
$$\phi(\boldsymbol{x}^*) = \mathop{\mathbb{E}}_{\boldsymbol{x} \sim \mu} \phi(\boldsymbol{x}).$$

**Warning: When $\phi$ is nonlinear, this depends on the choice of $\mu$**

⇒ "Kuhn's theorem fails when considering nonlinear deviations"

**A3:** When $\Phi$ consists only of linear maps, this doesn't matter (we can use sequence-form strategies + set $\phi(\boldsymbol{x}) = \boldsymbol{Mx}$

# No-Linear-Swap-Regret Learning

Pure strategy set $X \subseteq \{0,1\}^n$,

On each iteration:

- player outputs **mixed strategy** $\mu^t \in \Delta(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in \mathbb{R}^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\mathbb{E}_{x^t \sim \mu^t} \langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize $\Phi$-**regret** after $T$ timesteps

$$R_X^{\Phi}(T) := \max_{\boldsymbol{M} \in \Phi_{\text{LIN}}} \sum_{t=1}^{T} \mathbb{E}_{x^t \sim \mu^t} \langle \boldsymbol{u}^t, \boldsymbol{M}\boldsymbol{x}^t - \boldsymbol{x}^t \rangle$$

$$\Phi_{\text{LIN}} = \{\boldsymbol{M} : \boldsymbol{M}\boldsymbol{x} \in \text{conv}(X) \ \ \forall \boldsymbol{x} \in \text{conv}(X)\}$$

Advantages:
- Natural generalization of stochastic matrices for normal-form games
- GGM applies verbatim, and fixed points are easy (linear program: $\boldsymbol{M}\boldsymbol{x} = \boldsymbol{x}, \ \boldsymbol{x} \in \text{conv}(X)$)

# No-Linear-Swap-Regret Learning

Pure strategy set $X \subseteq \{0,1\}^n$,

On each iteration:

- player outputs **tree-form strategy** $x^t \in \text{conv}(X)$
- environment outputs (possibly adversarial) **utility vector** $u^t \in \mathbb{R}^n$
- player observes $u^t$ and gets reward $\langle u^t, x^t \rangle \in [-1,1]$

Goal: minimize $\Phi$-**regret** after $T$ timesteps

$$R_X^\Phi(T) := \max_{M \in \Phi_{\text{LIN}}} \sum_{t=1}^{T} \langle u^t, Mx^t - x^t \rangle$$

$$\Phi_{\text{LIN}} = \{M : Mx \in \text{conv}(X) \quad \forall x \in \text{conv}(X)\}$$

Advantages:
- Natural generalization of stochastic matrices for normal-form games
- GGM applies verbatim, and fixed points are easy (linear program: $Mx = x, \; x \in \text{conv}(X)$)
- We can still work with tree-form strategies (linearity of expectation)

# The GGM Framework

Gordon, Greenwald, Marks (*ICML* 2008)

## GGM requires two things.

- fixed point oracle fix : $\Phi_{\text{LIN}} \to \text{conv}(X)$, *i.e.*, $\boldsymbol{Mx} = \boldsymbol{x}$ if $\boldsymbol{x} = \text{fix}(\boldsymbol{M})$, and
  *Still easy! Use linear programming or power iteration*

- a regret minimizer $\mathcal{R}_\Phi$ on $\Phi_{\text{LIN}}$
  *How to characterize $\Phi_{\text{LIN}}$?*

# So what does $\Phi_{\mathrm{LIN}}$ look like?

Warm-up (Special case): What are the affine maps
$$\phi : [0,1]^n \to [0,1]?$$

- Constant functions:
$$\phi(\boldsymbol{x}) = 0, \qquad \phi(\boldsymbol{x}) = 1$$
- Functions that depend on one input coordinate:
$$\phi(\boldsymbol{x}) = x_i, \qquad \phi(\boldsymbol{x}) = 1 - x_i$$

**Claim:** Every affine $\phi : [0,1]^n \to [0,1]$ is a convex combination of these!

# So what does $\Phi_{\mathrm{LIN}}$ look like?
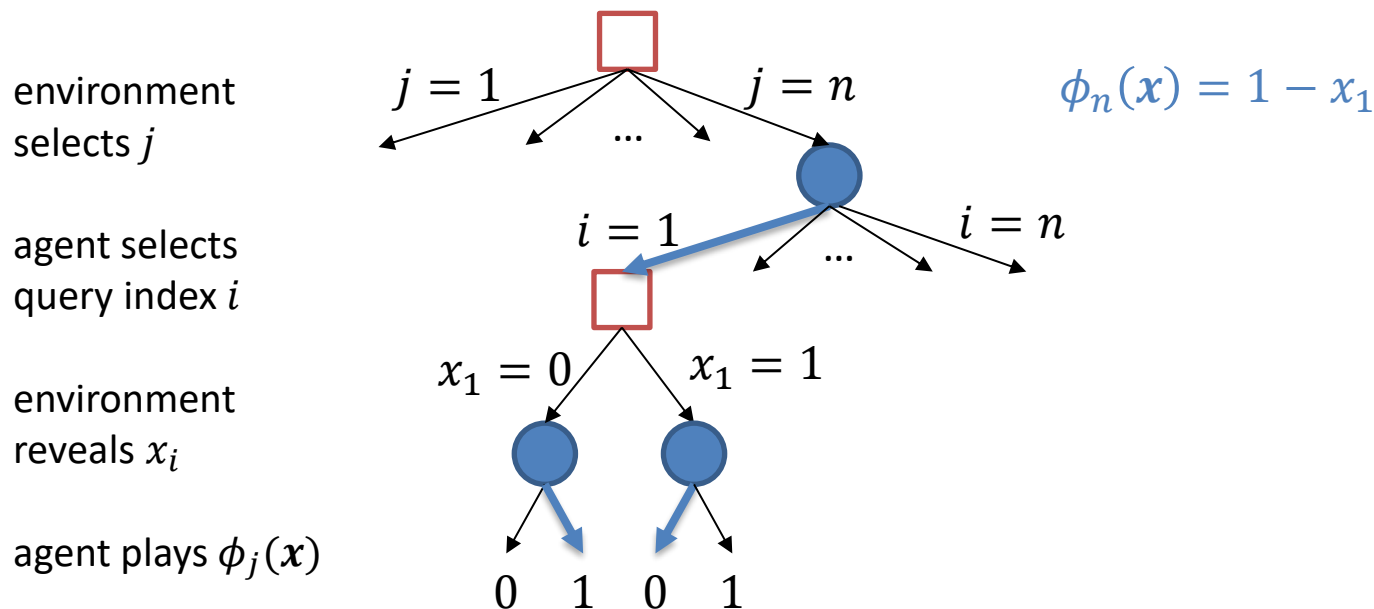
Warm-up (Special case): What are the affine maps
$$\phi : [0,1]^n \to [0,1]^{n}?$$

