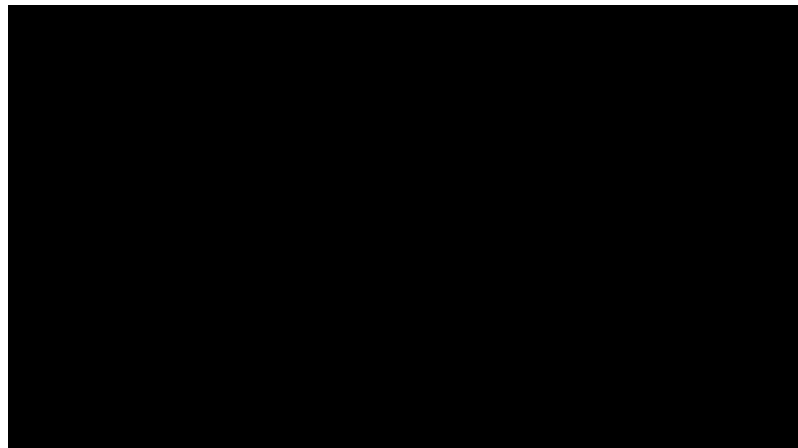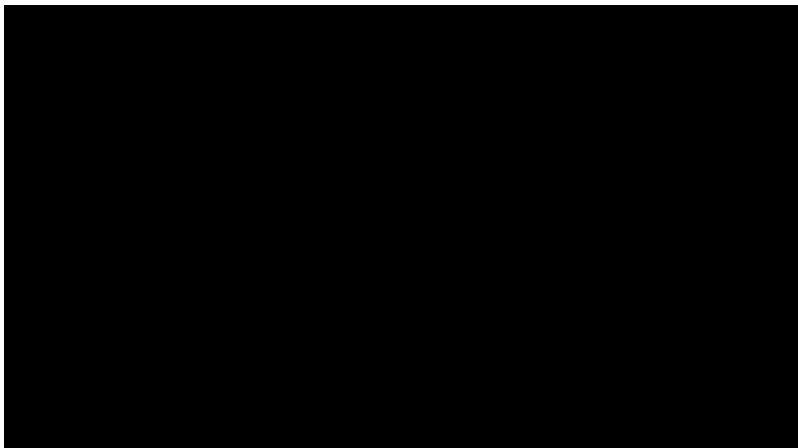# Adversarial Policies: Attacking Deep Reinforcement Learning
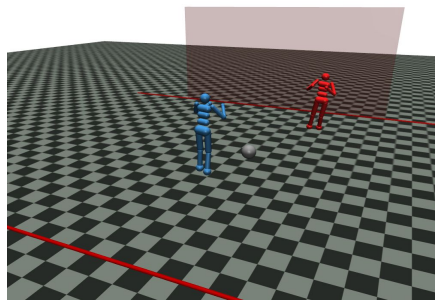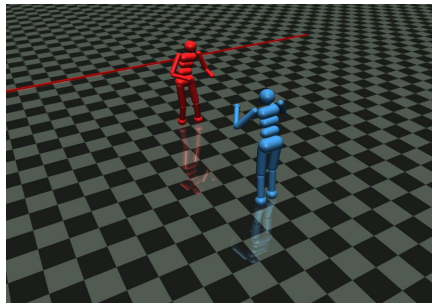
Presentation by Jarek Liesen and Lorenz Hufe
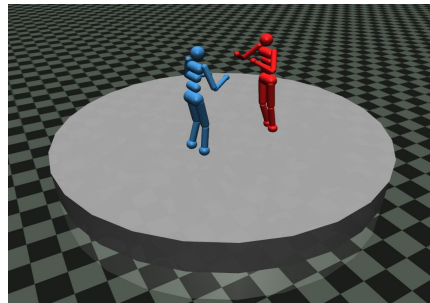
# Video

