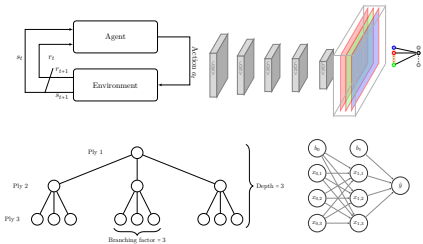
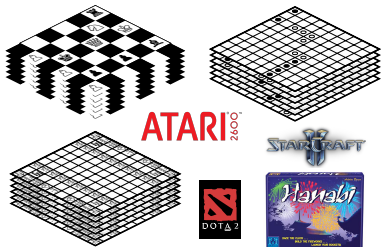


WHY GAMES

AI & Machine Learning



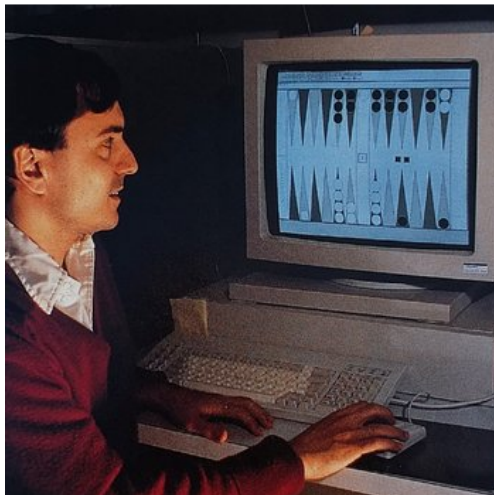
Games



Why Combine the two?

- ▶ AI & Games have a long history (Turing '50& Minsky 60')
- ▶ Simple to evaluate, reproducible, controllable, quick feedback loop
- ▶ Common benchmark for the research community

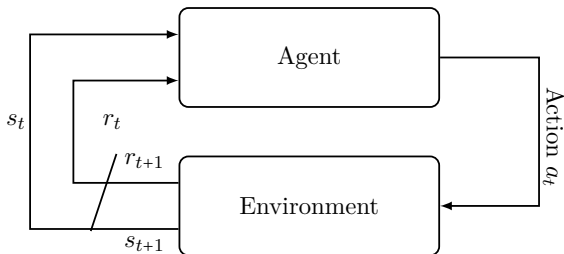
1992: TESAURO'S TD-GAMMON²



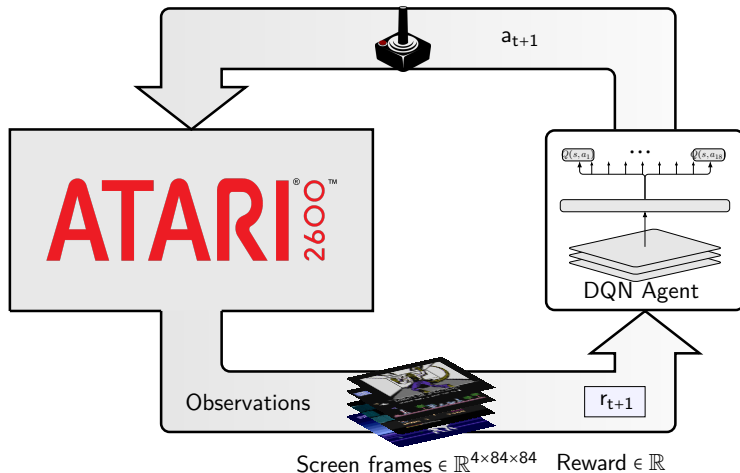
2em1² Gerald Tesauro. "TD-Gammon, a Self-Teaching Backgammon Program, Achieves Master-Level Play". In: *Neural Comput.* 6.2 (Mar. 1994), 215–219. ISSN: 0899-7667. DOI: 10.1162/neco.1994.6.2.215. URL: <https://doi.org/10.1162/neco.1994.6.2.215>