Each coordinate $j$ is an affine map $\phi_j : [0,1]^n \to [0,1]$
$\Rightarrow$ Each $\phi_j$ makes $\leq 1$ query to the input

environment selects $j$

agent selects query index $i$

environment reveals $x_i$

agent plays $\phi_j(\boldsymbol{x})$

$j = 1$     $j = n$

$\phi_n(\boldsymbol{x}) = 1 - x_1$

$i = 1$     $i = n$

$x_1 = 0$     $x_1 = 1$

0  1  0  1

# So what does $\Phi_{\mathrm{LIN}}$ look like?

**Insight:**

Affine maps
$\phi : [0,1]^n \rightarrow [0,1]^n$
$\equiv$
Tree-form strategies
with **one query**

environment
selects $j$

agent selects
query index $i$

environment
reveals $x_i$

agent plays $\phi_j(\boldsymbol{x})$

$j = 1$  $\quad$ $j = n$

$\cdots$

$i = 1$  $\quad$ $i = n$

$\cdots$

$x_1 = 0$  $\quad$ $x_1 = 1$

$0 \quad 1 \quad 0 \quad 1$

$\phi_n(\boldsymbol{x}) = 1 - x_1$

# Does this generalize?

What is the generalization of a "query" to an arbitrary tree-form strategy space?

**Mediator (holds $x$)**

**Real game (play $Mx$)**
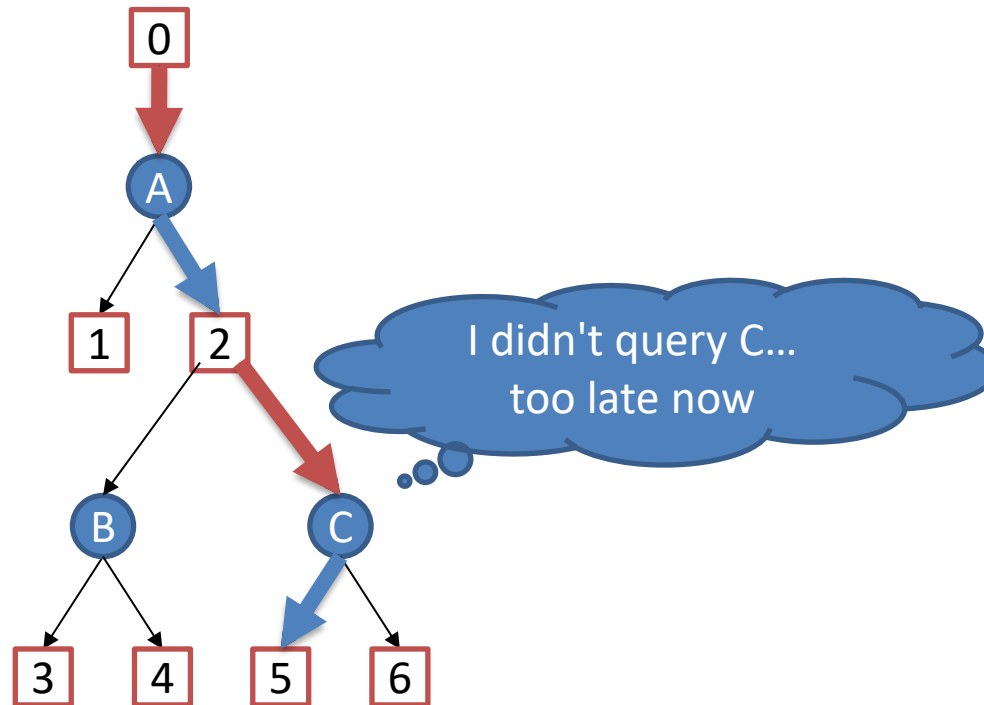
I didn't query C...
too late now
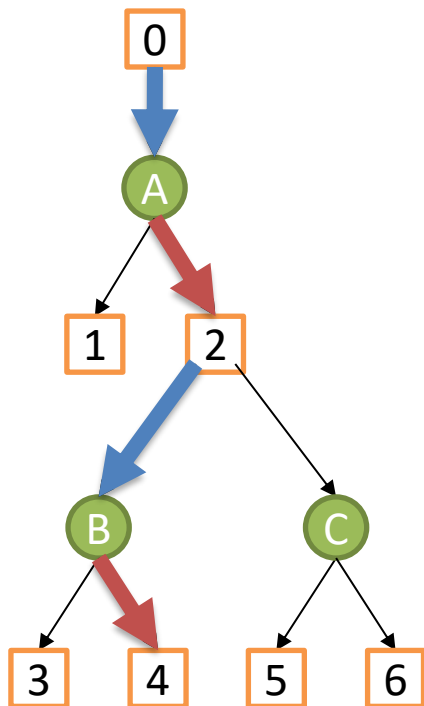
# Does this generalize?