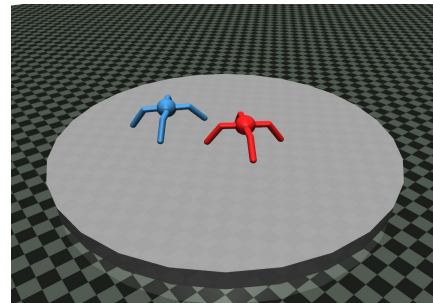# Our environments



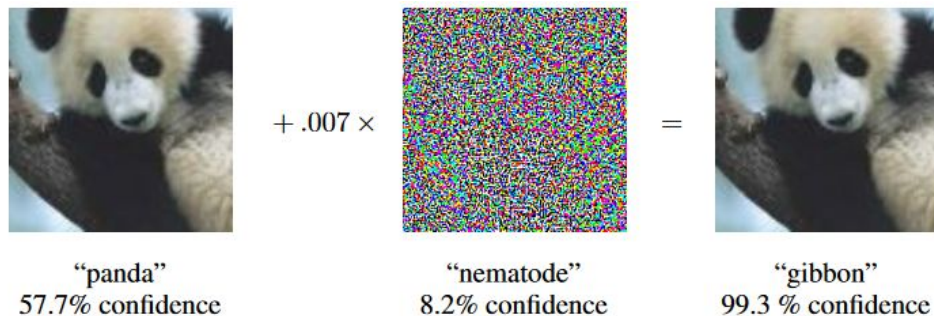Kick and defend      You shall not pass      Sumo humans      Sumo ants

- Blue and red: victim and opponent
- Observation: Each agent observes the configuration of both agents
- Action: The agents can move their joints
- Reward: Zero-sum games

# Adversarial attacks

- Adversarial Attacks [2]: generate an input that maximizes the probability of an error
- How can such an attack be built?



"panda"
57.7% confidence

+ .007 ×

"nematode"
8.2% confidence

=

"gibbon"
99.3 % confidence

| White-box attack (weights known) | Black-box attack (weights unknown) |
|---|---|
| Derive gradient from the loss function w.r.t. the pixel and then do gradient ascend | Different methods (gradient estimation, local search, …) |

# Limitations of adversarial attacks

-   The output of an adversarial attack is generally not observable in real world and therefore only of limited value



As what does the network classify this?

Stop

(a) Normal

Yield    Speed Limit

(b) Attack

-   We want natural observations -> Idea: adversarial policy in a multi agent environment
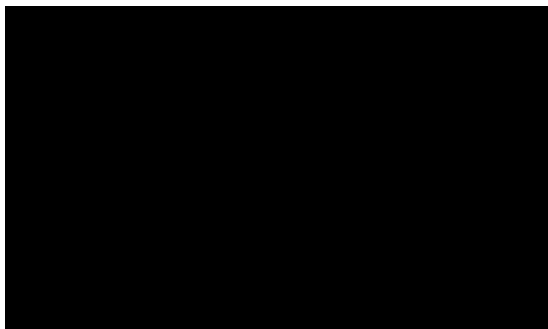
# Adversarial policies

- Playing against one fixed policy turns out to be a pretty good adversarial attack!
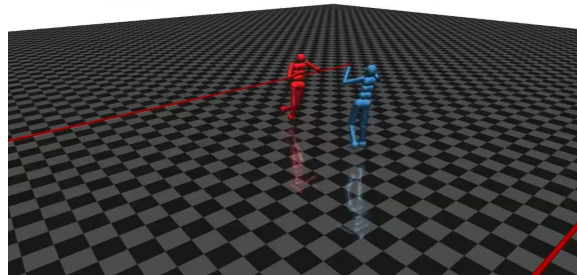- Fixing the other agent makes it single agent

$$R'_\alpha(s, a_\alpha, s') = R_\alpha(s, a_\alpha, a_\nu, s')$$
$$T_\alpha(s, a_\alpha) = T(s, a_\alpha, a_\nu)$$ if $a_v$ fixed