THE REINFORCEMENT LEARNING PROBLEM

- ▶ Notation; **policy**: π , **state**: s , **reward**: r , **action**: a
- ▶ Agent's goal: **maximize reward**, $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$ $0 \leq \gamma \leq 1$
- ▶ RL's goal, **find optimal policy** $\pi^* = \max_{\pi} \mathbb{E}[R|\pi]$



RL EXAMPLES: ATARI (Mnih '15)⁸



⁸2em1 Volodymyr Mnih et al. “Human-level control through deep reinforcement learning”. In: *Nature* 518.7540 (Feb. 2015), pp. 529–533. ISSN: 00280836. URL: <http://dx.doi.org/10.1038/nature14236>.

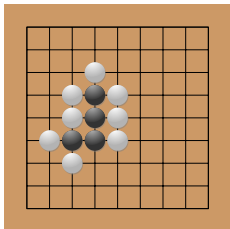
HOW TO ACT OPTIMALLY? (BELLMAN 57'¹⁰)

$$\begin{aligned}
 \text{optimal}(s_t) &= \max_{\pi} \mathbb{E} \left[\sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k} \mid s_t \right] \\
 &= \max_{\pi} \mathbb{E} \left[r_{t+1} \sum_{k=2}^{\infty} \gamma^{k-1} r_{t+k} \mid s_t \right] \\
 &= \max_{a_t} \mathbb{E} \left[r_{t+1} + \max_{\pi} \mathbb{E} \left[\sum_{k=2}^{\infty} \gamma^{k-1} r_{t+k} \mid s_{t+1} \right] \mid s_t \right] \\
 &= \max_{a_t} \mathbb{E} \left[r_{t+1} + \gamma \max_{\pi} \mathbb{E} \left[\sum_{k=2}^{\infty} \gamma^{k-2} r_{t+k} \mid s_{t+1} \right] \mid s_t \right] \\
 &= \max_{a_t} \mathbb{E} \left[r_{t+1} + \gamma \max_{\pi} \mathbb{E} \left[\sum_{k=2}^{\infty} \gamma^{k-2} r_{t+k} \mid s_{t+1} \right] \mid s_t \right] \\
 &= \max_{a_t} \mathbb{E} [r_{t+1} + \gamma \text{optimal}(s_{t+1}) \mid s_t]
 \end{aligned}$$

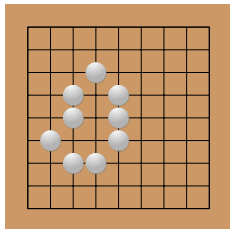
THE GAME OF GO

- ▶ The world's oldest game: 3000 years old, over 40M players world wide
- ▶ To win: **capture** the most territory on the board
 - ▶ Surrounded stones are captured and removed
- ▶ Why is it so hard for computers? 10^{170} **unique states!!**
 ≈ 250 **branching factor**
- ▶ High branching factor, large board (19×19), hard to evaluate etc..

(1)

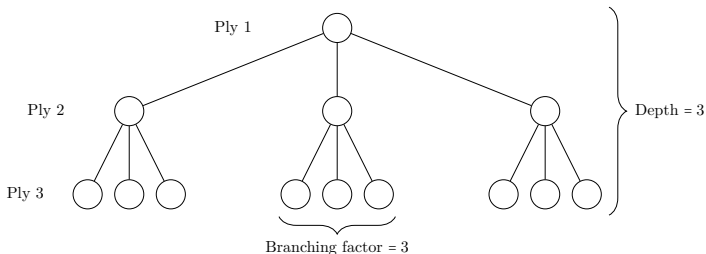


(2)

