

TD3 for Competitive Air Hockey

RL Course WS 2025/26 — Final Project

Julian Jurcevic · Team: alphabet-td3

February 2026

1. Clipped Double Q-Learning

Two critics; target uses the *minimum*:

$$y = r + \gamma(1-d) \min_{i=1,2} Q_{\phi'_i}(s', \tilde{a}')$$

2. Target Policy Smoothing

$$\tilde{a}' = \text{clip}(\pi_{\theta'}(s') + \epsilon, -c, c)$$

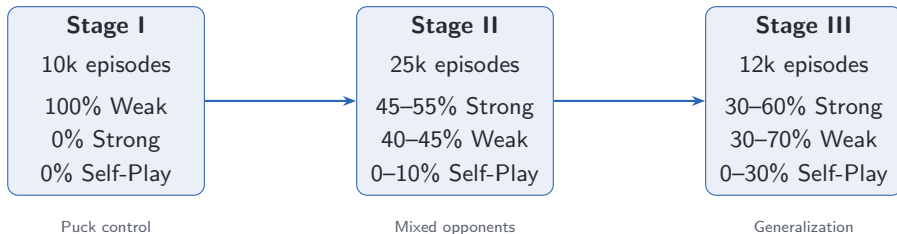
3. Delayed Actor Updates

- Critic: every step; Actor: every 2nd
- Polyak averaging:
 $\phi' \leftarrow \tau \phi + (1-\tau) \phi'$

Architecture

- Actor & twin critics: 2×256 (tanh)
- Actor: $a \in [-1, 1]^4$
- Critics: concatenated (s, a) input

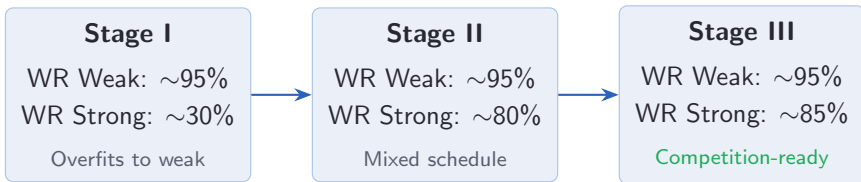
Three-Stage Curriculum with Self-Play



Self-Play Pool: Snapshots every $k=150$ eps · Pool $N_{\text{pool}}=25$ · Difficulty-weighted sampling ($\times 1.2$ loss, $\times 0.95$ win)

Curriculum Training Progression

Training curves for a representative single-seed run (see report Figure 1).



- Without curriculum: strong win rate stays at $\sim 30\%$ (pure weak training)
- Staged scheduling resolves this while retaining high weak win rates

Ablation: Noise Comparison (3 seeds, Stage II)

Noise Type	WR Weak (%)	WR Strong (%)	Ret. Weak	Ret. Strong
Gaussian	92.5 \pm 4.5	81.0 \pm 0.5	8.22 \pm 0.80	5.69 \pm 0.10
Ornstein–Uhlenbeck	94.7 \pm 0.6	89.0 \pm 2.7	8.56 \pm 0.13	7.06 \pm 0.46
Pink	92.6 \pm 4.3	86.1 \pm 2.8	8.17 \pm 0.57	6.40 \pm 0.34
Uniform	91.2 \pm 2.8	80.2 \pm 8.4	8.10 \pm 0.55	5.51 \pm 1.51

- OU best: **+8%** vs. Gaussian against strong (temporal correlation \rightarrow smoother trajectories)
- Pink noise second-best (also correlated) · Uniform: highest variance

Ablation: Self-Play & Prioritized Replay (3 seeds)

Variant	WR Weak (%)	WR Strong (%)	Ret. Weak	Ret. Strong
No PER, No SP	93.1 \pm 3.8	78.3 \pm 3.1	8.33 \pm 0.66	5.00 \pm 0.70
No PER, Self-Play	90.7 \pm 5.9	72.6 \pm 7.6	7.62 \pm 1.26	4.06 \pm 1.56
PER, No SP	75.8 \pm 9.2	66.1 \pm 4.7	4.22 \pm 2.00	1.99 \pm 1.04
PER, Self-Play	78.3 \pm 2.2	65.3 \pm 5.1	4.71 \pm 0.54	1.78 \pm 1.01

PER hurts performance

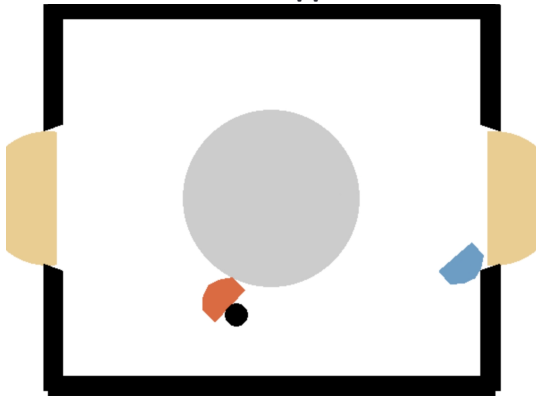
- Non-stationary opponents amplify priority variance
- \approx 15–20% WR drop \Rightarrow **not used**

Self-Play — trade-off

- Lower benchmark scores but **retained** for tournament robustness
- Pool diversity decreases as policies converge

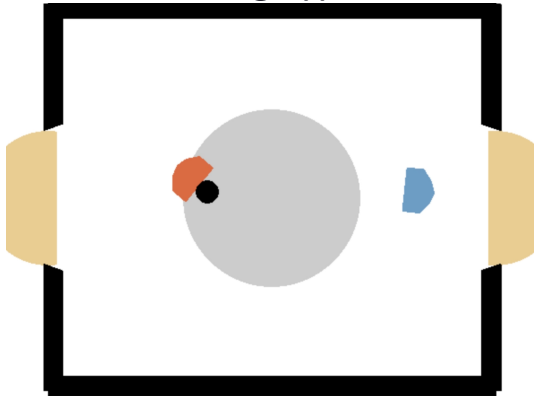
Gameplay Examples

vs. Weak Opponent



File: gameplay_weak.mov

vs. Strong Opponent



File: gameplay_strong.mov

Conclusion & Takeaways

What worked

- **Curriculum**: most impactful; prevents overfitting to single opponent
- **OU noise**: +8% WR (strong) via temporal correlation
- **Self-play**: retained for tournament generalization

What didn't work

- **PER**: ~15% drop under non-stationarity

Final Agent Config

Algorithm	TD3
Noise	OU (annealed)
Replay	Uniform (300k)
Curriculum	3 stages
Self-Play	Pool of 25
WR Weak	~95%
WR Strong	~85%

Limitations

- Single-seed training curves
- Manual curriculum tuning
- Self-play pool convergence