

Finding and Certifying (Near-)Optimal Strategies in Black-Box Extensive-Form Imperfect-Information Games

Tuomas Sandholm

Carnegie Mellon University

Strategic Machine, Inc.

Strategy Robot, Inc.

Optimized Markets, Inc.

Joint work with my PhD student Brian Hu Zhang
[NeurIPS-20, AAI-21 & draft]



**STRATEGY
ROBOT, INC.**

**STRATEGIC
MACHINE, INC.**



There has been amazing progress in game solving over the last 17 years.

Modules for solving games

- Automated abstraction
 - State abstraction [Gilpin & S., AAAI-06, AAMAS-07, JACM-07, AAAI-08; Gilpin, S. & Sørensen, AAAI-07; S. & Singh, EC-12; Kroer & S., EC-14, AAMAS-15, EC-16, NeurIPS-18; Ganzfried & S., AAAI-14; Brown & S., IJCAI-15; Brown, Ganzfried & S., AAMAS-15]
 - Action abstraction [Ganzfried & S., IJCAI-13; Brown & S., AAAI-14; Kroer & S., AAMAS-15]
- Real-time subgame solving [Ganzfried & S., AAMAS-15; Brown & S., NeurIPS-17, Science-18]
 - Depth-limited search [Brown, S. & Amos, NeurIPS-18; Brown & S., Science-19]
- Equilibrium-finding algorithms
 - Leading regret-minimization algorithms [Farina, Kroer & S., AAAI-19, ICML-19a,b, NeurIPS-19, ICML-20, AAAI-21; Brown & S., AAAI-19; Farina, Kroer, Brown & S., ICML-19; Farina & S., AAAI-21]
 - Incorporating deep learning [Brown, Lerer, Gross & S., ICML-19]
 - Leading first-order optimization methods [Hoda, Gilpin, Peña & S. *Mathematics of Operations Research-10*; Gilpin & S., AAMAS-10; Gilpin, Peña & S., *Mathematical Programming-12*; Kroer, Farina & S., NeurIPS-18; Kroer, Waugh, Kilinc-Karzan & S., EC-16, EC-17, *Mathematical Programming-20*]
 - Pruning [Brown & S., NeurIPS-15, ICML-17, Science-18, Science-19; Brown, Kroer & S., AAAI-17]
 - Sound warm starting [Brown & S., AAAI-14, AAAI-16]
 - Automatically sparsified LP for equilibrium finding [Zhang & S., ICML-20]
 - Computing equilibria by incorporating qualitative models [Ganzfried & S., AAMAS-10]
- Algorithms for equilibrium refinements [Kroer, Farina & S., IJCAI-17, AAAI-18; Farina, Kroer & S., ICML-17; Farina, Gatti & S., NeurIPS-18; Farina, Marchesi, Kroer, Gatti & S., IJCAI-18; Marchesi, Farina, Kroer, Gatti & S., AAAI-19]
- Self-improvement techniques [Brown & S., Science-18]
- Finding correlated and coarse correlation equilibria [Farina, Ling, Fang & S., NeurIPS-19a,b; Farina & S., NeurIPS-20]
- Algorithms for multi-player games [Berg & S., AAAI-17; Brown & S., Science-19]
- Solving team games with pre-game correlation in the team [Farina, Celli, Gatti & S., NeurIPS-18, draft-21]
- Opponent exploitation techniques [Ganzfried & S., AAMAS-11, TEAC-15; S. AAAI-15; Kroer & S., IJCAI-16, AIJ-20; Kroer, Farina & S., AAAI-18]

STRATEGY
ROBOT, INC.

STRATEGIC
MACHINE, INC.

What if the game model is inaccurate or unknown?

1. Sensitivity analysis
2. Lossy game abstraction techniques with ϵ -exploitability guarantees
[S. & Singh, EC-12; Kroer & S., EC-14, AAMAS-15, EC-16, NeurIPS-18]
apply to modeling also
3. **THIS TALK:** First techniques for computing provably (near-)equilibrium strategies while searching only a tiny fraction of the game tree
[Zhang & S., NeurIPS-20, AAI-21]
 - => algorithm with optimal $\tilde{O}(\#\text{nodes}/\sqrt{T})$ convergence in this setting
 - Prior methods (such as MCCFR) can be exponential in tree size

Black-Box Games

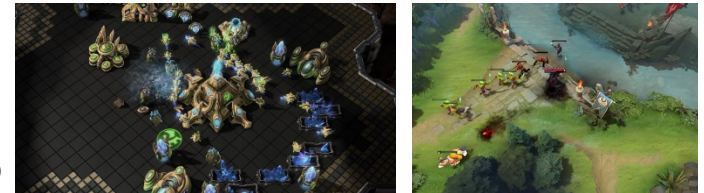
- Game is not explicitly given in the form of rules, but rather via access to playing it
 - We can control all players during the practice phase
- E.g., war games, strategy video games, and financial simulations