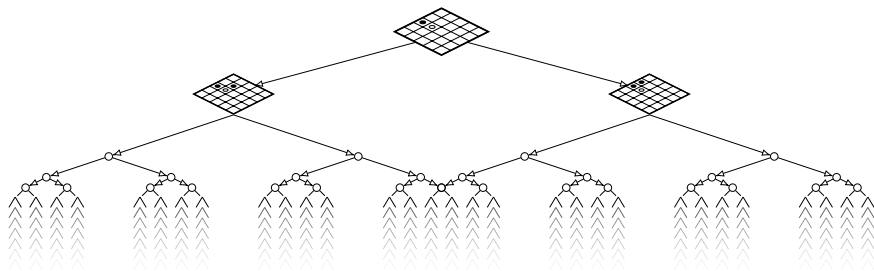


GAME TREES

- ▶ How do you program a computer to play a board game?
- ▶ Simplest approach:
 - ▶ (1) Program a game tree; (2) Assume opponent think like you; (3) Look-ahead and evaluate each move
 - ▶ Requires Knowledge of game rules and evaluation function

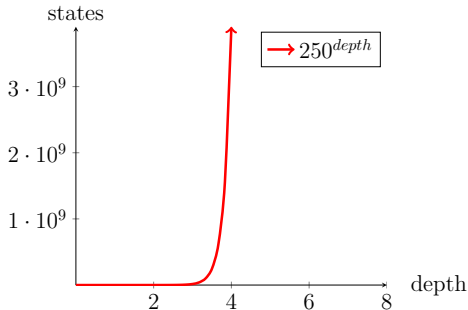


SEARCH + GO = 



SOME NUMBERS

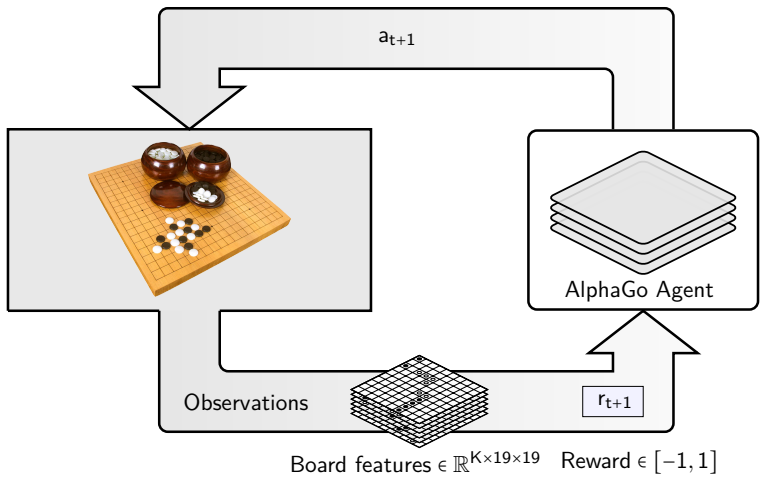
- ▶ **Atoms in the universe**
 - ▶ $\approx 10^{80}$
- ▶ **States**
 - ▶ Go: 10^{170} , Chess: 10^{47}
- ▶ **Game tree complexity**
 - ▶ Go: 10^{360} , Chess: 10^{123}
- ▶ **Average branching factor**
 - ▶ Go: 250, Chess: 35
- ▶ **Board size (positions)**
 - ▶ Go: 361, Chess: 64



ALPHA GO'S APPROACH

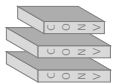
- ▶ Brute-Force Search does not work
 - ▶ At least not until hardware has improved **a lot**.
- ▶ Human Go professionals rely on small search guided by intuition/experience
- ▶ **AlphaGo's Approach:** Complement MCTS with “artificial intuition”
 - ▶ Artificial intuition provided by two neural networks: value network and policy network

COMPUTER GO AS AN RL PROBLEM

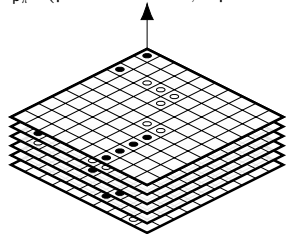


ALPHAGO TRAINING PIPELINE (1/2)

Supervised **Rollout** Policy Network
 $p_{\pi}(a|s)$



Classification
 $\min_{p_{\pi}} L(\text{predicted move, expert move})$

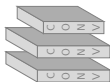


Human Expert Moves D_1

ALPHA Go TRAINING PIPELINE (1/2)

