

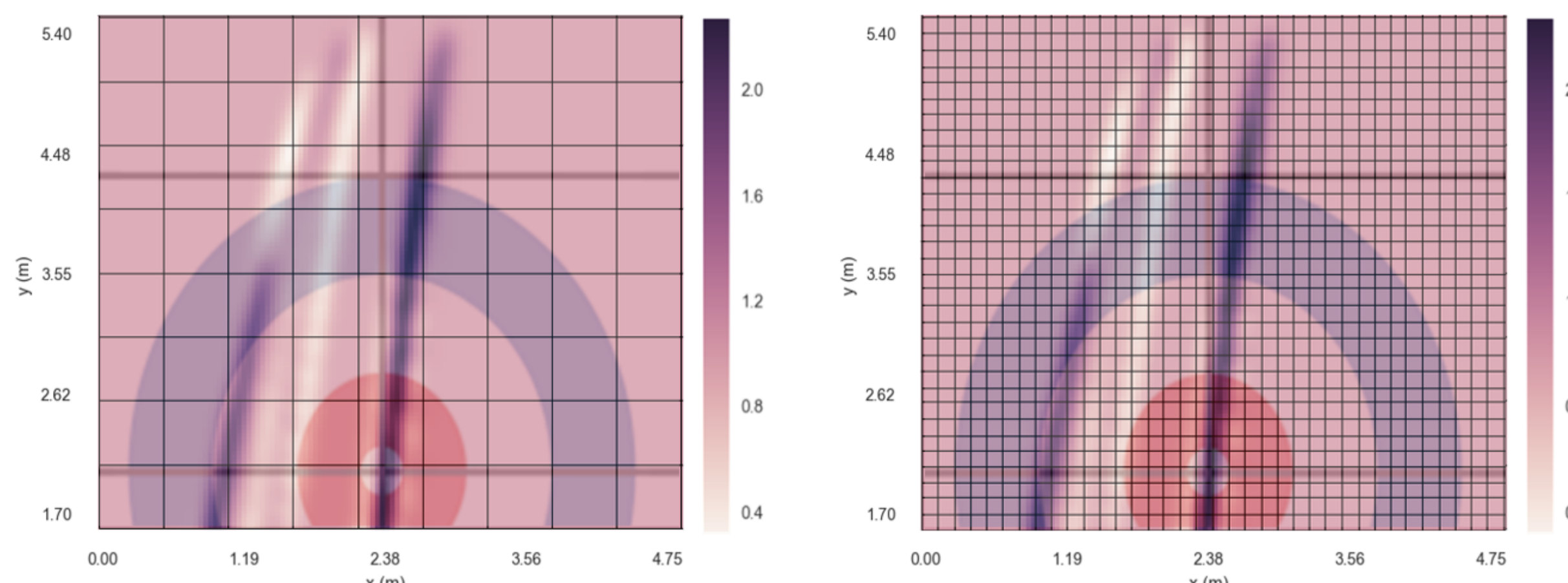
Introduction - Deep Reinforcement Learning in Continuous Action Space

- ▶ **Goal:** Training an **agent** to learn a improved policy in the **continuous space**
- ▶ **Input:**
 - ▶ Current **state** from the environment
 - ▶ **Reward** from the environment
- ▶ **Output:** The **best action** given constraints
- ▶ **Case study in the game of simulated curling:**
 - ▶ Two dimensional **continuous action space** with two kinds of curl directions
 - ▶ **Execution uncertainty** modeled by asymmetric Gaussian noise.