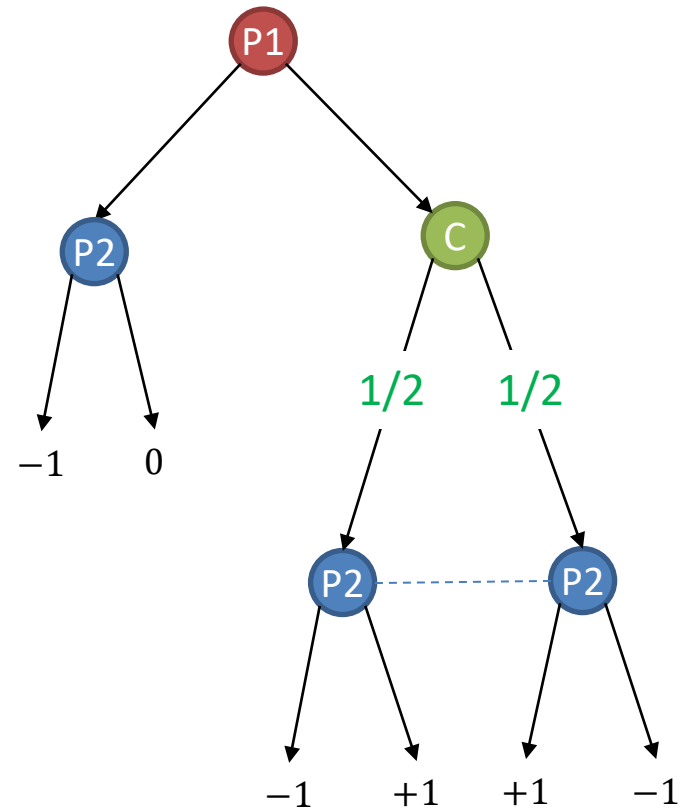
**STRATEGY
ROBOT, INC.**

Learning to Play Black-Box Games

- **Deep Reinforcement Learning** (e.g., *AlphaStar* [Vinyals et al., 2019], *OpenAI Five* [Berner et al., 2019])
 - Great practical performance for a while
 - **Issue:** No exploitability bounds
 - Leads to strategies that can be beaten in practice also
- **Bandit Regret Minimization** [Farina & Sandholm AAAI-21]
 - Converges to ϵ -equilibrium after $\text{poly}(N, 1/\epsilon)$ game samples (N = size of game)
 - **Issues (online MCCFR [Lanctot et al. 2009] has these issues also and other issues):**
 - Worst-case exploitability bounds are trivial until number of iterations is much larger than N
 - Need to expand rest of game tree to compute *ex-post* exploitability guarantee
- **Certificates [This work]**
 - Compute Nash equilibrium by incrementally expanding game tree
 - Exploitability bounds always computable *ex post* without expanding remainder of tree!

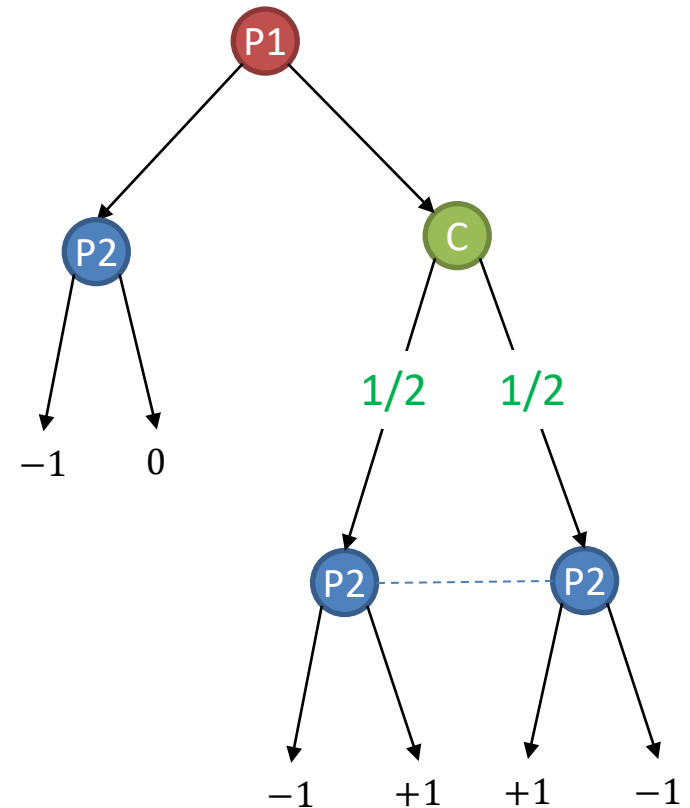


Extensive-Form Games



Pseudogames and Certificates

Pseudogame: Game without known utilities on all terminal nodes



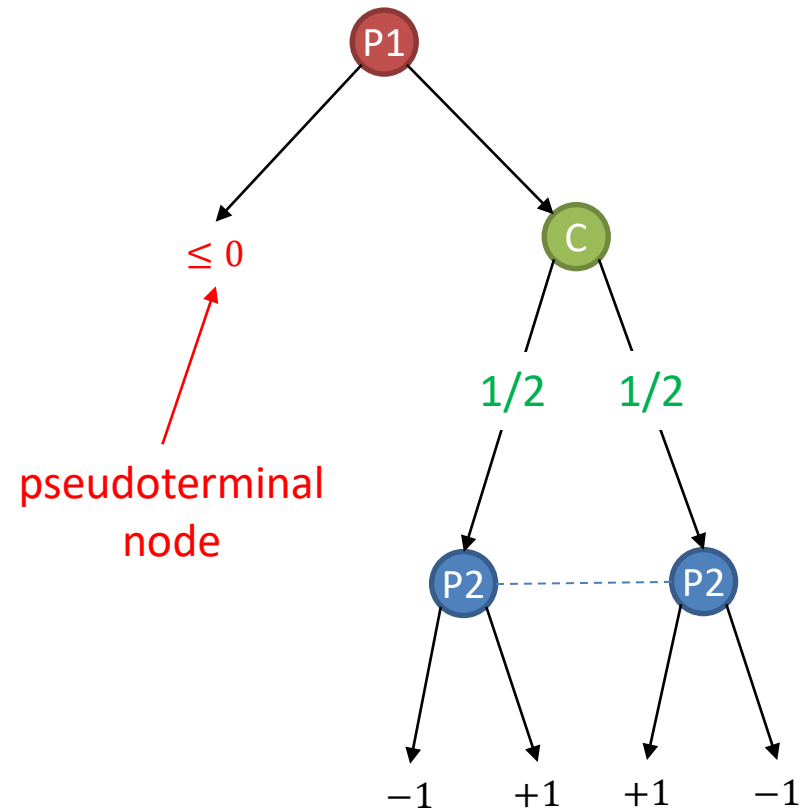
Pseudogames and Certificates

Pseudogame: Game without known utilities on all terminal nodes

Think: partially-expanded game tree, “alpha-beta” style

In zero-sum land, gives rise to **two** games:

