

Adversarial Policies Beat Superhuman Go AIs

Tony Wang* Adam Gleave* Tom Tseng Kellin Pelrine
Nora Belrose Joseph Miller Michael Dennis Yawen Duan
Viktor Pogrebniak Sergey Levine Stuart Russell

twang6@mit.edu

goattack.far.ai

Link:

https://docs.google.com/presentation/d/1pvDdWQ3lbHZIz4dltoFdBcpaog_zqb_jC1obRZecl0Y/edit#slide=id.p

Hi, my name is Tony and I want to teach you about how adversarial policies beat superhuman Go AIs.

[click]



Go is an ancient Chinese board game invented over 2500 years ago. Two players take turns placing black and white stones on a square board, trying to surround territory, and kill their opponent's stones. At the end of the game, whoever controls more of the board wins.

[click]



AlphaGo def. Lee Sedol (4-1)

The game you just saw was part of a match between Lee Sedol and the AI AlphaGo. Lee (on the right here) was one of the strongest humans to ever play the game. However, AlphaGo pulled off an upset, and beat Lee 4 to 1. AlphaGo's victory was a great demonstration of the power of deep learning.

[click]

AlphaGo
(2016)



AlphaZero
(2017)



Katago
(2023)

~99.97%
winrate

~98%
winrate

Progress did not stop with AlphaGo. [click] A year afterwards, Deepmind published AlphaZero, a system more general than AlphaGo and much much stronger. The current state of the art is even further along. [click] We estimate that KataGo, currently the strongest open-source Go AI, beats AlphaZero 98% of the time. However, it turns out that almost all of these AIs have a hidden weakness.

[click, pause]

Sources:

1. AlphaGo vs. AlphaZero (<https://www.nature.com/articles/nature24270>), 5185 - 3739 = 1446 => 99.97% winrate
2. From our paper, AlphaZero_s800 has 3813 elo of goratings.org. cp505_s1 has 2738 elo on goratings.org. cp505_s800 has 4500 elo on goratings.org. So 700 elo gap => [98% winrate](#).
3. From [reddit.com/r/baduk/comments/hma3nx/unified_elo_rating_for_ais/](https://www.reddit.com/r/baduk/comments/hma3nx/unified_elo_rating_for_ais/), AlphaZero is 2065 - 1330 = 735 elo weaker than cp505.



Adversarial Policies Beat Superhuman Go AIs

Yoon T. Kim¹, Adam Cheung^{2,3}, Don Tarras⁴, Yoon Young Kwon^{5,6}, Kevin Holme⁷, Sami Shadmehr⁸, Joseph Miller⁹, Michael Hoeser¹⁰, Yoonhan Wu¹¹, Yuhua Peng^{12,13}, Sergey Levine¹⁴, Stuart Bonald¹⁵

Abstract

We attack the state-of-the-art Go-playing AI system KataGo by creating adversarial policies. KataGo is using a supervised learning approach to improve its performance. We use our adversarial policies to train a new Go-playing AI system, which we call AdvKataGo. We show that AdvKataGo can beat KataGo in a head-to-head match. Our adversarial policies are able to find Go moves that KataGo is not able to find. We show that AdvKataGo can beat KataGo in a head-to-head match. Our adversarial policies are able to find Go moves that KataGo is not able to find. We show that AdvKataGo can beat KataGo in a head-to-head match. Our adversarial policies are able to find Go moves that KataGo is not able to find.

1. Introduction

The state-of-the-art performance of Go systems has greatly improved in recent years. The AI system that currently holds the record for performance is KataGo (Silver et al., 2017). KataGo is a Go-playing AI system that uses a supervised learning approach to improve its performance. We use our adversarial policies to train a new Go-playing AI system, which we call AdvKataGo. We show that AdvKataGo can beat KataGo in a head-to-head match. Our adversarial policies are able to find Go moves that KataGo is not able to find. We show that AdvKataGo can beat KataGo in a head-to-head match. Our adversarial policies are able to find Go moves that KataGo is not able to find.

