

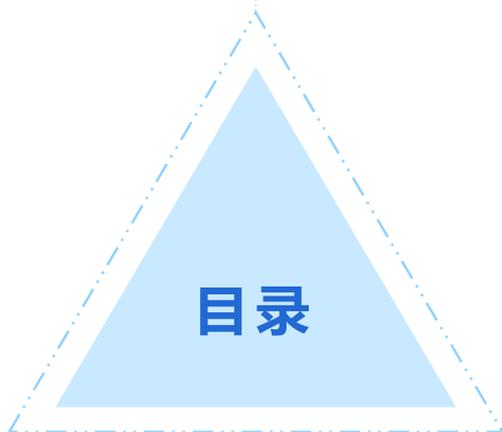
AlphaGo的胜利

(理工类)

东南大学 计算机科学与工程学院



目录



- 1 / 围棋简介
- 2 / 传统方法
- 3 / AlphaGo
- 4 / 后AlphaGo时代
- 5 / 总结与展望

围棋简介



“琴棋书画”之棋

- 围棋起源于中国，称之为
“弈”
- 经日本传入欧洲被译成
“Go”



尧造围棋，丹朱善之

——先秦典籍《世本》

围棋基本介绍



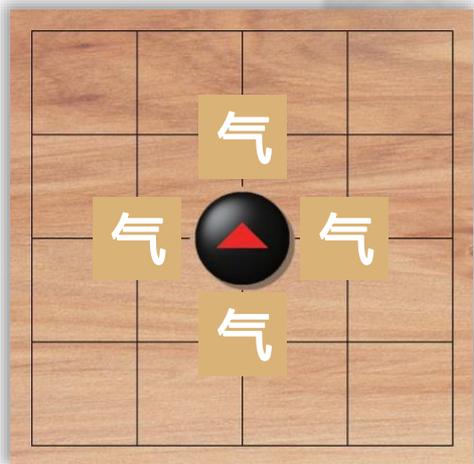
棋盘与棋子

行棋规则

胜负判定

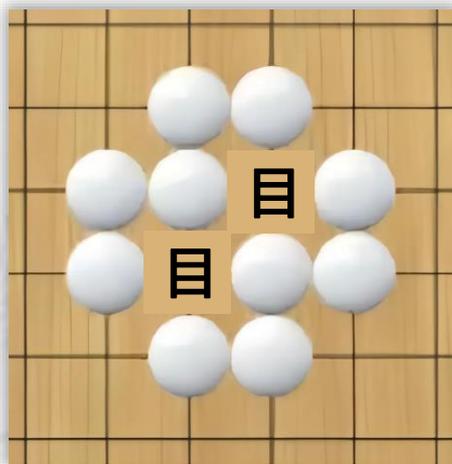
特殊规则

围棋基本介绍



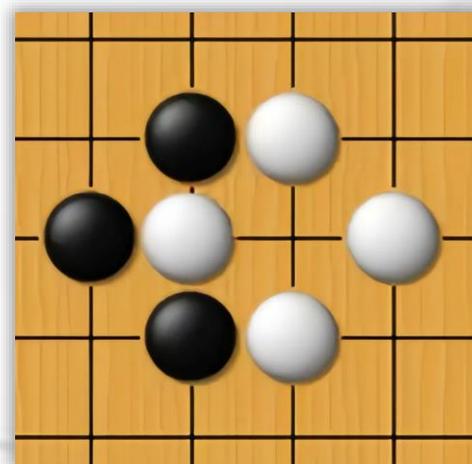
气

与棋子直线相邻空交叉点



目

被一方棋子所围地域的
空交叉点



打劫

双方可以轮流提取对方棋
子的情况

围棋棋手的养成



一名优秀的棋手的成长路径

基本的布局理论

死活、对杀常识

基本的布局理论

数万死活题

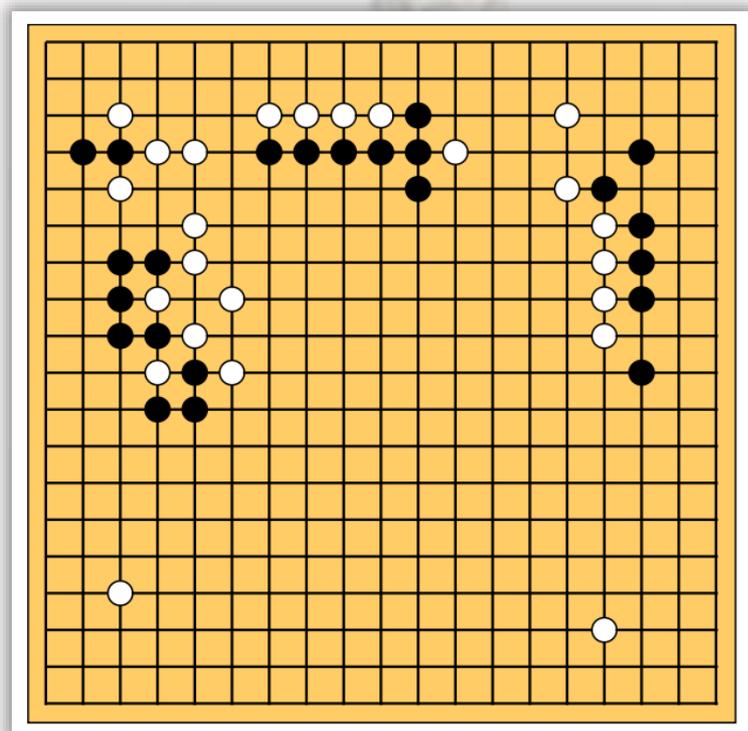
对局中磨练

短期/局部战术

长期/全局战略



围棋复杂性



➤ 棋盘规模

19x19 的棋盘拥有 361 个交叉点

➤ 棋局数量

大约为 10^{170} ，超越宇宙原子数

➤ 决策深度

通常每盘棋需要250步左右

围棋的特点



典型的完全信息博弈问题

信息对称

零和

纳什平衡

任何一种完全信息博弈的棋类游戏，在当前状态下都存在一个最优策略以及与之对应的最优状态价值函数



恩斯特·策梅洛

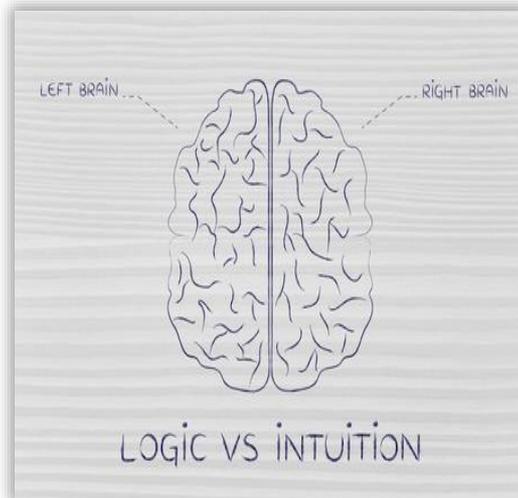
围棋与人工智能



围棋长期以来被认为是人工智能领域的“圣杯”