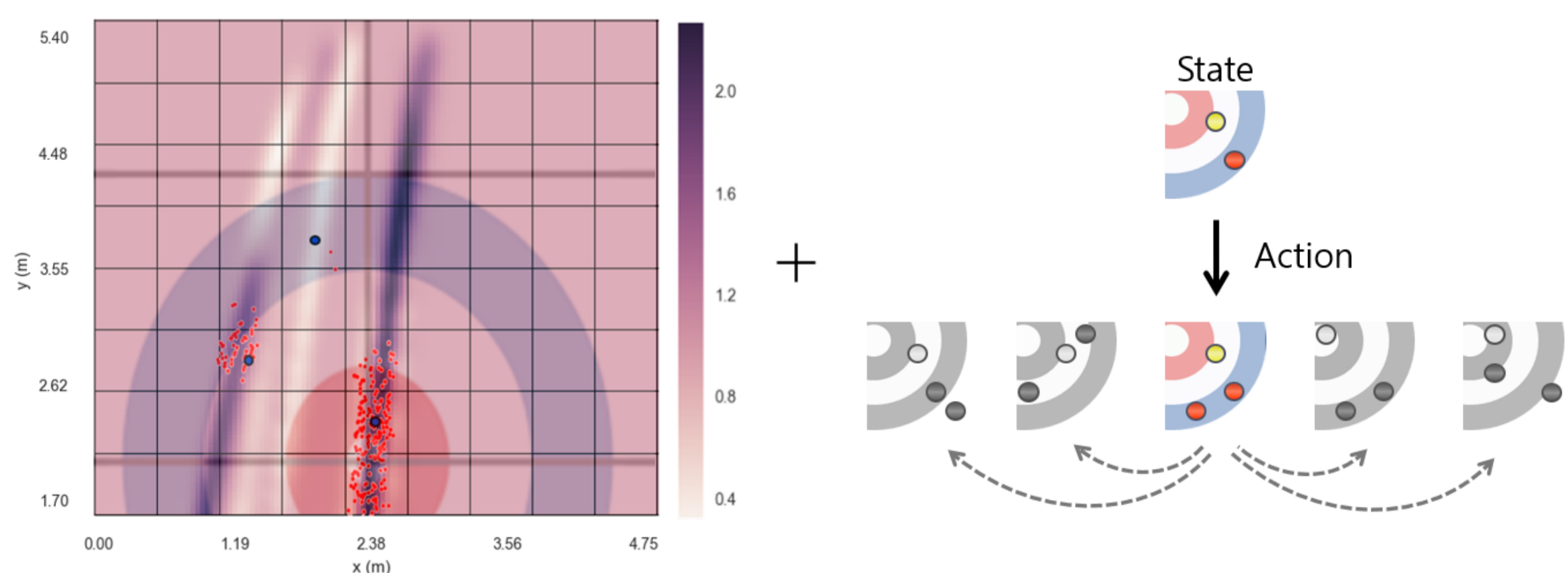


Motivation

- ▶ Deep neural networks for the discrete actions are not suitable for devising strategies for games in which a very small change in an action can dramatically affect the outcome.
- ▶ **Challenges of learning in the discretized action**



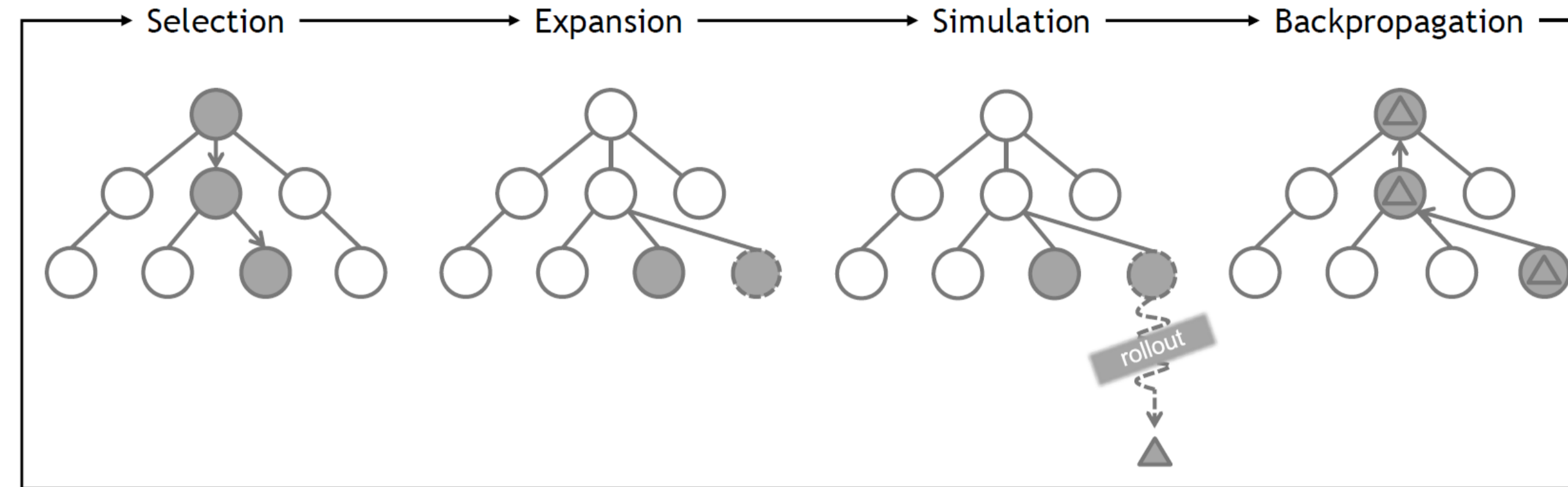
- ▶ Deterministic discretization has problems: (1) low resolution → **strong bias** in policy evaluation and improvement (2) high resolution → **slow searching and learning speed** and **exponential growth** in the number of actions to explore
- ▶ **Learning in the discretized action space with Kernel Meth-**