What is the generalization of a "query" to an arbitrary tree-form strategy space?

**Mediator (holds $x$)**



**Real game (play $Mx$)**



These are the
**untimed communication
(UTC) deviations**

**Communication:** Player has two-way communication with mediator to gain information

**Untimed:** Player can send zero, one, or multiple queries between real game actions

# Untimed communication deviations as tree-form decision problems



**Mediator (holds $x$)**

**Real game (play $Mx$)**

# Untimed communication deviations as ~~tree~~-form decision problems
## DAG

Strategy in DAG $\Rightarrow$ function $\phi : X \rightarrow X$

Size of DAG: $O(n^2)$

**Mediator (holds $x$)**

**Real game (play $Mx$)**

# Untimed communication deviations as ~~tree~~-form decision problems

## DAG

$$\phi(\boldsymbol{x})[\sigma] = \sum_{\sigma'} M[\sigma, \sigma'] \; x[\sigma']$$

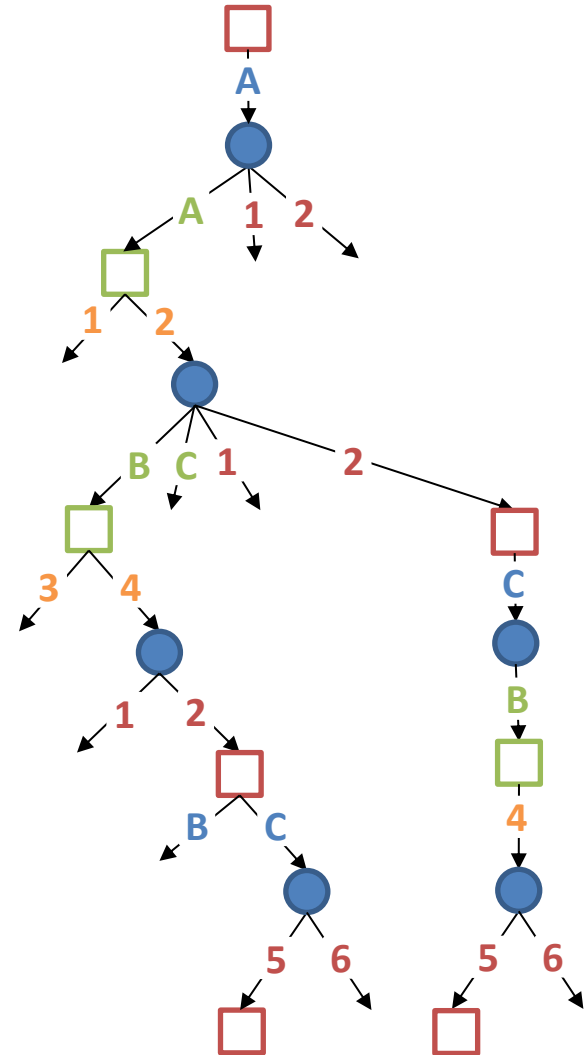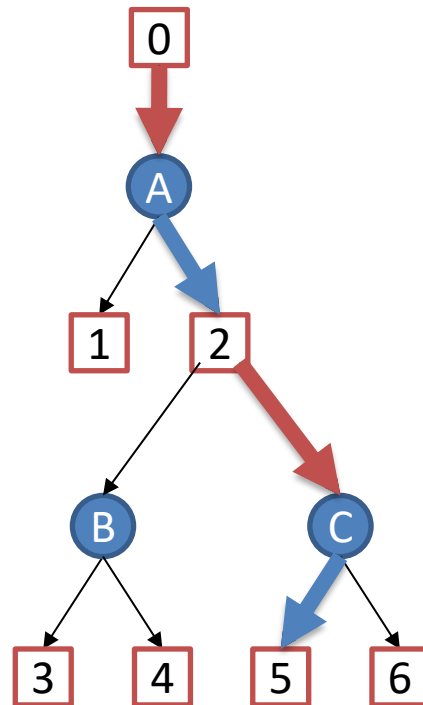$\underbrace{\phantom{\phi(\boldsymbol{x})[\sigma]}}$
$\phi(x)$ plays to $\sigma$

$\underbrace{\phantom{M[\sigma,\sigma']}}$
$\phi(x)$ plays to $\sigma$
if $x$ plays to $\sigma'$

$\underbrace{\phantom{x[\sigma']}}$
$x$ plays to $\sigma'$

$$\phi(\boldsymbol{x}) = \boldsymbol{Mx}$$

$$\Rightarrow \Phi_{\mathrm{UTC}} \subseteq \Phi_{\mathrm{LIN}}$$

**THEOREM**
[Zhang, Farina, Sandholm ICLR'24]
$\Phi_{\mathrm{UTC}} = \Phi_{\mathrm{LIN}}.$

Size of DAG:
$O(n^2)$

The UTC functions are
exactly the linear functions
[Zhang, Farina, Sandholm *ICLR*'24]