Supervised **Rollout** Policy Network

$$p_{\pi}(a|s)$$

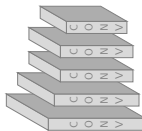


Classification

$$\min_{p_{\pi}} L(\text{predicted move, expert move})$$

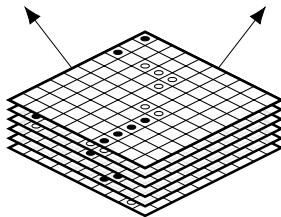
Supervised Policy Network

$$p_{\sigma}(a|s)$$



Classification

$$\min_{p_{\sigma}} L(\text{predicted move, expert move})$$



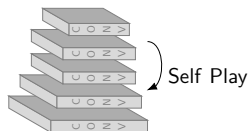
Human Expert Moves D_1

ALPHAGO TRAINING PIPELINE (2/2)

Reinforcement Learning Policy Network

$$p_{\rho}(a|s)$$

Initialize with p_{σ} weights



PolicyGradient

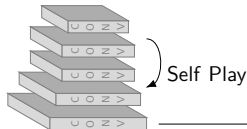
$$J(\mathbf{p}_{\rho}) = \mathbb{E}_{\mathbf{p}_{\rho}} [\sum_{t=0}^{\infty} r_t]$$

$$\rho \leftarrow \rho + \alpha \nabla_{\rho} J(\mathbf{p}_{\rho})$$

ALPHA Go TRAINING PIPELINE (2/2)

Reinforcement Learning Policy Network

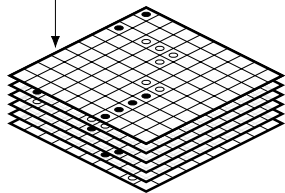
$$p_{\rho}(a|s)$$

Initialize with p_{σ} weights

PolicyGradient

$$J(p_{\rho}) = \mathbb{E}_{p_{\rho}} [\sum_{t=0}^{\infty} r_t]$$

$$\rho \leftarrow \rho + \alpha \nabla_{\rho} J(p_{\rho})$$

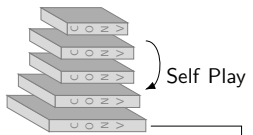
Self Play Dataset D_2

ALPHAGO TRAINING PIPELINE (2/2)

Reinforcement Learning Policy Network

$$p_{\rho}(a|s)$$

Initialize with p_{σ} weights



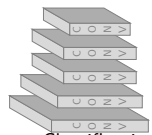
PolicyGradient

$$J(\rho) = \mathbb{E}_{p_{\rho}} [\sum_{t=0}^{\infty} r_t]$$

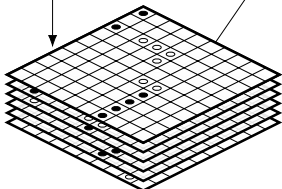
$$\rho \leftarrow \rho + \alpha \nabla_{\rho} J(\rho)$$