Had Lee Sedol known about this weakness, his challenge match might have turned out very differently. Today I'm going to teach you what this weakness is, how we discovered it, and what implications this has. [click]

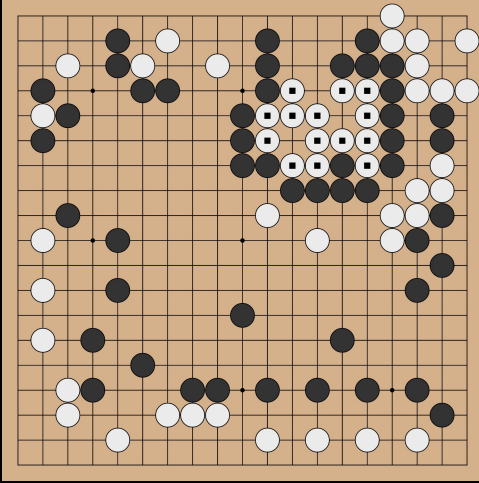
The Cyclic-Exploit

 KataGo AI

 Kellin

First, let's start with the weakness. What you're seeing right now is a game between the superhuman Go AI KataGo, and a member of our research team, Kellin Pelrine ["pell-rin"]. KataGo is playing as black, and Kellin as white. Let's see how Kellin manages to beat a superhuman AI. [click, click]

The Cyclic-Exploit

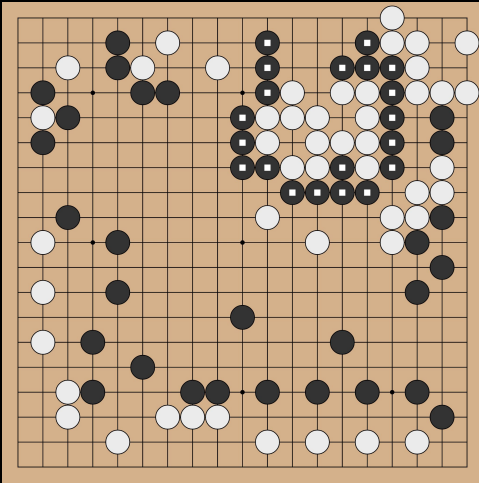


● KataGo AI

● Kellin

At this point, roughly a hundred moves in, Kellin has constructed a small white group in the top right. [click click, pause]

The Cyclic-Exploit

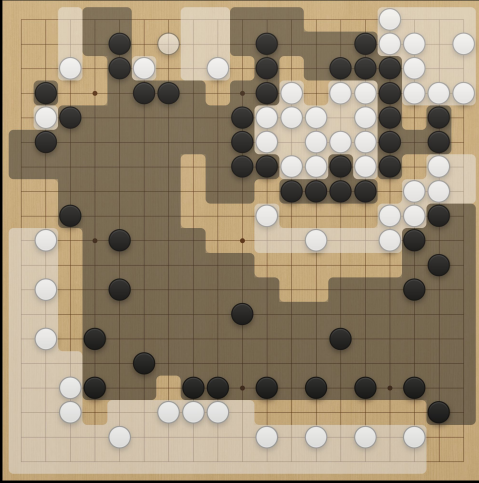


● KataGo AI

● Kellin

KataGo has contained this group with a circle of black stones [click]

The Cyclic-Exploit



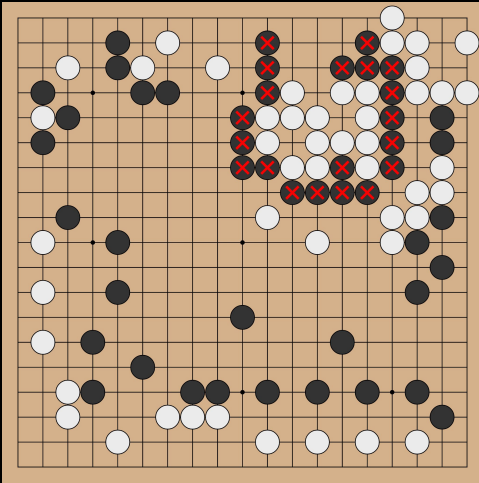
● KataGo AI

● Kellin

This is a losing position for Kellin at the moment. [\[click\]](#)

This is because black controls roughly 50 more squares than white, and in Go, whoever controls more area wins. However, Kellin has other plans. [\[click\]](#)

The Cyclic-Exploit



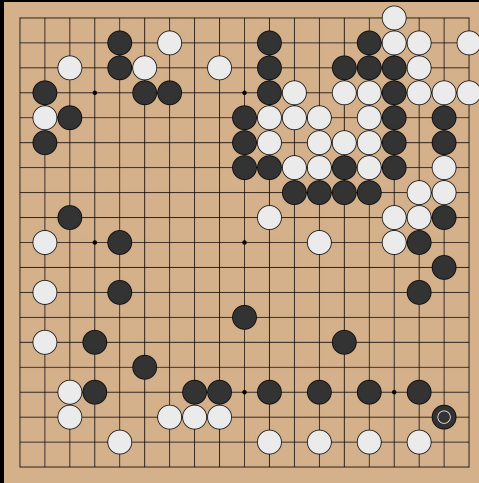
● KataGo AI

● Kellin

Namely, Kellin's plan will be to kill the cyclic-group of black stones that is encircling his white group. He will do this, by slowly re-encircling the black stones from the outside. Let's watch this in action.