**+**

Regret minimization on DAGs of size $m = n^2$
is possible with regret $m\sqrt{T}$ using CFR + scaled extensions
[Zhang, Farina, Sandholm *ICML*'23]

**+**

Fixed-point solving using LP or power iteration

GGM

**COROLLARY**
[Zhang, Farina, Sandholm *ICLR*'24]
$\Phi_{\text{LIN}}$-regret minimization on tree-form decision
problems is possible with regret $n^2\sqrt{T}$

# Beyond Linear Deviations

Pure strategy set $X \subseteq \{0,1\}^n$, set of deviations $\Phi \subseteq X^X$

On each iteration:

- player outputs **mixed strategy** $\mu^t \in \Delta(X)$
- environment outputs (possibly adversarial) **utility vector** $\boldsymbol{u}^t \in \mathbb{R}^n$
- player observes $\boldsymbol{u}^t$ and gets reward $\underset{\boldsymbol{x}^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \boldsymbol{x}^t \rangle \in [-1,1]$

Goal: minimize $\Phi$-**regret** after $T$ timesteps

$$R_X^{\Phi}(T) := \max_{\phi \in \Phi} \sum_{t=1}^{T} \underset{\boldsymbol{x}^t \sim \mu^t}{\mathbb{E}} \langle \boldsymbol{u}^t, \phi(\boldsymbol{x}^t) - \boldsymbol{x}^t \rangle$$

# The GGM Framework

Gordon, Greenwald, Marks (*ICML* 2008)

Pure strategy set $X \subseteq \{0,1\}^n$ , set of deviations $\Phi \subseteq X^X$

**GGM requires two things.**

- Fixed point oracle $\mathrm{fix} : \Phi \to \mathrm{conv}(X)$, *i.e.,* $\phi(\boldsymbol{x}) = \boldsymbol{x}$ if $\boldsymbol{x} = \mathrm{fix}(\phi)$

  **Problem:** *$\phi : X \to X$ is a discrete function!*

  - *It may not have a fixed point*

  - *Even if we make some assumption like $\phi$ being continuous, fixed points are PPAD-hard to compute*

- Regret minimizer $\mathcal{R}_\Phi$ on $\Phi$

  **Problem:** *if $X = \{0,1\}^n$ then $|\Phi| > 2^{n \cdot 2^n}$. How can we hope to minimize regret efficiently?*

# The GGM Framework: Upgraded

Zhang, Anagnostides, Farina, Sandholm (arXiv 2024)

Pure strategy set $X \subseteq \{0,1\}^n$ , set of deviations $\Phi \subseteq X^X$

**GGM requires two things.**

- **Expected** fixed point oracle fix $: \Phi \to \Delta(X)$, *i.e.,* $\displaystyle\mathop{\mathbb{E}}_{x \sim \mu} x = \mathop{\mathbb{E}}_{x \sim \mu} \phi(x)$ if $\mu = \text{fix}(\phi)$

  - Always exist

  - Easy to compute! $\mu := \text{Unif}\{x, \phi(x), \phi^2(x), \dots, \phi^{L-1}(x)\}$ satisfies

$$\mathop{\mathbb{E}}_{x \sim \mu} [\phi(x) - x] = \frac{1}{L} \sum_{\ell=0}^{L-1} [\phi^{\ell+1}(x) - \phi^\ell(x)] = \frac{1}{L} [\phi^L(x) - x] \to 0$$

- Regret minimizer $\mathcal{R}_\Phi$ on $\Phi$

    When $\Phi = \{\text{degree-}k \text{ polynomials}\}$ and the game tree is balanced, regret minimizers with regret $\exp(\text{poly}(k, \log n)) \sqrt{T}$ exist

**Theorem:** There exist efficient regret minimizers with regret $\exp(\text{poly}(k, \log n)) \sqrt{T}$ against the set $\Phi_k$ of degree-$k$ polynomials.

# Swap Regret in Extensive-Form Games

**Q:** Can **swap regret** be efficiently minimized in *extensive-form* games?

### Theorem
[Corollary of Blum-Mansour]

There exists a swap regret minimizer for tree-form strategy sets whose swap regret is $\epsilon T$ after $\mathcal{O}\left(n \cdot 2^n/\epsilon^2\right)$ **iterations.**

Bad per-iteration complexity and convergence rate

### Theorem
[*Special case of* Peng & Rubinstein *STOC'*24; Dagan, Daskalakis, Fishelson, Golowich *STOC'*24]

There exists a swap regret minimizer for tree-form strategy sets* whose swap regret is $\epsilon T$ after $n^{\widetilde{\mathcal{O}}(1/\epsilon)}$ **iterations.**

*or, indeed, any set $X \subset \mathbb{R}^n$ for which *external* regret is minimizable

$\Rightarrow$ For **constant** $\epsilon$, an $\epsilon$-CE can be computed in **polynomial time!**

### Theorem
[Daskalakis, Farina, Golowich, Sandholm, Zhang *arXiv'*24]

There is a constant $c > 0$ such that achieving swap regret $\epsilon T$ in tree-form strategy sets requires $\exp\left(\Omega\left(\min\{n, 1/\epsilon\}^c\right)\right)$ **iterations.**

**Open question:** Can $\epsilon$-CE be computed in time $\text{poly}(n, 1/\epsilon)$ or even $\text{poly}(n, \log(1/\epsilon))$?