庞大的搜索空间



长远战略与直觉

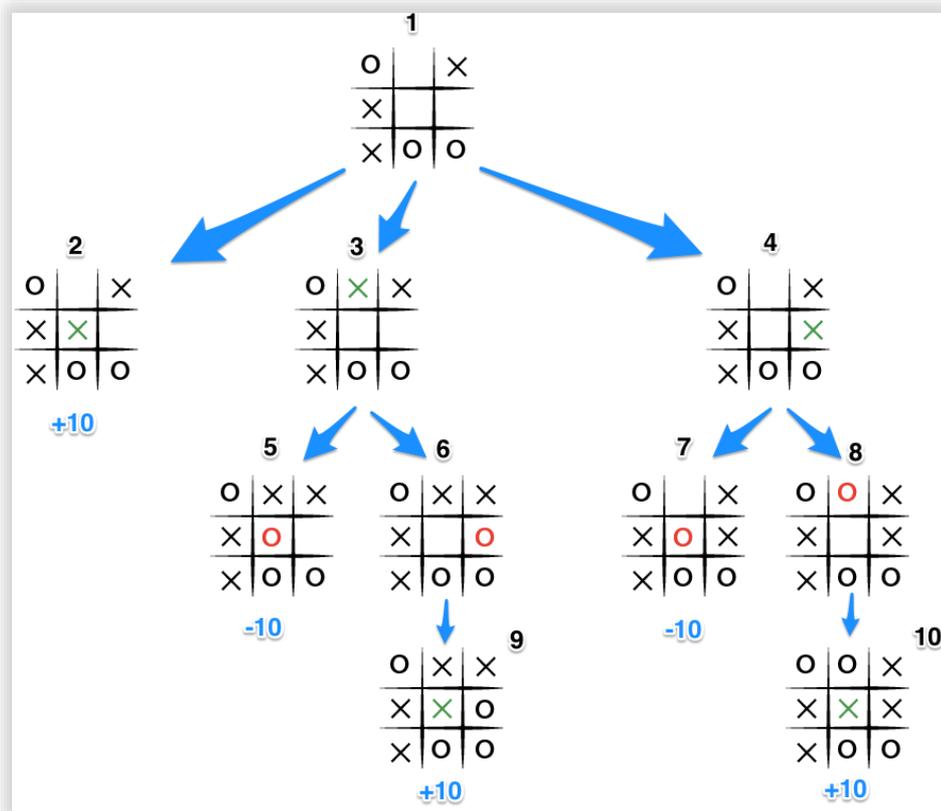
$$\begin{aligned}
 & \mathcal{L}_{\text{StandardModel}} = \\
 & -\frac{1}{2} \partial_\mu g_\nu^\alpha \partial_\mu g_\nu^\alpha - g_\mu f^{abc} \partial_\mu g_\nu^\alpha g_\nu^\beta g_\nu^\gamma - \frac{1}{2} g^2 f^{abc} f^{ade} g_\nu^\alpha g_\nu^\beta g_\nu^\gamma g_\nu^\delta + \\
 & \frac{1}{2} i g^2 (\bar{\psi}^\alpha \gamma^\mu \psi^\alpha) g_\mu^\alpha + C^a \partial^a C^a + g_s f^{abc} \partial_\mu C^a G^b G^c - \partial_\mu W_\nu^\alpha \partial_\mu W_\nu^\alpha - \\
 & M^2 W_\mu^+ W_\mu^- - \frac{1}{2} \partial_\mu Z_\nu^\alpha \partial_\mu Z_\nu^\alpha - \frac{1}{2} M^2 Z_\nu^\alpha Z_\nu^\alpha - \frac{1}{2} \partial_\mu A_\nu \partial_\mu A_\nu - \frac{1}{2} \partial_\mu H \partial_\mu H - \\
 & \frac{1}{2} m_H^2 H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - M^2 \phi^+ \phi^- - \frac{1}{2} \partial_\mu \phi^0 \partial_\mu \phi^0 - \frac{1}{2} M^2 \phi^0 \phi^0 - \beta_h \left[\frac{23}{9} g^2 + \right. \\
 & \left. \frac{20}{9} H + \frac{1}{2} (H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right] + \frac{23}{9} g^2 \alpha_h - i g_{\text{CW}} [\partial_\mu Z_\nu^\alpha (W_\mu^+ W_\nu^- - \\
 & W_\mu^- W_\nu^+) - Z_\nu^\alpha (W_\mu^+ \partial_\mu W_\nu^- - W_\mu^- \partial_\mu W_\nu^+) + Z_\nu^\alpha (W_\mu^- \partial_\mu W_\nu^+ - \\
 & W_\mu^+ \partial_\mu W_\nu^-)] - i g_{\text{SW}} [\partial_\mu A_\nu (W_\mu^+ W_\nu^- - W_\mu^- W_\nu^+) - A_\nu (W_\mu^+ \partial_\mu W_\nu^- - \\