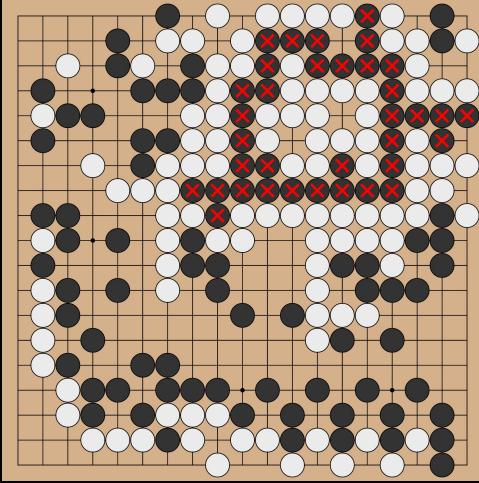
[click]

The Cyclic-Exploit



[pause] KataGo, actually has numerous opportunities to stop this re-encirclement. However it does nothing. This is the hidden weakness of most modern AlphaZero-style Go AIs. When these AIs see a cyclic-group on the board, they think it is invulnerable, even when it is not. By the time KataGo realizes something is wrong, it is too late. Indeed, KataGo resigned in this position, [click, click]

The Cyclic-Exploit

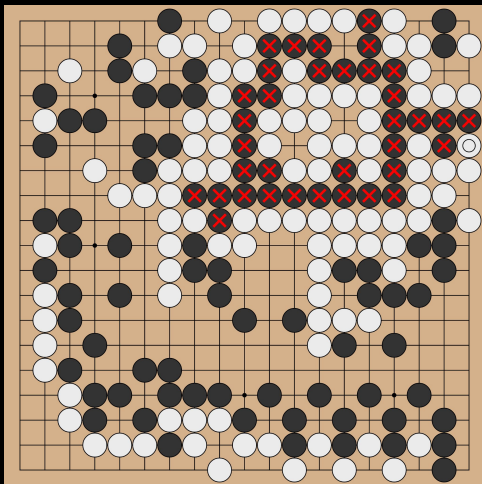


● KataGo AI

● Kellin

Its black cycle is guaranteed to die, as after Kellin plays on the right [click, pause]

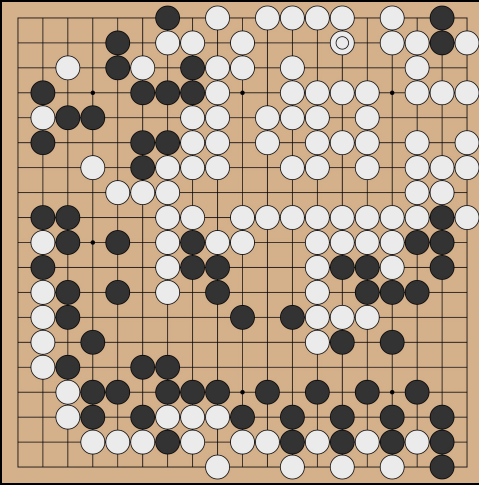
The Cyclic-Exploit



- KataGo AI
- Kellin

and then the top [click, pause]