Supervised Value Network

$$v_{\theta}(s')$$

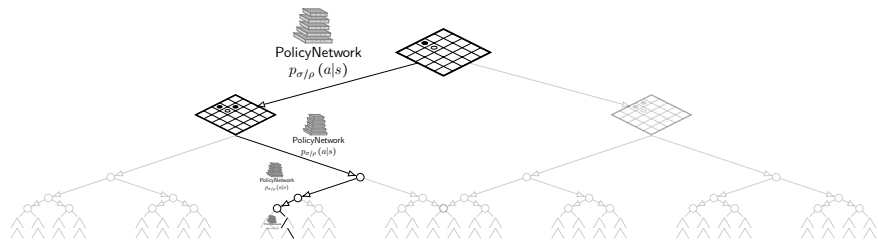


$$\min_{v_{\theta}} L(\text{predicted outcome, actual outcome})$$

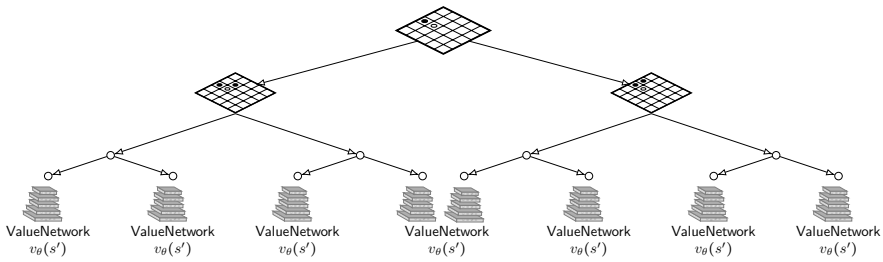


Self Play Dataset D_2

GUIDED SEARCH SEARCH USING THE POLICY NETWORK

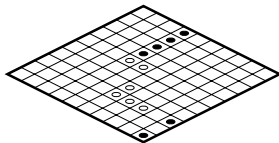


DEPTH-LIMITED SEARCH USING THE VALUE NETWORK



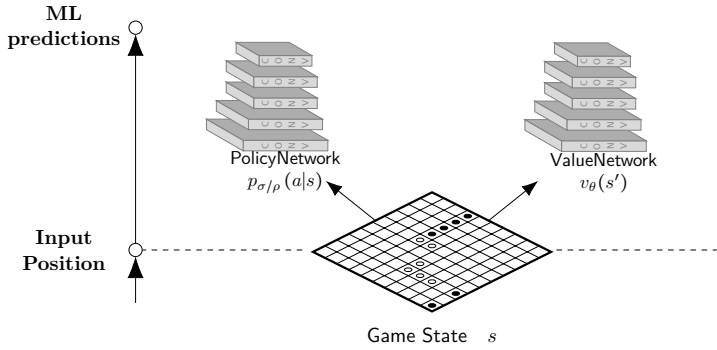
ALPHAGO PREDICTION PIPELINE

Input
Position

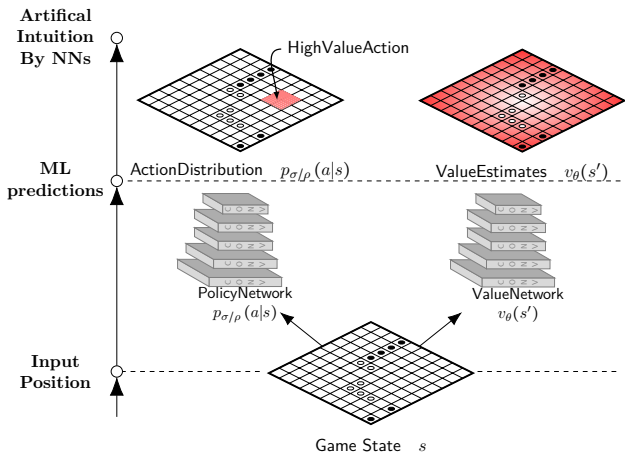


Game State s

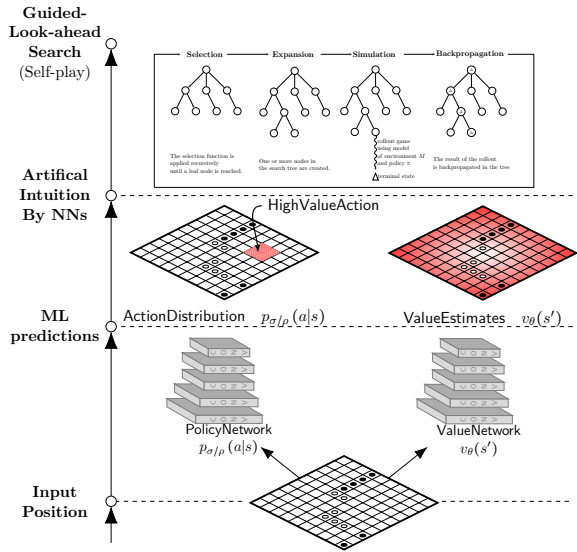
ALPHA Go PREDICTION PIPELINE



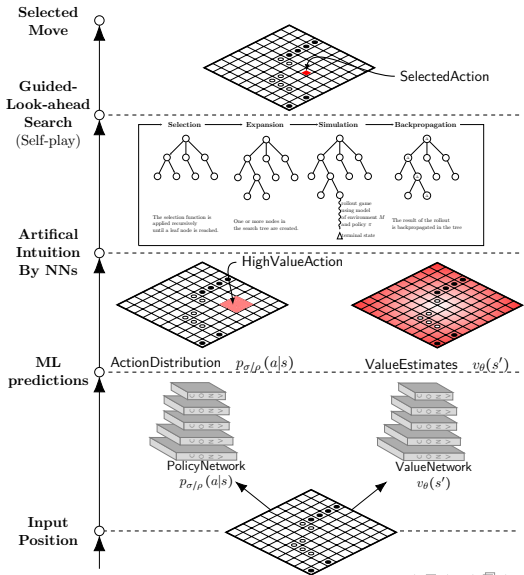
ALPHAGO PREDICTION PIPELINE



ALPHAGO PREDICTION PIPELINE



ALPHA GO PREDICTION PIPELINE



ALPHAGo VS LEE SEDOL



- ▶ In March 2016, Alpha Go won against Lee Sedol 4-1
- ▶ Lee Sedol was 18-time World Champion prior to the game
- ▶ Two famous moves: Move 37 by AlphaGo and Move 78 by Sedol

ALPHAGo: KEY TAKEAWAYS


1. Different AI techniques can be complementary
 - ▶ Supervised learning
 - ▶ Reinforcement learning
 - ▶ Search
 - ▶ Rules/Domain Knowledge