- ▶ Conducts **local search** with **continuous** action samples generated from a deep convolutional neural network (CNN).
- ▶ **Generalizes the information** between similar actions through kernel methods.

Monte Carlo Tree Search and Kernel Regression

- ▶ **Monte Carlo Tree Search (MCTS)** is a simulation-based search approach to planning in finite-horizon sequential decision-making settings.



- ▶ **Upper Confidence Bound applied to Trees (UCT)** is a commonly used MCTS algorithm using an Upper Confidence Bound (UCB) selection function.

$$\operatorname{argmax}_a \bar{v}_a + C \sqrt{\frac{\log \sum_b n_b}{n_a}}$$

- ▶ **Kernel Regression** is a non-parametric method which uses a kernel function as a weight for estimating the conditional expectation of a random variable.

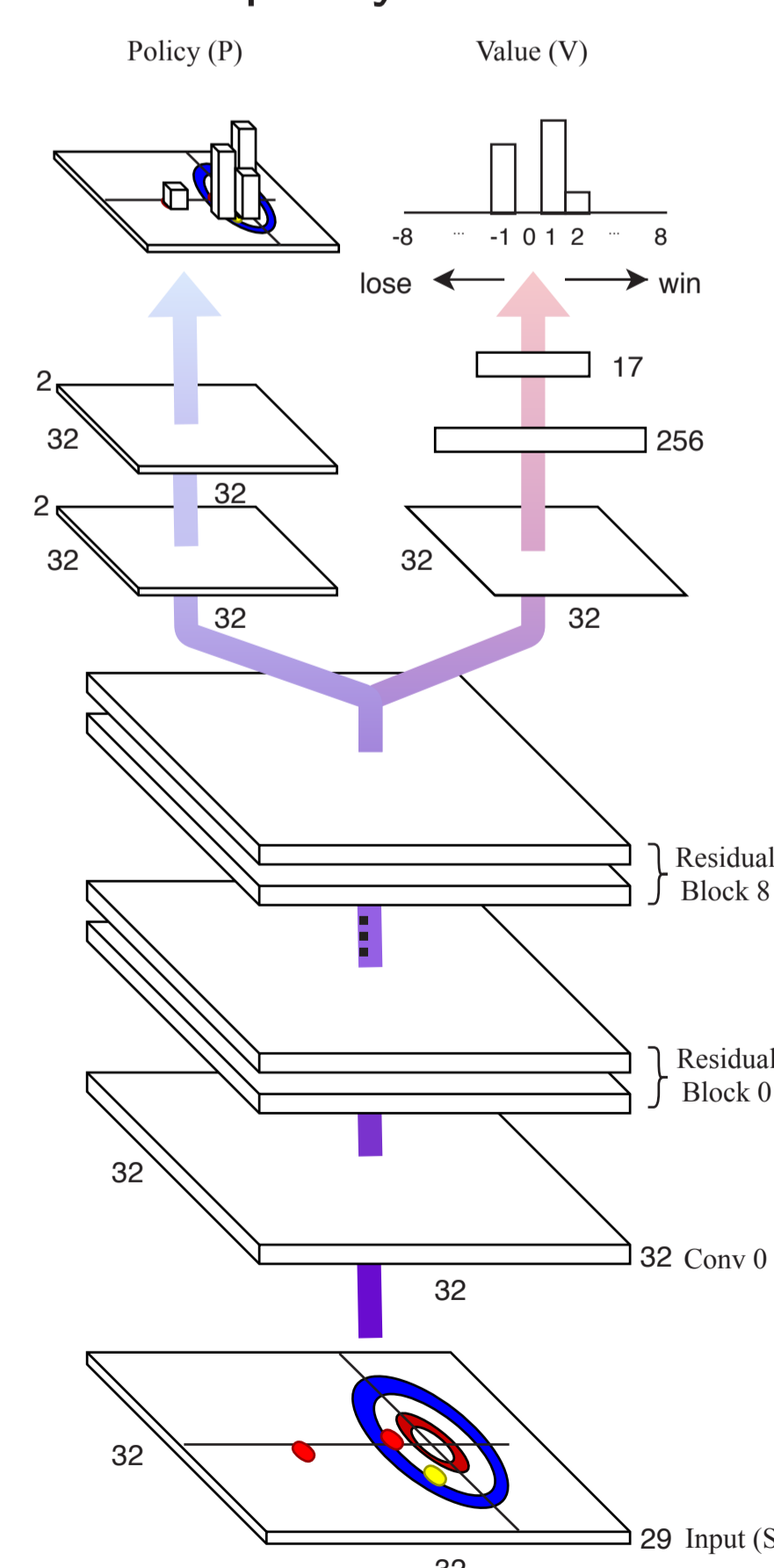
$$E[y|x] = \frac{\sum_{i=0}^n K(x, x_i) y_i}{\sum_{i=0}^n K(x, x_i)}$$

- ▶ The denominator of kernel regression is related to **Kernel Density Estimation** which is method for estimating the probability density function of random variable.

$$W(x) = \sum_{i=0}^n K(x, x_i)$$

Kernel Regression Deep Learning UCT

- ▶ The policy-value network



Algorithm 1 KR-DL-UCT

- 1: $\mathbf{p}_\theta \leftarrow$ the policy network