(using something other than adversarial no-swap-regret learning)

# TreeSwap

Peng & Rubinstein (*STOC*'24); Dagan, Daskalakis, Fishelson, Golowich (*STOC*'24)

**Given:** External regret minimizer $R_X$ on $X \subset [0,1]^n$ achieving $\epsilon K$ regret after $K$ steps (e.g., for extensive-form games, CFR gives $K = n^2/\epsilon^2$)

**Goal:** Build a **swap regret minimizer** on $X$

**Idea:**



depth = $D = 1/\epsilon$

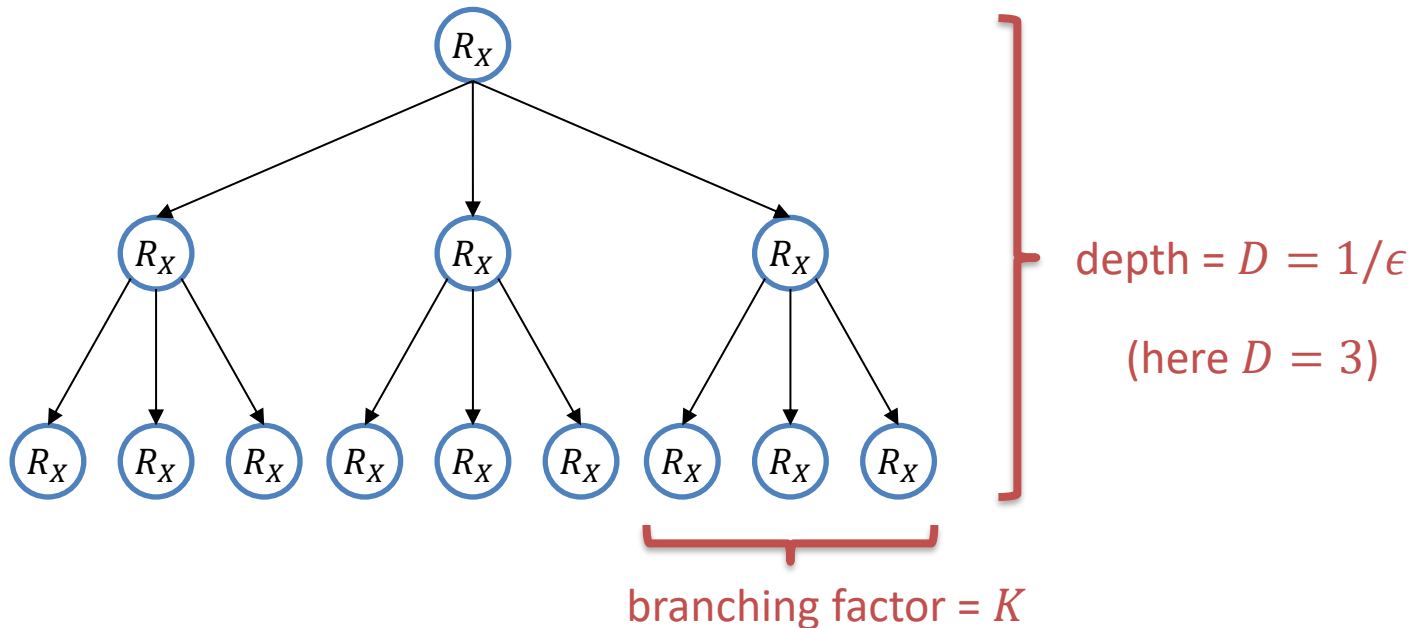(here $D = 3$)

branching factor = $K$

# TreeSwap

Peng & Rubinstein (*STOC*'24); Dagan, Daskalakis, Fishelson, Golowich (*STOC*'24)

**Given:** External regret minimizer $R_X$ on $X \subset [0,1]^n$ achieving $\epsilon K$ regret after $K$ steps (e.g., for extensive-form games, CFR gives $K = n^2/\epsilon^2$)

**Goal:** Build a **swap regret minimizer** on $X$

**Idea:**



Time: $t = 1$

Play $\mu^1 := \text{Unif}\{x_1^1, \ldots, x_{D-1}^1, x_D^1\}$
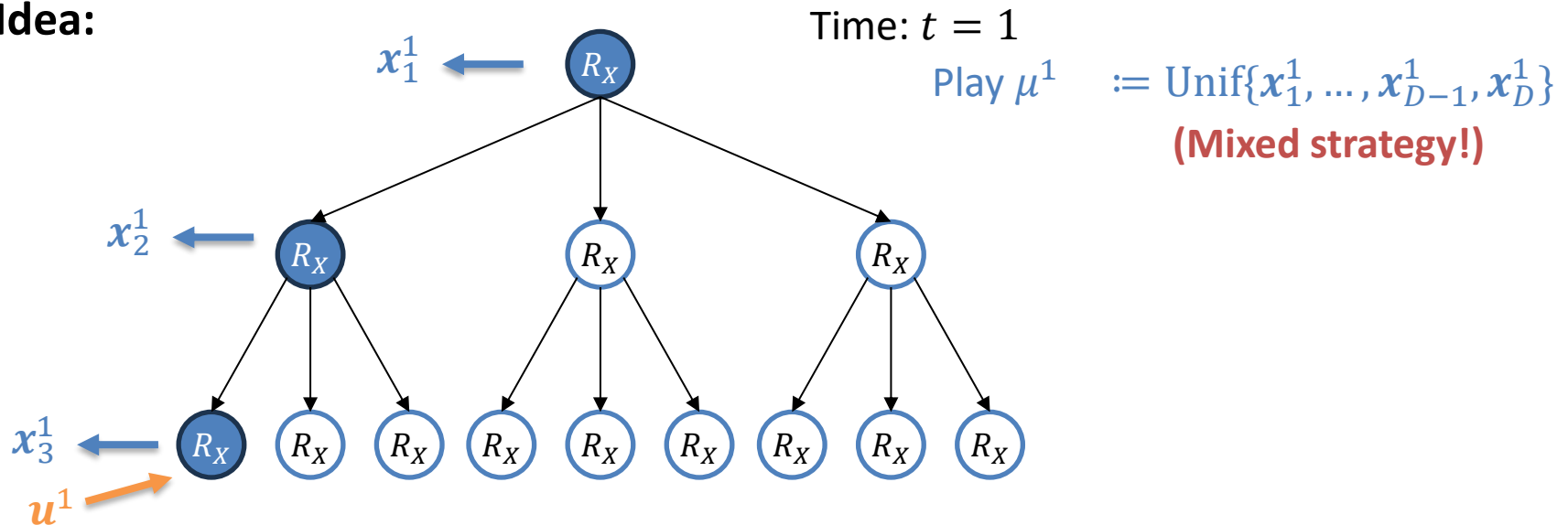
**(Mixed strategy!)**