- a *lower-bound game* in which rewards are optimistic for P2, and
- an *upper-bound game* in which rewards are optimistic for P1

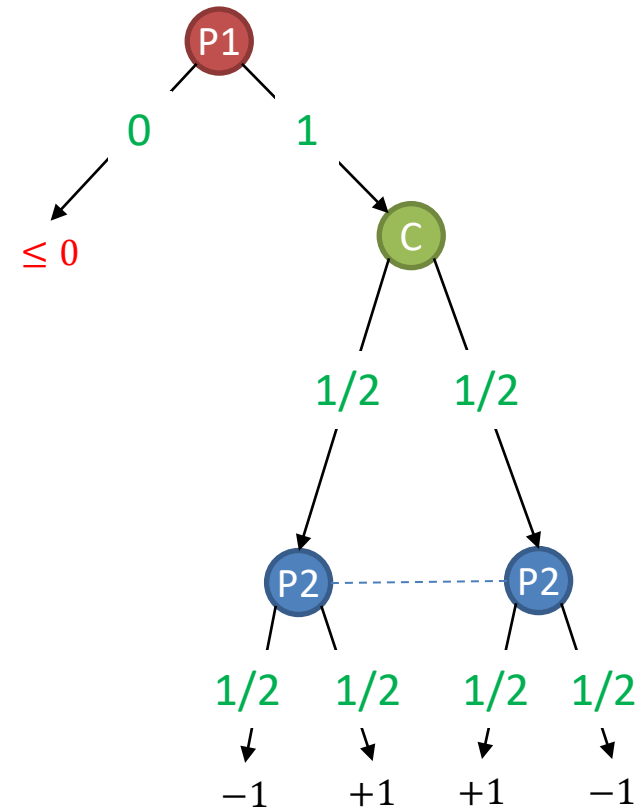


Pseudogames and Certificates

(Approximate) Nash equilibrium in a pseudogame: strategy profile in which every player is *provably* playing an (approximate) best response (irrespective of what happens at pseudoterminal nodes)

Results in Nash equilibrium regardless of what the pseudoterminal node hides!

(Approximate) Certificate:
Pseudogame created from partial expansion of a full game + (approximate) Nash equilibrium of that pseudogame



Small Certificates

Question: When do small ε -certificates exist?
Specifically, size $O(N^c \text{poly}(1/\varepsilon))$ for some $c < 1$

Again, N is the number of nodes.

When do Small Certificates Exist?

- **Answer #1:** They exist in **perfect-information zero-sum games with no nature randomness**,
...under reasonable assumptions about the game tree (e.g., uniform branching factor and depth, alternating moves)
 - **Proof:** The optimal alpha-beta search tree is a certificate of size $\approx \sqrt{N}$.

Small Certificates

Answer #2: They exist in (squarish) **normal-form games**.

Proof:

Lipton et al., 2003:

ε -Nash equilibrium exists where each player mixes between $\log(m) / \varepsilon^2$ pure strategies.

	A	B	C	D	E	F	G
T							
U							
V							
W							
X							
Y							
Z							

Small Certificates

Answer #2: They exist in (squarish) **normal-form games**.

Proof:

Lipton et al., 2003:

ε -Nash equilibrium exists where each player mixes between $\log(m) / \varepsilon^2$ pure strategies.

	A	B	C	D	E	F	G
T							
U							
V							
W							
X							
Y							
Z							

Small Certificates

Answer #2: They exist in (squarish) **normal-form games**.

Proof:

We only need those rows and columns!

$\Rightarrow O(m \log(m) / \varepsilon^2)$ -size certificate

	A	B	C	D	E	F	G
T							
U							
V							
W							
X							
Y							
Z							

Small Certificates

So... small certificates exist in games where the players have **perfect information** or **no information**.

What about in between?

A: Unfortunately, not in all games.

Small Certificates

Counterexample: Consider this game:

- Matching pennies
- repeated T times, each round worth $1/T$ points.
- After each round, P2 learns what P1 played, but P1 doesn't learn what P2 played.

Game tree size: 4^T

Theorem: Any ε -certificate of this game must have size

$$\Omega\left(4^{T(1-O(\varepsilon))}\right).$$

Proof Sketch: P1's strategy **must have high entropy**, but this is not possible unless lots of nodes get expanded

Bad News

Theorem: It is NP-hard to approximate the smallest certificate of an extensive-form zero-sum game, to better than an $O(\log N)$ multiplicative factor.

Proof Idea: Reduction from set cover.

Simulators

Assume access to a **simulator**:

- Allows us to play through the game **from the perspective of all players at once**
- Gives **bounds** (not necessarily tight) on future utilities
- Allows us to control nature actions (I'll relax this later..)

Goal:

- *Compute and verify “ex-post” approximate equilibria with only black-box access*
- Output both an equilibrium strategy **and** a bound ε on exploitability

More Bad News

Theorem: With only a black-box simulator of an extensive-form zero-sum game, there is no equilibrium-finding algorithm that runs in time polynomial in the size of the smallest certificate.

Proof: One-player “SAT” games: certificate of size $O(\log N)$ exists, but clearly no sublinear-time algorithm.