- 2: $\mathbf{v}_\theta \leftarrow$ the value network
- 3: $s_t \leftarrow$ the current state
- 4: $A_t \leftarrow$ a set of visited actions in s_t
- 5: $expanded \leftarrow$ false
- 6: **if** s_t is terminal **then**
- 7: **return** Score(s_t), false
- 8: **end if**
- 9: $a_t \leftarrow \operatorname{argmax}_{a \in A_t} \mathbb{E}[\bar{v}_a | a] + C \sqrt{\frac{\log \sum_{b \in A_t} W(b)}{W(a)}} \triangleright$ Selection
- 10: **if** $\sqrt{\sum_{a \in A_t} n_a} < |A_t|$ **then**
- 11: $s_{t+1} \leftarrow$ TakeAction(s_t, a_t)
- 12: $reward, expanded \leftarrow$ KR-DL-UCT(s_t)
- 13: **end if**
- 14: **if not expanded then**
- 15: $a'_t \leftarrow \operatorname{argmin}_{K(a_t, a) > \gamma} W(a) \triangleright$ Expansion
- 16: $A_t \leftarrow A_t \cup a'_t$
- 17: $s_{t+1} \leftarrow$ TakeAction(s_t, a'_t)
- 18: $A_{t+1} \leftarrow \cup_{j=1}^k \{a_{init}^{(j)}\}$ s.t. $a_{init}^{(j)} \sim \pi_{a|s_{t+1}}$ // Policy net
- 19: $reward \leftarrow \mathbf{v}_\theta(s_{t+1} | s_t, a'_t)$ // Value net \triangleright Simulation
- 20: **end if**
- 21: $\bar{v}_{a_t} \leftarrow \frac{1}{n_{a_t} + 1} (n_{a_t} \bar{v}_{a_t} + reward) \triangleright$ Backpropagation
- 22: $n_{a_t} \leftarrow n_{a_t} + 1$
- 23: **return** $reward$, true