 - ▶ **What ever it takes to win!!** 😊
2. Self-play
3. Vast computation still required for training and inference
 - ▶ AlphaGo used 1200 CPUs and 176 GPUs

AlphaGo Zero '2017¹²

nature

Article | [Published: 19 October 2017](#)

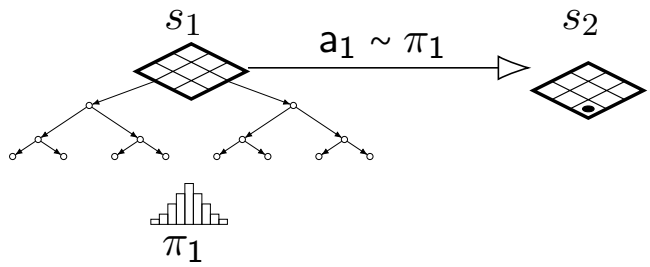
Mastering the game of Go without human knowledge

[David Silver](#) , [Julian Schrittwieser](#), [Karen Simonyan](#), [Ioannis Antonoglou](#), [Aja Huang](#), [Arthur Guez](#), [Thomas Hubert](#), [Lucas Baker](#), [Matthew Lai](#), [Adrian Bolton](#), [Yutian Chen](#), [Timothy Lillicrap](#), [Fan Hui](#), [Laurent Sifre](#), [George van den Driessche](#), [Thore Graepel](#) & [Demis Hassabis](#)

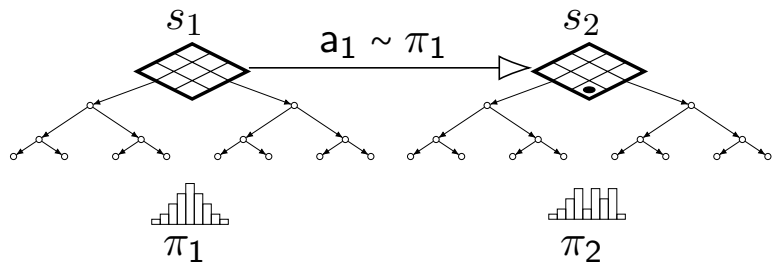
ALPHA GO ZERO

- ▶ AlphaGo Zero is a successor to AlphaGo
- ▶ AlphaGo Zero is **simpler** and **stronger** than AlphaGo
 - ▶ AlphaGo Zero beats AlphaGo 100 – 0 in matches
- ▶ AlphaGo Zero starts from **Zero** domain knowledge
 - ▶ Uses a **single neural network** (compared to 4 NNs in AlphaGo)
 - ▶ Learns by **Self-Play only** (No supervised learning like in AlphaGo)