Let's Do It Anyway

Repeat until satisfied:

- **Solve** both the upper- and lower-bound pseudogames *exactly* (e.g., using an LP solver)
- **Create** the next pseudogame by expanding all pseudoterminal nodes in the support of the **optimistic profile** (in which the max-player plays her equilibrium strategy in the upper-bound game, and the min-player plays her strategy in the lower-bound game)

Output: Pessimistic profile, and ε = difference in values between upper- and lower-bound pseudogames

Intuition: In the perfect-information setting with no nature randomness, it's just **alpha-beta search**

Let's Do It Anyway

Repeat until satisfied:

- **Solve** both the upper- and lower-bound pseudogames *exactly* (e.g., using an LP solver)
- **Create** the next pseudogame by expanding all pseudoterminal nodes in the support of the **optimistic profile** (in which the max-player plays her equilibrium strategy in the upper-bound game, and the min-player plays her strategy in the lower-bound game)

Output: Pessimistic profile, and $\varepsilon =$ difference in values between upper- and lower-bound pseudogames

Theorem: The expansion in the second step expands a node if and only if the game is not already solved.

Let's Do It Anyway

Repeat until satisfied:

- **Solve** both the upper- and lower-bound pseudogames *exactly* (e.g., using an LP solver)
- **Create** the next pseudogame by expanding all pseudoterminal nodes in the support of the **optimistic profile** (in which the max-player plays her equilibrium strategy in the upper-bound game, and the min-player plays her strategy in the lower-bound game)

Output: Pessimistic profile, and ε = difference in values between upper- and lower-bound pseudogames

Works even on games that have unbounded rewards!

Experiments

game	size of game		size of certificate			
	nodes	infosets	nodes		infosets	
search game	234,705	11,890	13,682	5.8%	532	4.5%
4-rank PI Goofspiel	2,229	1,653	275	12.3%	110	6.7%
5-rank PI Goofspiel	55,731	41,331	2,593	4.7%	957	2.3%
6-rank PI Goofspiel	2,006,323	1,487,923	21,948	1.1%	7,584	0.5%
4-rank Goofspiel	2,229	738	614	27.5%	117	15.9%
5-rank Goofspiel	55,731	9,948	11,415	20.5%	2,160	21.7%
6-rank Goofspiel	2,006,323	166,002	266,756	13.3%	15,776	9.5%
3-rank random Goofspiel	1,066	426	309	29.0%	92	21.6%
4-rank random Goofspiel	68,245	17,432	16,416	24.1%	3,270	18.8%
5-rank random Goofspiel	8,530,656	1,175,330	1,854,858	21.7%	241,985	20.6%
5-rank Leduc	∞	∞	26,306	—	2,406	—
9-rank Leduc	∞	∞	137,662	—	6,811	—
13-rank Leduc	∞	∞	337,312	—	12,171	—

Simulators

Assume access to a **simulator**:

- Allows us to play through the game **from the perspective of all players at once**
- Gives **bounds** (not necessarily tight) on future utilities
- Allows us to control nature actions

Goal:

- *Compute and verify “ex-post” approximate equilibria with only black-box access*
- Output both an equilibrium strategy **and** a bound ε on exploitability

Simulators

Assume access to a **simulator**:

- Allows us to play through the game **from the perspective of all players at once**
- Gives **bounds** (not necessarily tight) on future utilities
- ~~Allows us to control nature actions~~

Goal:

- *Compute and verify “ex-post” approximate equilibria with only black-box access*
- Output both an equilibrium strategy **and** a bound ε on exploitability

Simulators

Assume access to a **simulator**:


- Allows us to play through the game **from the perspective of all players at once**
- Gives **bounds** (not necessarily tight) on future utilities
- ~~Allows us to control nature actions~~

Goal:

- *Compute and verify “ex-post” approximate equilibria with only black-box access*
- Output both an equilibrium strategy **and** a bound ε on exploitability
- **Want:** correctness with high probability, say, $1 - T^{-\gamma}$ for some $\gamma > 0$ after T iterations.

Roadmap for the Rest of the Talk: Certificate-Finding in Zero-Sum Games

	Sampling-limited	Compute-limited
Visiting a chance node gives the full distribution at that node	Our NeurIPS-20 paper	
Visiting a chance node gives an action sample at that node		



Usable even in general-sum games
(computes coarse-correlated equilibrium)

Lower Bound

Theorem: Consider any algorithm with the following guarantee.

For some constant $\gamma > 0$,

given a zero-sum game in our black-box setting,

with T game samples,

the algorithm outputs a pair of strategies (x, y) **and a bound** ε_T such that, with probability $1 - O(T^{-\gamma})$,

(x, y) is an ε_T -Nash equilibrium.