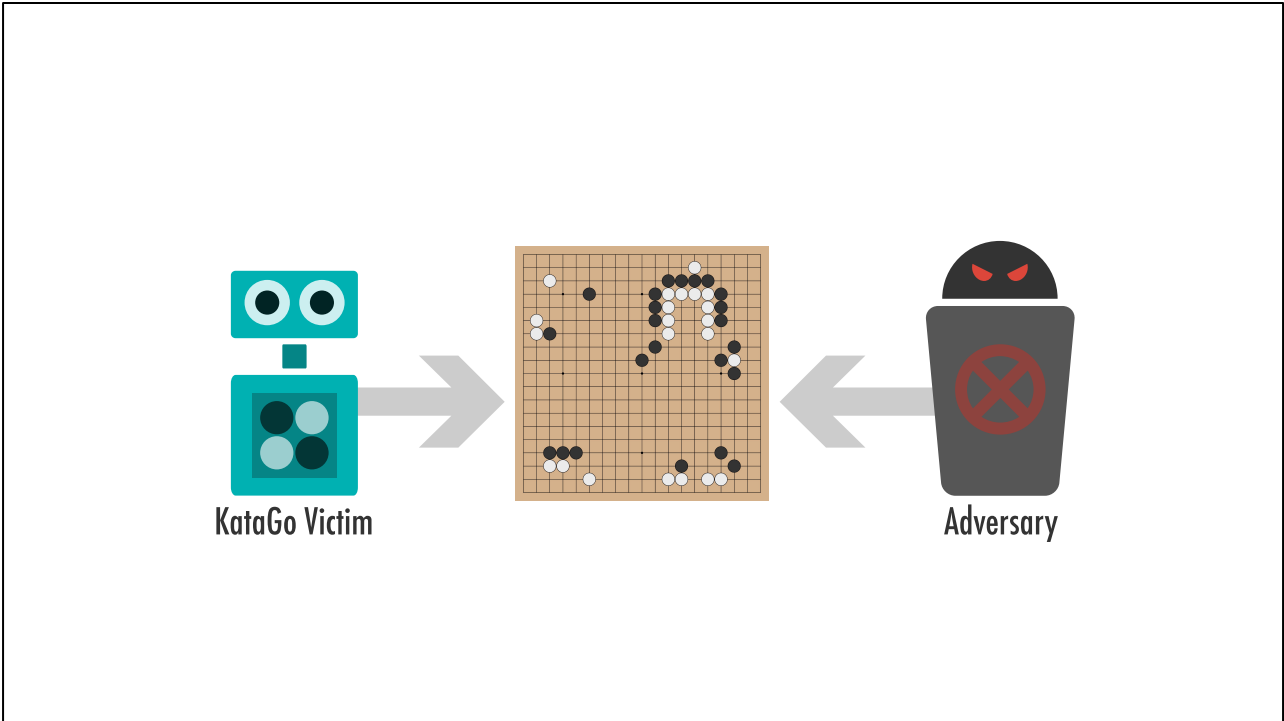
The Cyclic-Exploit



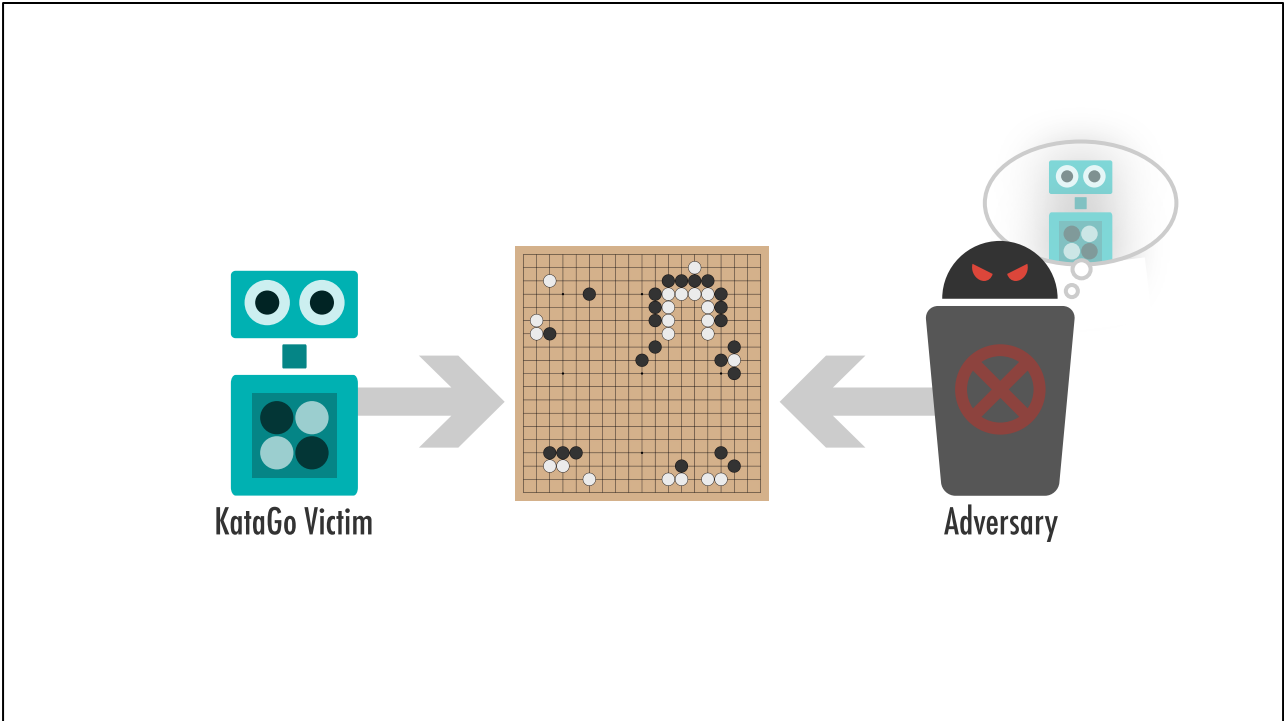
● KataGo AI

● Kellin

the entire group is killed and Kellin now controls the majority of the board.
[click]