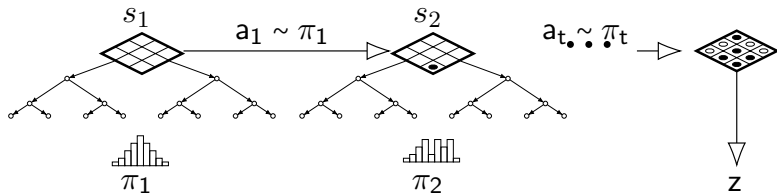
ALPHAGO ZERO SELF-PLAY TRAINING ALGORITHM



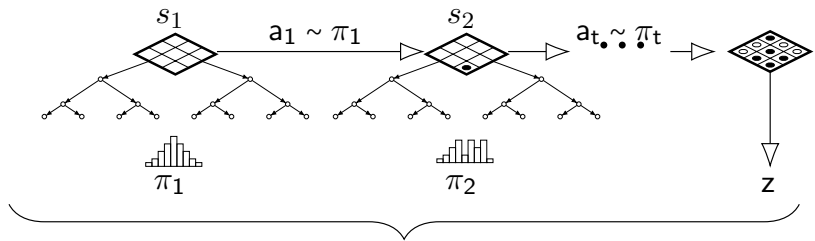
ALPHAGO ZERO SELF-PLAY TRAINING ALGORITHM