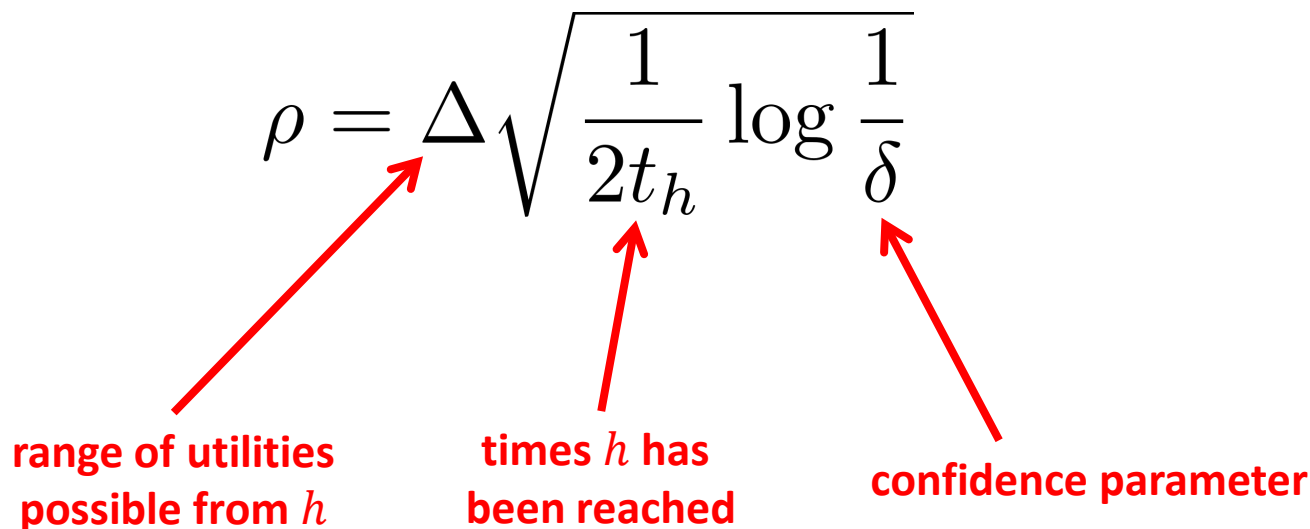
Then

$$\varepsilon_T = \Omega \left(\sqrt{\frac{\log T}{T}} \right)$$

Our goal: Match this bound.

Main Tool: Pseudogames as Confidence Bounds

- At **nodes that have not yet been expanded**, use bounds given by simulator
- At **nature nodes h** , give each player reward bounded by $[-\rho, \rho]$, where

$$\rho = \Delta \sqrt{\frac{1}{2t_h} \log \frac{1}{\delta}}$$


range of utilities possible from h

times h has been reached

confidence parameter

Main Tool: Pseudogames as Confidence Bounds

- At **nodes that have not yet been expanded**, use bounds given by simulator.
- At **nature nodes h** , give each player reward bounded by $[-\rho, \rho]$, where

$$\rho = \Delta \sqrt{\frac{1}{2t_h} \log \frac{1}{\delta}}$$

Intuition: ρ represents the **uncertainty** in the nature distribution at h

Main Tool: Pseudogames as Confidence Bounds

- At **nodes that have not yet been expanded**, use bounds given by simulator.
- At **nature nodes** h , give each player reward bounded by $[-\rho, \rho]$, where

$$\rho = \Delta \sqrt{\frac{1}{2t_h} \log \frac{1}{\delta}}$$

Intuition: It looks like UCB. That is not a coincidence, as I'll discuss.

Choice of Confidence Bound

During equilibrium computation, values of children are changing, so we need to use a Hoeffding bound to be robust:

$$\rho = \Delta \sqrt{\frac{1}{2t_h} \log \frac{1}{\delta}}$$

where Δ is the range of possible utilities from h

NEW IDEA SINCE OUR AAAI-21 PAPER:

During best response computation, strategy profiles after h are fixed by induction, so we can use a tighter empirical Bernstein bound [Maurer & Pontil '09]:

$$\rho = S \sqrt{\frac{2}{t_h} \log \frac{2}{\delta}} + \frac{7\Delta'}{3(t_h - 1)} \log \frac{2}{\delta}$$

where S is the unbiased sample standard deviation, and Δ' is the range of possible utilities from h **under the fixed strategy profile**, which may be much smaller than Δ

Main Tool: Pseudogames as Confidence Bounds

Theorem: For appropriate choice of δ , with high probability, at every time, for every strategy profile, for every player, the true reward of the player is bounded by the pessimistic and optimistic rewards achieved in the confidence bound pseudogame.

(i.e., “confidence bounds are actually bounds”)

Zero-Sum LP-Based Algorithm

Repeat T times:

- **Solve** both the upper- and lower-bound pseudogames *exactly* (e.g., using an LP solver)
- **Sample** one play-through from the optimistic profile (in which the max-player plays her equilibrium strategy in the upper-bound game, and the min-player plays her strategy in the lower-bound game)
- **Create** the next pseudogame:
 - **Expand** all encountered nodes
 - **Update** empirical nature distributions of nature nodes sampled during play

Output: Pessimistic profile, and ε_T = difference in values between upper- and lower-bound pseudogames

Intuition: In the perfect-information setting with no nature randomness, it's just **alpha-beta search**

Zero-Sum LP-Based Algorithm

Repeat T times:

- **Solve** both the upper- and lower-bound pseudogames *exactly* (e.g., using an LP solver)
- **Sample** one play-through from the optimistic profile (in which the max-player plays her equilibrium strategy in the upper-bound game, and the min-player plays her strategy in the lower-bound game)
- **Create** the next pseudogame:
 - **Expand** all encountered nodes
 - **Update** empirical nature distributions of nature nodes sampled during play

Output: Pessimistic profile, and ε_T = difference in values between upper- and lower-bound pseudogames

Intuition: In the one-player “multi-armed bandit” setting, it’s **UCB** (except algorithm has a different constant in the upper confidence bound term, and so does the regret bound).

Zero-Sum LP-Based Algorithm

Advantage: Sample-efficient

Disadvantage: Expensive iterations (requires game re-solve on each iteration)

- We warm start from the previous LP, whose values typically change very little based on the one new sample

Theorem: The *best iterate* of the algorithm

converges at rate $\mathbb{E} \varepsilon_T \leq O \left(N_T \sqrt{\frac{\log T}{T}} \right)$

number of nodes in current pseudogame
(may be \ll total number of nodes!)

Regret-Based Algorithm

Idea: Just use a regret minimizer, like CFR, for each player

Regret-Based Algorithm

Repeat T times:

- **Query** the regret minimizers for all players to obtain a strategy profile
- **Sample** one play-through from that strategy profile
- **Pass** each player's regret minimizer that player's *optimistic* reward
- **Create** the next pseudogame:
 - **Expand** all encountered nodes
 - **Update** empirical nature distributions of nature nodes sampled during play

Output: Average strategy profile

Several problems!