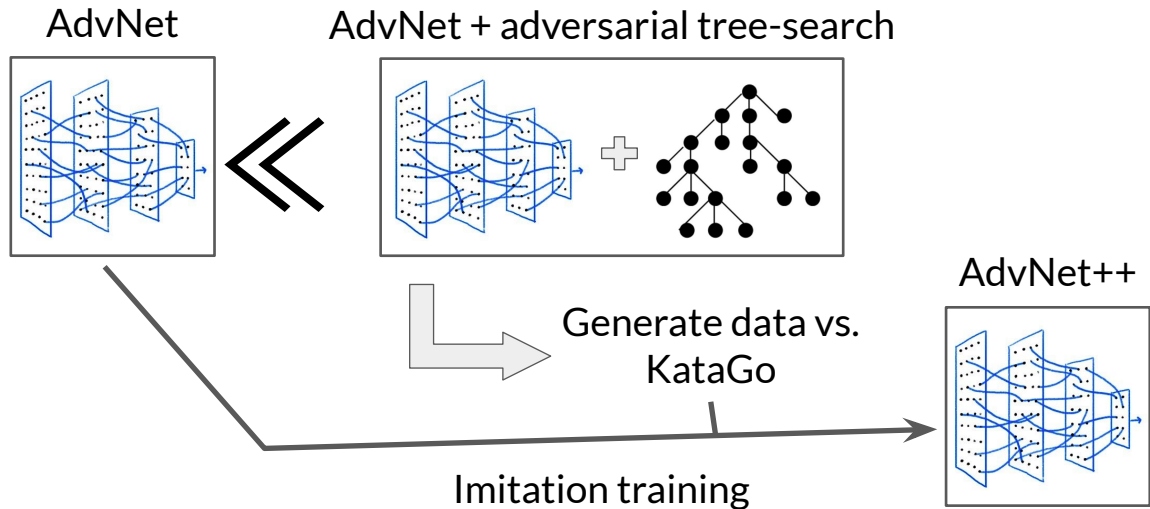


How did we discover this exploit? Well, we did so by training an adversary AI to defeat KataGo. [\[click\]](#)



Our adversary has a special ability. Namely, it can simulate the victim's behavior when it searches over future moves.

Adversarial AlphaZero



To train our adversary, we use an adversarial variant of the AlphaZero algorithm. Here's how it works. [click] We start with a randomly initialized adversary neural network [click] Next, we augment this network with an adversarial variant of Monte-Carlo Tree Search. Tree-search is a policy improvement operator, meaning the network with tree search is a stronger adversary than the network alone. This adversarial tree-search is also where the adversary simulates possible victim responses. [click] We then pit our search-augmented adversary against KataGo, generating a dataset of behavior. [click] Finally, we train the adversary network to mimic the behavior of the search-augmented adversary. This imitation training yields a slightly stronger network.

Repeating this process, we eventually get an adversary that is able to reliably defeat KataGo via the cyclic-exploit I showed previously.

1. Superhuman performance and planning are not sufficient for robustness.
2. Adversarial optimization helps find hidden failure modes.

In summary, we showed that even superhuman Go AIs can have unexpected failure modes. Here are two key takeaways from this result. [click] The first, is that superhuman performance and planning, are not sufficient for robustness. [click] The second, is that adversarial optimization is a very useful technique for finding hidden failure modes, and should be used to improve or validate the robustness of safety-critical systems.

Website: goattack.far.ai

Paper:

arxiv.org/abs/2211.00241



Adversarial Tree Search,
Curriculum Learning,
Transfer, Interpretability,
Defending, and more!



Adam Gleave



Tom Tseng



Kellin Pelrine



Nora Belrose



Joseph Miller



Michael Dennis



Yawen Duan



Viktor Pogrebniak



Sergey Levine



Stuart Russell

For more information, including many more details on our methods and results, check out our website or come find us in person at ICML 2023. Thanks for listening!