

Rollout policy

SL policy network

RL policy network

Value network

p_{π}

p_{σ}

p_{ρ}

v_{θ}



Policy gradient

Classification

Classification

Self Play

Regression

Human expert positions

Self-play positions

Neural network

Data