Regret-Based Algorithm

Repeat T times:

- **Query** the regret minimizers for all players to obtain a strategy profile
- **Sample** one play-through from that strategy profile
- **Pass** each player's regret minimizer that player's *optimistic* reward
- **Create** the next pseudogame:
 - **Expand** all encountered nodes
 - **Update** empirical nature distributions of nature nodes sampled during play

Output: Average strategy profile

Problem 1: The strategy space of each player is changing over time

Solution: CFR “handles it naturally”. *Formalization:* “Extendable” regret minimizers

Regret-Based Algorithm

Repeat T times:

- **Query** the regret minimizers for all players to obtain a strategy profile
- **Sample** one play-through from that strategy profile
- **Pass** each player's regret minimizer that player's *optimistic* reward
- **Create** the next pseudogame:
 - **Expand** all encountered nodes
 - **Update** empirical nature distributions of nature nodes sampled during play

Output: Average strategy profile

Problem 2: We don't want to run a full CFR iterate on every sample; that is expensive

Solution: Use MCCFR + outcome sampling. Nothing breaks

Regret-Based Algorithm

Repeat T times:

- **Query** the regret minimizers for all players to obtain a strategy profile
- **Sample** one play-through from that strategy profile
- **Pass** each player's regret minimizer that player's *optimistic* reward
- **Create** the next pseudogame:
 - **Expand** all encountered nodes
 - **Update** empirical nature distributions of nature nodes sampled during play

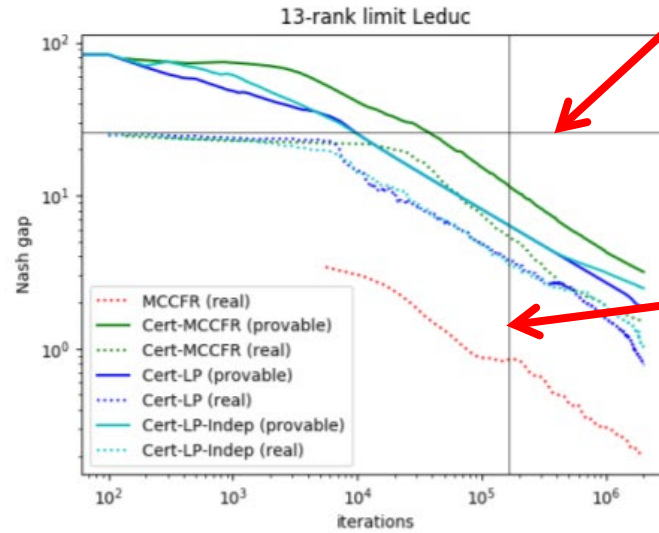
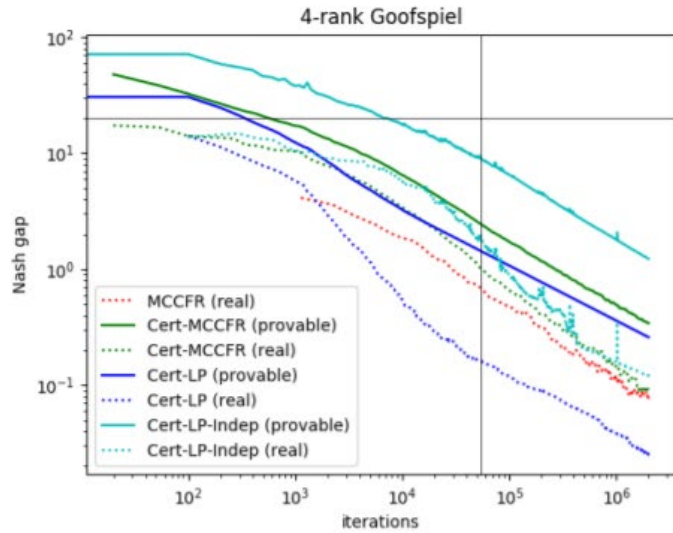
Output: Average strategy profile

Problem 3: What equilibrium gap bound can we compute?

What Equilibrium Gap Bound Can We Compute?

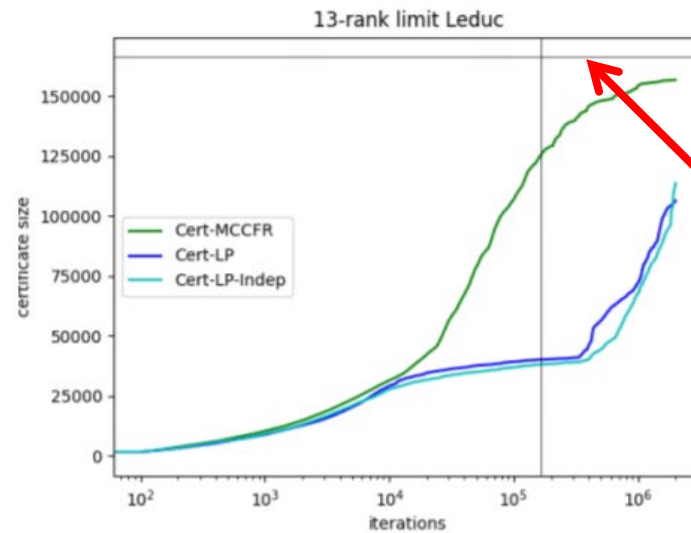
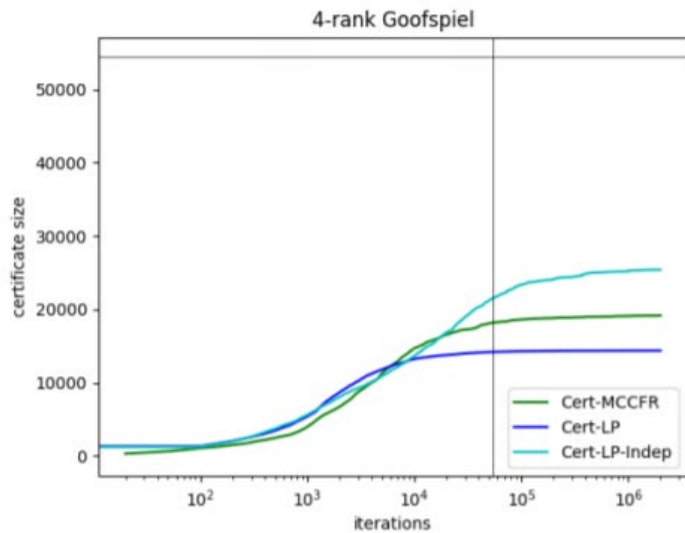
- The natural game-specific equilibrium gap bound —used in our exact LP-based algorithm— (difference in optimistic best response values using the final pseudogame) **doesn't converge as $\tilde{O}(1/\sqrt{T})$** in the worst case
- ...but, we know that the *worst-case-over-games* equilibrium gap bound of the algorithm *does* converge as $\tilde{O}(1/\sqrt{T})$ (for the same reason that MCCFR does)
- **Solution:** In practice, take the former; it's basically always smaller. In theory, take the minimum of the two