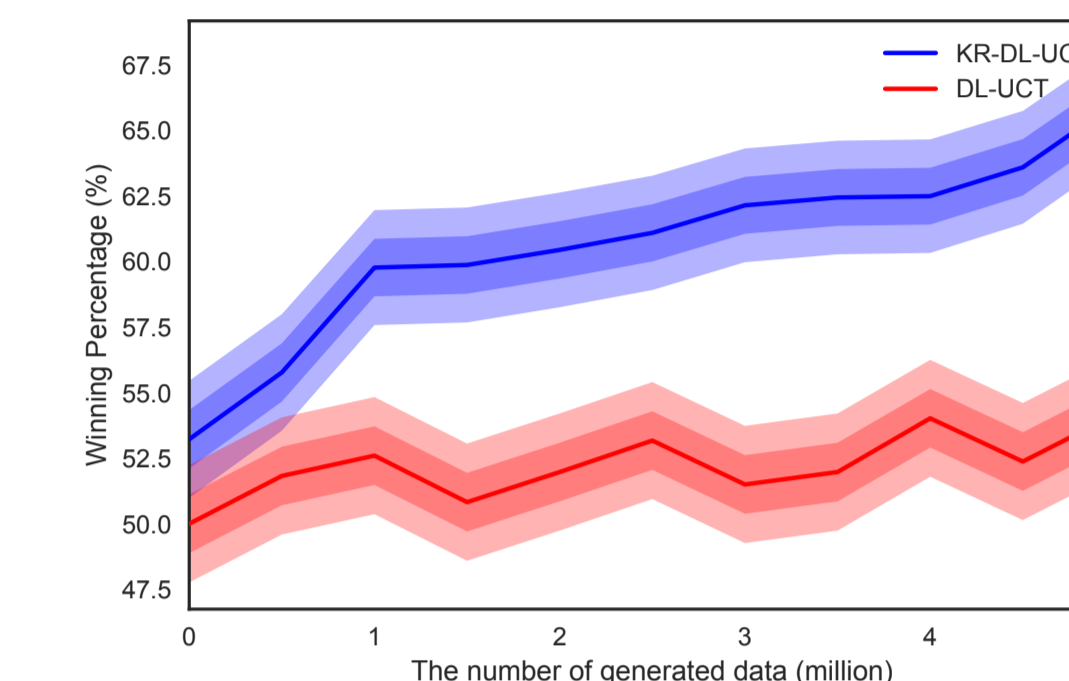
- ▶ Source codes will be available at <https://github.com/leekwoon/KR-DL-UCT>