- Victim agents are pretrained agents from Bansal et al [3]
- Trained 20 Million timesteps versus fixed policy
- PPO Implementation by Stable Baselines [4]
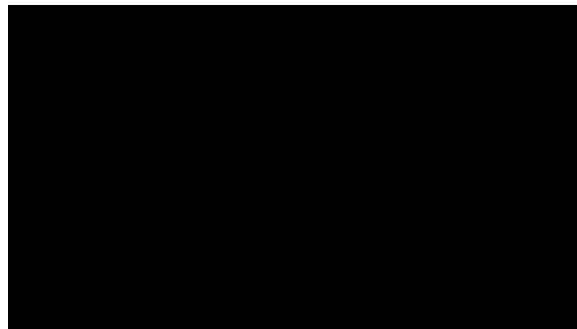- Sparse positive reward on win/negative on loss and ties
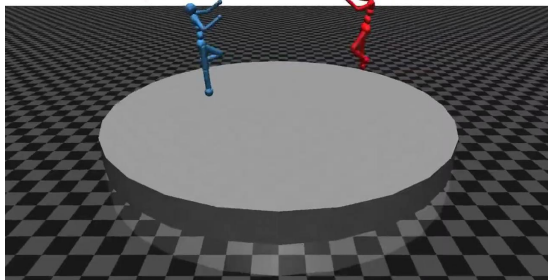
# Result Videos
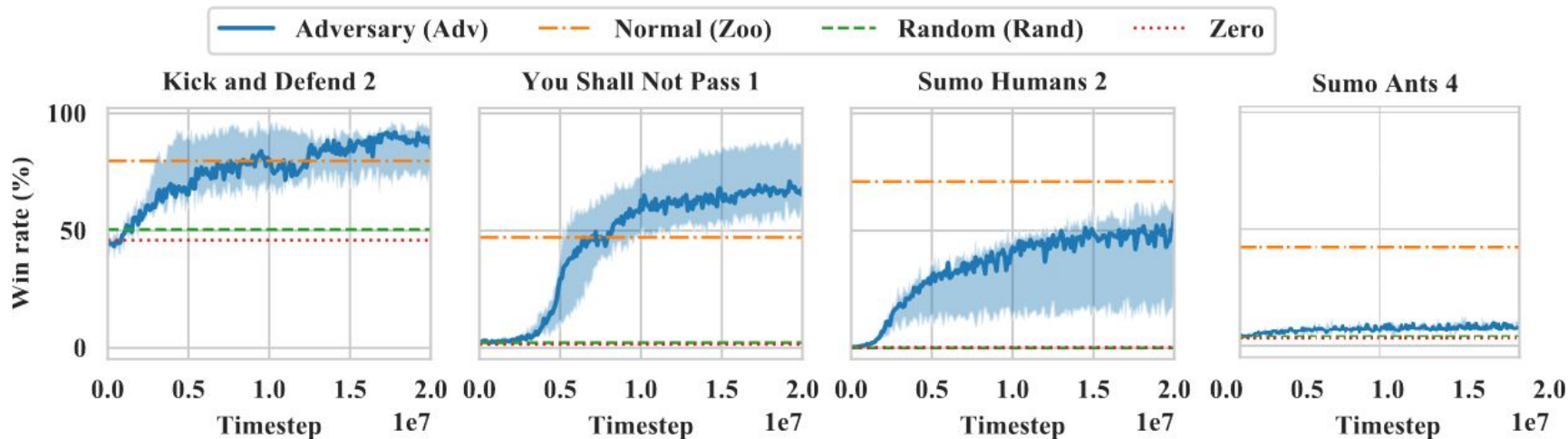


Opponent = 0    Ties = 0    Victim = 0
Adversary (Adv1)            Normal (ZooV1)

Opponent = 0    Ties = 0    Victim = 0
Adversary (Adv2)            Normal (Zoo2)

# Winrate while training



Dimensionality matters:

Adversarial attack works best in high dimensional spaces

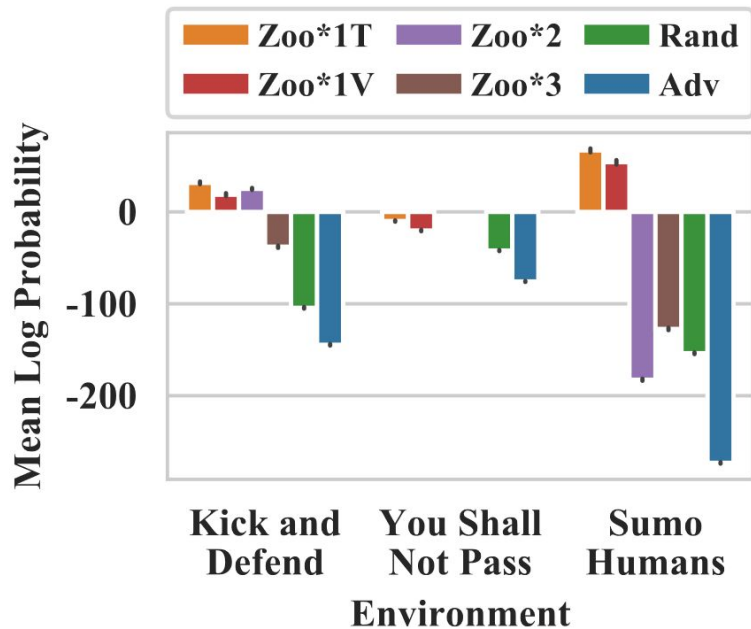    -> Bad results on low dimensional Sumo Ant task