Experiments



Horizontal line:
range of a
player's reward
in full game

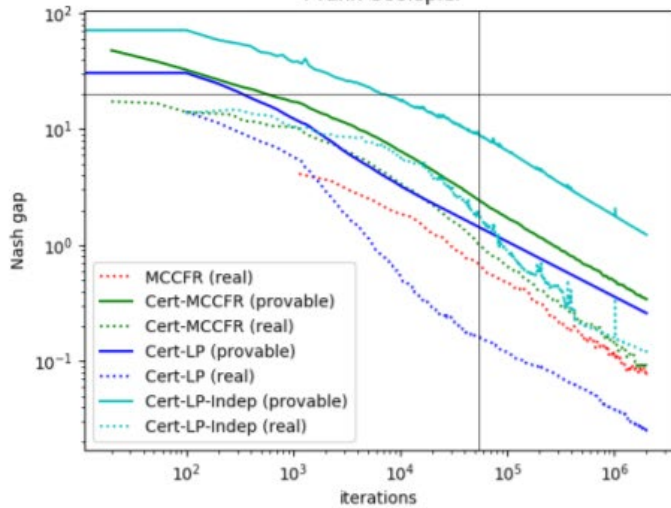
Vertical line:
number of
nodes in full
game



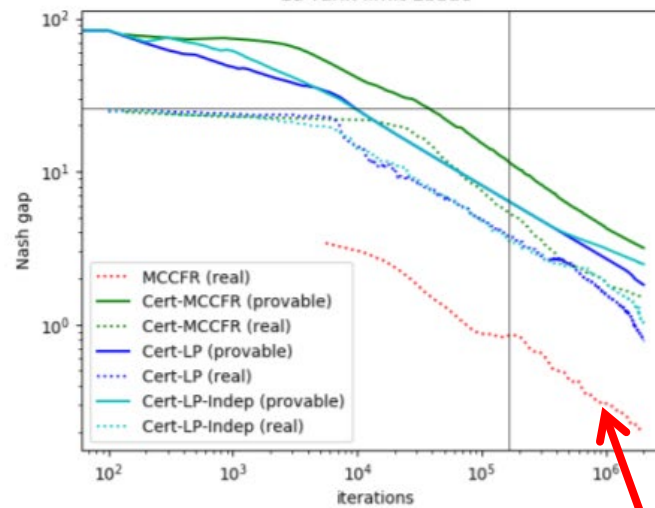
Horizontal line:
number of
nodes in full
game

Experiments

4-rank Goofspiel

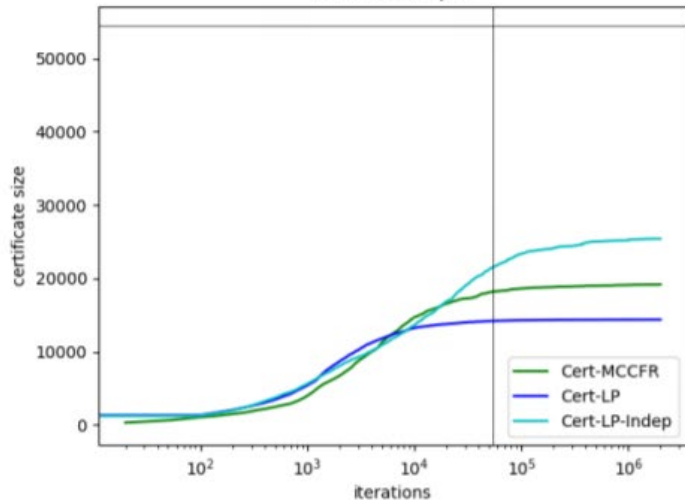


13-rank limit Leduc

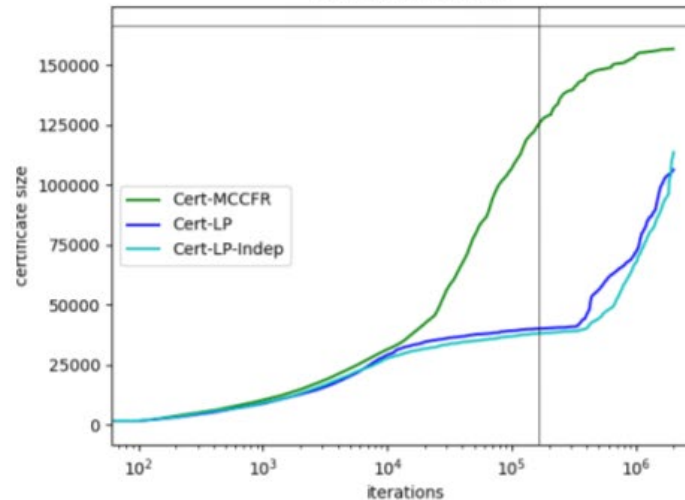


In all games, with all algorithms, nontrivial certificates are found **without expanding the full game tree**, in fact, **with fewer game samples than there are game tree nodes**