# TreeSwap

Peng & Rubinstein (*STOC*'24); Dagan, Daskalakis, Fishelson, Golowich (*STOC*'24)

**Given:** External regret minimizer $R_X$ on $X \subset [0,1]^n$ achieving $\epsilon K$ regret after $K$ steps (e.g., for extensive-form games, CFR gives $K = n^2/\epsilon^2$)

**Goal:** Build a **swap regret minimizer** on $X$

**Idea:**



Time: $t = K$

Play $\mu^K := \text{Unif}\{x_1^1, \dots, x_{D-1}^1, x_D^K\}$

**(Mixed strategy!)**

$x_1^1$

$x_2^1$

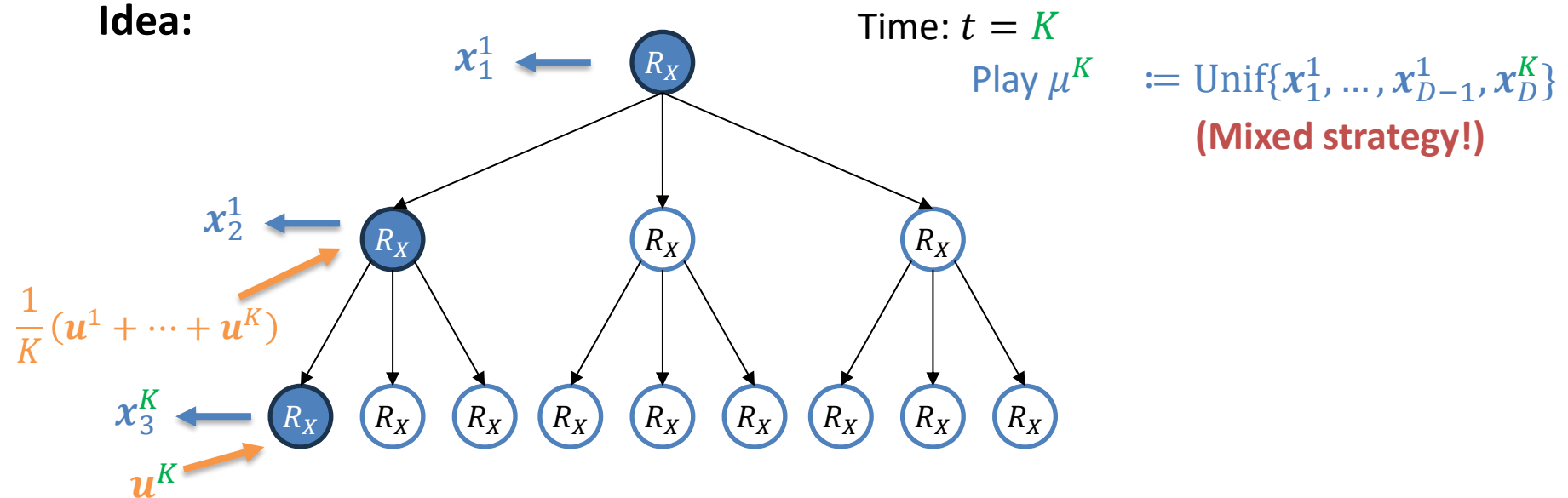$\frac{1}{K}(u^1 + \cdots + u^K)$

$x_3^K$

$u^K$

# TreeSwap

Peng & Rubinstein (*STOC*'24); Dagan, Daskalakis, Fishelson, Golowich (*STOC*'24)

**Given:** External regret minimizer $R_X$ on $X \subset [0,1]^n$ achieving $\epsilon K$ regret after $K$ steps (e.g., for extensive-form games, CFR gives $K = n^2/\epsilon^2$)