ALPHAGO ZERO SELF-PLAY TRAINING ALGORITHM



ALPHAGo ZERO SELF-PLAY TRAINING ALGORITHM

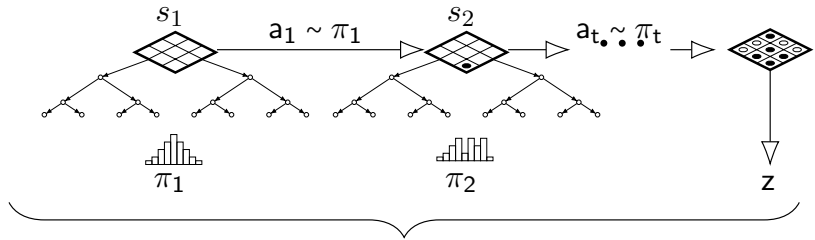


Self Play Training Data

$$f_{\theta}(s_t) = (p_t, v_t)$$

$$\theta' = \theta - \alpha \nabla_{\theta} L((p_t, v_t), (\pi_t, z))$$

ALPHA Go ZERO SELF-PLAY TRAINING ALGORITHM

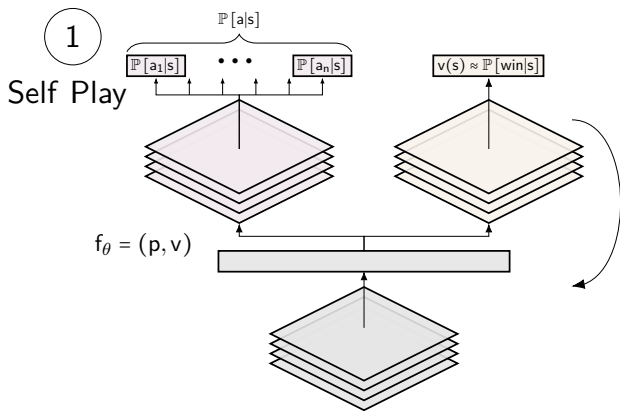


Self Play Training Data

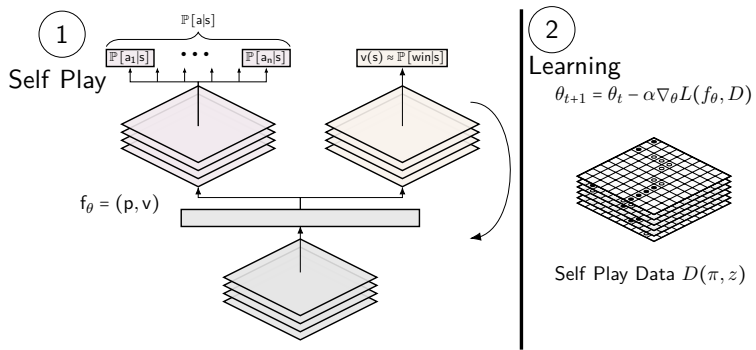
$$f_{\theta}(s_t) = (p_t, v_t) \qquad \theta' = \theta - \alpha \nabla_{\theta} L((p_t, v_t), (\pi_t, z))$$

$$L(f_{\theta}(s_t), (\pi_t, z)) = \underbrace{(z - v_t)^2}_{\text{MSE}} - \underbrace{\pi_t^T \log p_t + c \|\theta\|^2}_{\text{Cross-entropy loss}}$$

ALPHA GO ZERO SELF-PLAY TRAINING PIPELINE



ALPHAGO ZERO SELF-PLAY TRAINING PIPELINE



ALPHAZERO

- ▶ AlphaGo Zero is able to reach superhuman level at Go without any domain knowledge...
- ▶ As AlphaGo Zero is not dependent on Go, **can the same algorithm play other games?**
- ▶ AlphaZero extends AlphaGo to play **not only Go but also Chess and Shogi**
 - ▶ The same algorithm achieves superhuman performance on all three games

ALPHAZERO: KEY TAKEAWAYS

1. Being able to play three games at a superhuman level, AlphaZero is one step closer to **general AI**
2. Massive compute power still required for training AlphaZero
 - ▶ 5000 TPU_s
3. Self-play

PRESENTATION SUMMARY

- ▶ Sometimes a simpler system can be more powerful than a complex one
- ▶ Universal research principle: **strive for generality**, simplicity, Occam's Razor
- ▶ **Self-play: no human bias, learn from first principles**
- ▶ **Deep RL is still in its infancy**, a lot to more to be expected in the next few years
- ▶ Open challenges: Sample efficiency, data efficiency
 - ▶ Yes, AlphaGo can learn to play Go after hundreds of game years, but a human can reach a decent level of play in only a couple of hours
 - ▶ How can we make reinforcement learning more efficient? Model-based learning is a research area with increasing attention

REFERENCES¹⁹

- ▶ DQN¹⁴
- ▶ AlphaGo¹⁵
- ▶ AlphaGo Zero¹⁶
- ▶ AlphaZero¹⁷
- ▶ AlphaStar¹⁸

2em1¹⁴ Volodymyr Mnih et al. “Human-level control through deep reinforcement learning”. In: *Nature* 518.7540 (Feb. 2015), pp. 529–533. ISSN: 00280836. URL: <http://dx.doi.org/10.1038/nature14236>.

2em1¹⁵ David Silver et al. “Mastering the game of Go without human knowledge”. In: *Nature* 550 (Oct. 2017), pp. 354–. URL: <http://dx.doi.org/10.1038/nature24270>.

2em1¹⁶ David Silver et al. “Mastering the game of Go without human knowledge”. In: *Nature* 550 (Oct. 2017), pp. 354–. URL: <http://dx.doi.org/10.1038/nature24270>.

2em1¹⁷ David Silver et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play”. In: *Science* 362.6419 (2018), pp. 1140–1144. URL: <http://science.sciencemag.org/content/362/6419/1140/tab-pdf>.

2em1¹⁸ Oriol Vinyals et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. In: *Nature* 575 (Nov. 2019). DOI: 10.1038/s41586-019-1724-z.

2em1¹⁹ Thanks to Rolf Stadler for Reviewing and discussing drafts of this presentation