4-rank Goofspiel

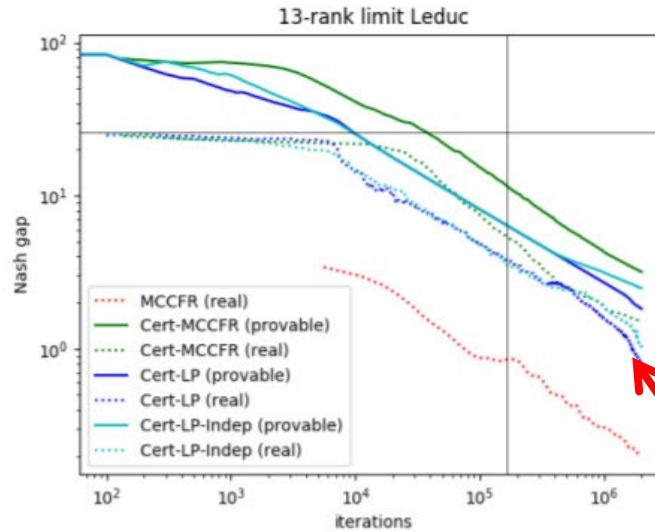
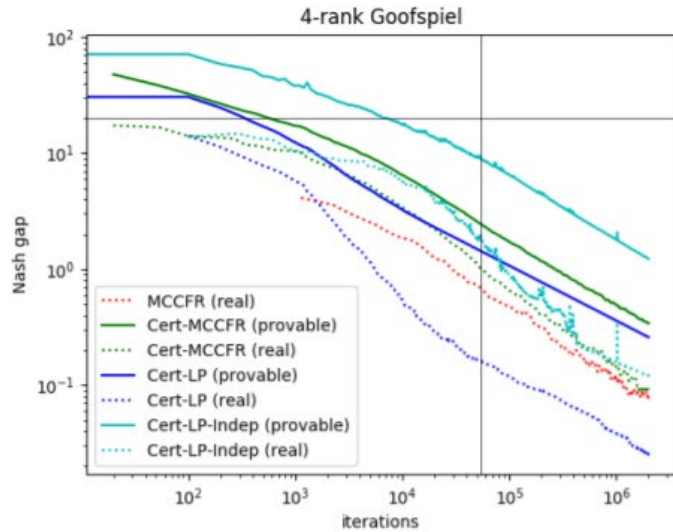


13-rank limit Leduc



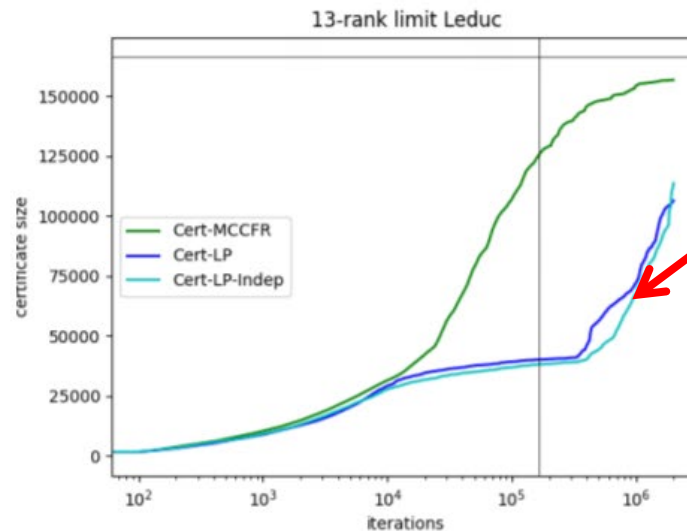
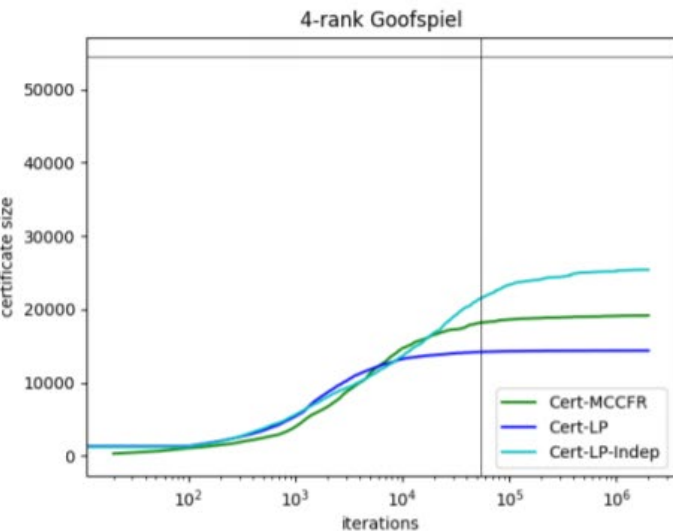
MCCFR converges quickly in reality, but this cannot be verified without expanding the rest of the game tree (or knowing something else that enables best-response computation)

Experiments



In all games, with all algorithms, nontrivial certificates are found without expanding the full game tree, in fact, with fewer game samples than there are game tree nodes

LP-based certificate finding has better sample efficiency and final certificate size than regret-based, but (not shown) runs slower



Conclusion

Black-box imperfect-information games
(of at least moderate size)
can now be **solved**

i.e., we can get the non-exploitability
guarantee of game theory

This talk covered parts of the following papers and a new concentration result

- Small Nash Equilibrium Certificates in Very Large Games, *NeurIPS-20*
<https://arxiv.org/abs/2006.16387>
- Finding and Certifying (Near-)Optimal Strategies in Black-Box Extensive-Form Games, *AAAI-21*
<https://arxiv.org/abs/2009.07384>