Experimental Results

Datasets

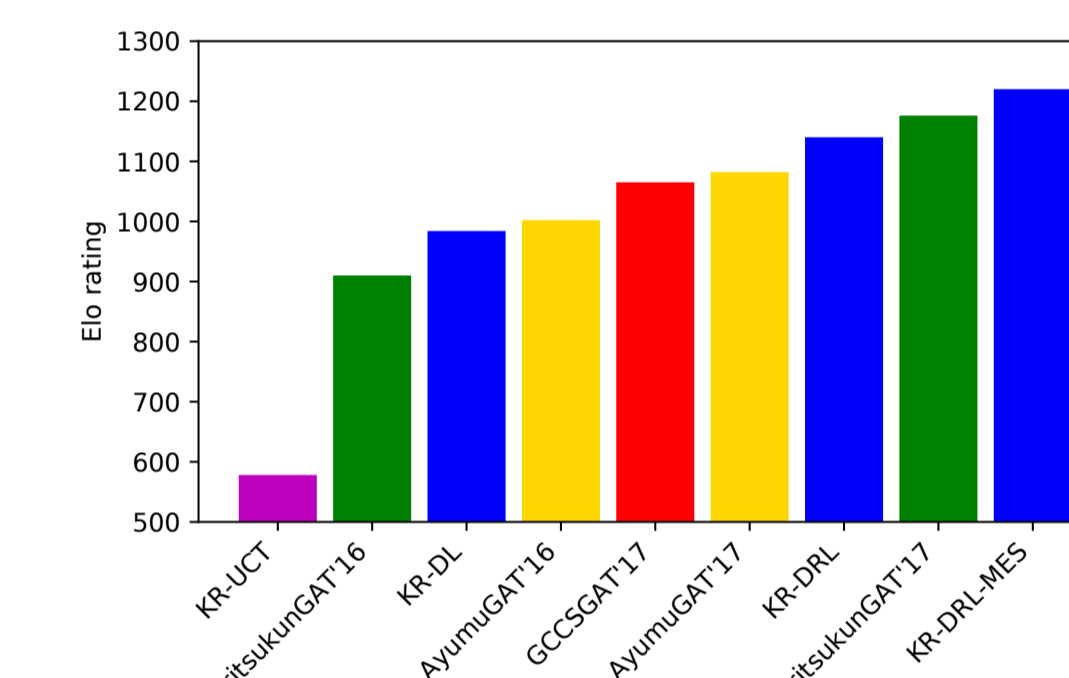
- ▶ **Supervised learning:** 0.4 million of the play data from the champion program (AyumuGAT'16) of Game AI Tournaments (GAT) digital curling championship in 2016.
- ▶ **Self-play Reinforcement Learning:** 5 million of the play data from self-play matches of KR-DL-UCT, executing 400 simulations per move.

Quantitative Results



# OF DATA (MILLION)	DL-UCT (1)	KR-DL-UCT (2)	(2)-(1)
0.0	50.0%	53.2%	3.2%
1.0	52.6%	60.0%	7.2%
2.0	51.9%	60.5%	8.5%
3.0	51.5%	62.2%	10.7%
4.0	54.0%	62.5%	8.5%
5.0	54.2%	66.0%	11.9%

- ▶ KR-DL-UCT (blue) outperforms (53.23%) DL-UCT (red) even without the self-play RL.
- ▶ After gathering 5 million shots from self-play, KR-DL-UCT wins 66.05% which is significantly higher than DL-UCT case.



PROGRAM	WINNING PERCENTAGE
GCCSGAT'17	74.0 ± 6.22%
AYUMUGAT'16	66.5 ± 6.69%
AYUMUGAT'17	62.3 ± 6.87%
JIRITSUKUNGAT'16	86.3 ± 4.88%
JIRITSUKUNGAT'17	55.5 ± 7.04%

- ▶ KR-DL : KR-DL-UCT with supervised learning
- ▶ KR-DRL : KR-DL with self-play RL
- ▶ KR-DRL-MES: KR-DRL with winning percentage table (multi-end strategy)
- ▶ **Our program KR-DRL-MES won in the international digital curling competition, GAT-2018.**

Conclusion

- ▶ We provide a new framework which incorporates a neural network for learning strategy with a kernel based Monte Carlo tree search in the continuous action space.
- ▶ The developed method is applied to the game of Simulated Curling and achieves the state-of-the-art performance.

Acknowledgements

- ▶ This work was supported by Institute for Information and communications Technology Promotion (IITP) grant funded by the Korea government (MSIT) (No.2017-0-01779, A machine learning and statistical inference framework for explainable artificial intelligence, and No.2017-0-00521, AI curling robot which can establish game strategies and perform games).

References

- ▶ Yee, T., Lisý, V., and Bowling, M. Monte carlo tree search in continuous action spaces with execution uncertainty. In Proceedings of the International Joint Conference on Artificial Intelligence, IJCAI, pp. 690697, 2016.
- ▶ Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., van den Driessche, G., Graepel, T., and Hassabis, D. Mastering the game of go without human knowledge. Nature, pp. 354359, 2017.
- ▶ Ito, T. and Kitasei, Y. Proposal and implementation of digital curling. In Proceedings of the IEEE Conference on Computational Intelligence and Games, CIG, pp. 469-473, 2015.