**Goal:** Build a **swap regret minimizer** on $X$

**Idea:**



Time: $t = K + 1$

Play $\mu^{K+1} := \mathrm{Unif}\{x_1^1, \ldots, x_{D-1}^2, x_D^{K+1}\}$
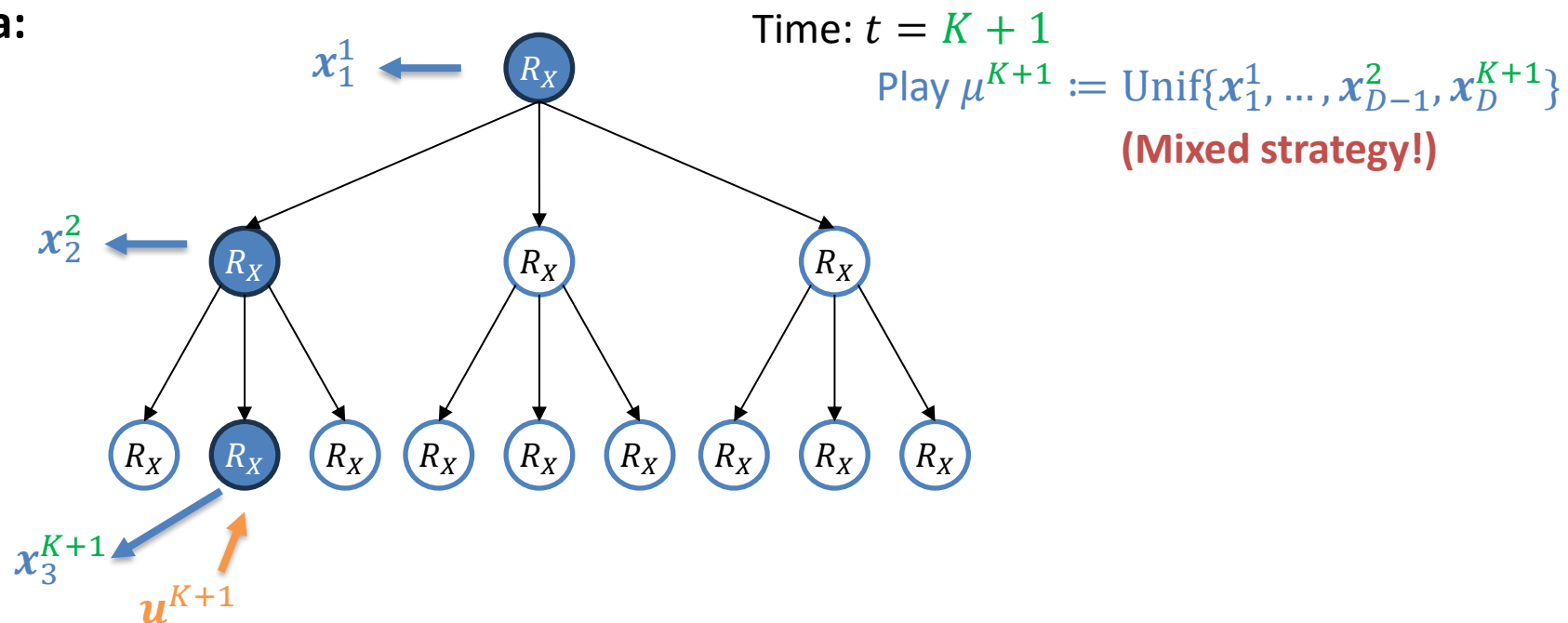
**(Mixed strategy!)**

# TreeSwap

Peng & Rubinstein (*STOC*'24); Dagan, Daskalakis, Fishelson, Golowich (*STOC*'24)

**Given:** External regret minimizer $R_X$ on $X \subset [0,1]^n$ achieving $\epsilon K$ regret after $K$ steps (e.g., for extensive-form games, CFR gives $K = n^2/\epsilon^2$)

**Goal:** Build a **swap regret minimizer** on $X$

**Idea:**



Time: $t = T = K^d$

Play $\mu^{K^d} := \text{Unif}\{\boldsymbol{x}_1^K, \dots, \boldsymbol{x}_{D-1}^{K^{d-1}}, \boldsymbol{x}_D^{K^d}\}$
**(Mixed strategy!)**

**Intuition:** In the GGM framework, if $\mu^t = \text{Unif}\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_D\}$ let $\phi^t$ be the "map" that takes $\boldsymbol{x}_1 \mapsto \boldsymbol{x}_2 \mapsto \cdots \mapsto \boldsymbol{x}_D$

- $\boldsymbol{\mu}^t$ is an expected fixed point of $\phi^t$
- each value of $\phi^t$ is being picked by regret minimizer $\Rightarrow \Phi$-regret is small!
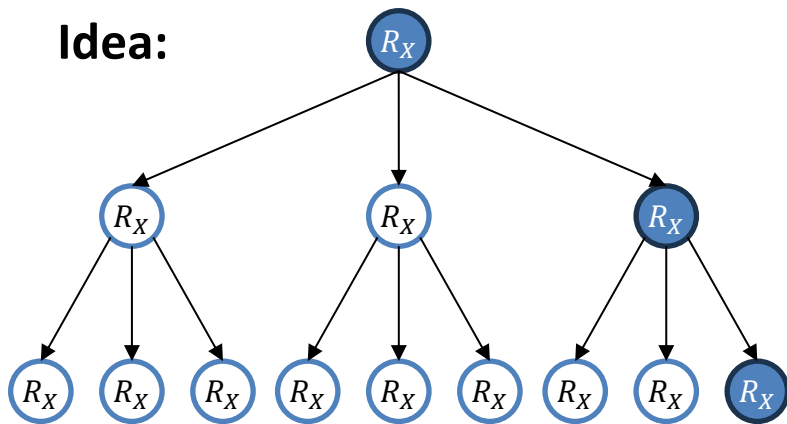
# TreeSwap

Peng & Rubinstein (*STOC*'24); Dagan, Daskalakis, Fishelson, Golowich (*STOC*'24)

**Given:** External regret minimizer $R_X$ on $X \subset [0,1]^n$ achieving $\epsilon K$ regret after $K$ steps (e.g., for extensive-form games, CFR gives $K = n^2/\epsilon^2$)