# Winrate against other agents

The adversarial policy performs well against it's training partner, but does not learn a general method to win the game



Win Matrix for Kick and Defend

# Understanding adversarial policies



Gaussian mixture model fit to activity of one victim, mean log probabilities for opponent policies



t-SNE embedding of the victim activations for different opponent policies, victim is ZooV2 for Kick and Defend

# Defending against Adversarial Attacks

- Masking the configuration of the opponent significantly increases its performance (13% -> 99%)

- Fine-tuning against adversarial agents makes victims more robust (13% -> ~90% WR)
- Adversarial attack may be reapplied, but the resulting adversarial policy is conceptually different (1% WR)

# Summary and Outlook

- Agents in multi-agent settings are susceptible to adversarial policies
- The adversarial policies induce abnormal activation patterns in the victims
- Masking and fine-tuning can help against adversarial agents
- It would be interesting to…
  - analyse the adversarial policies after repeated fine-tuning
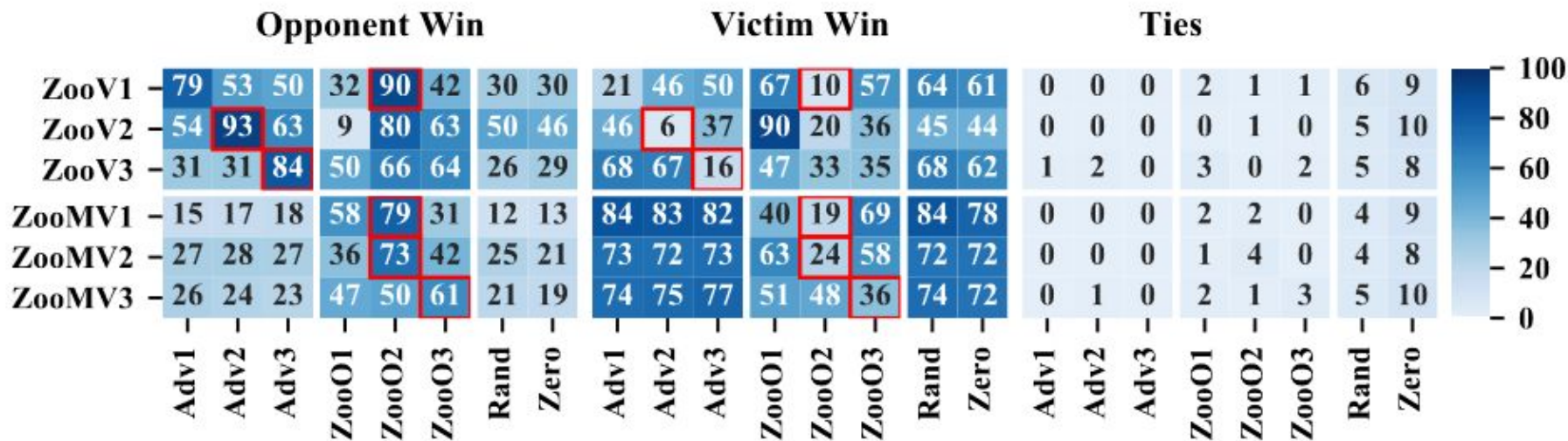  - train policies with adversarial attackers in the pool

# Remarks

- Interesting and direct method for applying adversarial attacks in RL
- Some findings are not fully explained or left out
- Great analysis of the network activations
- Great additional material [Source code, Videos]

# References

[1] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine & S. Russell, *Adversarial Policies: Attacking Deep Reinforcement Learning* in ICLR 2020
https://adversarialpolicies.github.io
[2] I. Goodfellow, J.Shlens & C. Szegedy, *Explaining and Harnessing Adversarial Example* in ICLR 2015
[3] T. Bansal, J. Pachocki, S. Sidor, I. Sutskever & I. Mordatch, *Emergent Complexity via Multi-Agent Competition* in ICLR 2018
[4] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor & Y. Wu. *Stable Baselines*
https://github.com/hill-a/stable-baselines
[5] L. van der Maaten, G. Hinton, *Visualizing Data using t-SNE* in Journal of Machine Learning Research
[6] S. Bhambri, S. Muku, A. Tulasi & A. Buduru *A Survey of Black-Box Adversarial Attacks on Computer Vision Models*
[7] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu *Towards Deep Learning Models Resistant to Adversarial Attacks*

# Appendix A:



(a) Kick and Defend

# Appendix B