**Goal:** Build a **swap regret minimizer** on $X$

Time: $t = T = K^d$

Play $\mu^{K^d} := \text{Unif}\{\boldsymbol{x}_1^K, \dots, \boldsymbol{x}_{D-1}^{K^{d-1}}, \boldsymbol{x}_D^{K^d}\}$

**(Mixed strategy!)**

**Theorem:**

$$R_X^{\text{Swap}}(T) \leq T\left(\epsilon + \frac{1}{D}\right) \leq 2\epsilon T$$

from regret of each $R_X$

from expected fixed point error

**Intuition:** In the GGM framework, if $\mu^t = \text{Unif}\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_D\}$ let $\phi^t$ be the "map" that takes $\boldsymbol{x}_1 \mapsto \boldsymbol{x}_2 \mapsto \cdots \mapsto \boldsymbol{x}_D$

- $\mu^t$ is an expected fixed point of $\phi^t$
- each value of $\phi^t$ is being picked by regret minimizer $\Rightarrow \Phi$-regret is small!

# Summary + some further references

*What equilibrium concepts can be reached by **efficient learning algorithms**?*

<span style="color:red">Previously believed to be
the limit of GGM</span>

| Correlated equilibrium concept | | Normal-form coarse-correlated | Extensive-form correlated | Linear-swap correlated | Low-degree swap correlated | Normal-form correlated |
|---|---|---|---|---|---|---|
| Set of deviations $\Phi$ | | Constant functions | "Trigger" functions | Linear functions | Degree-$k$ polynomials | All functions |
| **Best-known algorithm** | # iterations for $\epsilon T$ regret | $n/\epsilon^2$ | $nbd/\epsilon^2$ | $n^4/\epsilon^2$ | $n^{\mathcal{O}(kd\log b)^3}/\epsilon^2$ | $n^{\tilde{\mathcal{O}}(1/\epsilon)}$ |
| | Per-iteration complexity | $n$ | $\mathrm{FP}(n)$ | $\mathrm{FP}(n)$ | $n^{\mathcal{O}(kd\log b)^3}/\epsilon$ | $n/\epsilon$ |
| | Citation | Farina, Lee, Luo, Kroer *ICML*'22 | Farina, Celli, Marchesi, Gatti *JACM*'22 | Zhang, Farina, Sandholm *ICLR*'24 | Zhang, Anagnostides Farina, Sandholm *arXiv*'24 | Peng & Rubinstein *STOC*'24; Dagan, Daskalakis, Fishelson, Golowich *STOC*'24 |

Notation:
$b$ = branching factor of game
$d$ = depth of game
$\mathrm{FP}(n)$ = time complexity of computing a fixed point of an $n \times n$ matrix
$\mathrm{QP}(n)$ = time complexity of solving an $n$-variable convex quadratic program

<span style="color:#4a86c8">Tighter equilibrium concepts
Larger sets $\Phi$
Harder to learn</span>

# Summary + some further references

*What equilibrium concepts can be reached by **efficient learning algorithms**?*

**Previously believed to be the limit of GGM**

| Correlated equilibrium concept | Normal-form coarse-correlated | Extensive-form correlated | Linear-swap correlated | Low-degree swap correlated | Normal-form correlated |
|---|---|---|---|---|---|
| **Set of deviations $\Phi$** | Constant functions | "Trigger" functions | Linear functions | Degree-$k$ polynomials | All functions |
| **# iterations for $\epsilon T$ regret** | $n/\epsilon^2$ | | $n^4/\epsilon^2$ | $n^{\mathcal{O}(kd\log b)^3}/\epsilon^2$ | $n^{\tilde{\mathcal{O}}(1/\epsilon)}$ |
| **Per-iteration complexity** | | | $\mathrm{FP}(n)$ | $n^{\mathcal{O}(kd\log b)^3}/\epsilon$ | $n/\epsilon$ |
| **Citation** | Farina, Lee, Luo, Kroer *ICML*'22 | Farina, Celli, | Zhang, | Zhang, | Peng & Rubinstein *STOC*'24; Daskalakis, |

One UTC mediator

$O(kd\log b)^3$ UTC mediators + "expected fixed points" to circumvent PPAD-hard computation

...concepts
Larger sets $\Phi$
Harder to learn

Notation:
$b$ = branching factor of game
$d$ = depth of game
$\mathrm{FP}(n)$ = time complexity of computing a fixed point of an $n \times n$ matrix
$\mathrm{QP}(n)$ = time complexity of solving an $n$-variable convex quadratic program

# References

- A Blum, Y Mansour (*JMLR* 2007), "From external to internal regret"

- GJ Gordon, A Greenwald, C Marks (*ICML* 2008), "No-regret learning in convex games"

- BH Zhang, G Farina, T Sandholm (*ICML* 2023), "Team belief DAG: generalizing the sequence form to team games for fast computation of correlated team max-min equilibria via regret minimization"

- BH Zhang, G Farina, T Sandholm (*ICLR* 2024), "Mediator Interpretation and Faster Learning Algorithms for Linear Correlated Equilibria in General Extensive-Form Games"

- BH Zhang, I Anagnostides, G Farina, T Sandholm (*arXiv* 2024), "Efficient Φ-Regret Minimization with Low-Degree Swap Deviations in Extensive-Form Games"

- C Daskalakis, G Farina, N Golowich, T Sandholm, BH Zhang (*arXiv* 2024), "A Lower Bound on Swap Regret in Extensive-Form Games"