 & W_\mu^- \partial_\mu W_\nu^+) + A_\nu (W_\mu^- \partial_\mu W_\nu^+ - W_\mu^+ \partial_\mu W_\nu^-)] - \frac{1}{2} g^2 W_\mu^+ W_\mu^- W_\nu^+ W_\nu^- + \\
 & \frac{1}{2} g^2 W_\mu^+ W_\nu^- W_\nu^+ W_\mu^- + g^2 c_w^2 (Z_\nu^\alpha W_\mu^+ Z_\nu^\alpha W_\mu^- - Z_\nu^\alpha W_\mu^+ W_\nu^+ W_\nu^-) + \\
 & g^2 s_w^2 (A_\mu W_\mu^+ A_\nu W_\nu^- - A_\mu A_\nu W_\mu^+ W_\nu^-) + g^2 s_w c_w [A_\mu Z_\nu^\alpha (W_\mu^+ W_\nu^- - \\
 & W_\mu^- W_\nu^+) - 2A_\mu Z_\nu^\alpha W_\mu^+ W_\nu^-] - g_0 [H^3 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-] - \\
 & \frac{1}{8} g^2 \alpha_h [H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2] - \\
 & g M W_\mu^+ W_\nu^- H - \frac{1}{2} g \frac{M}{c_w^2} Z_\nu^\alpha Z_\nu^\alpha H - \frac{1}{2} i g [W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^- \partial_\mu \phi^0) - \\
 & W_\mu^- (\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)] + \frac{1}{2} g [W_\mu^+ (H \partial_\mu \phi^- - \phi^- \partial_\mu H) - W_\mu^- (H \partial_\mu \phi^+ - \\
 & \phi^+ \partial_\mu H)] + \frac{1}{2} g \frac{1}{c_w} (Z_\nu^\alpha (H \partial_\mu \phi^0 - \phi^0 \partial_\mu H) - i g \frac{2c_w}{2c_w} Z_\nu^\alpha (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) + \\
 & i g s_w M A_\mu (W_\mu^+ \phi^- - W_\mu^- \phi^+) - i g \frac{1-2c_w^2}{2c_w} Z_\nu^\alpha (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) + \\
 & i g s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \frac{1}{4} g^2 W_\mu^+ W_\mu^- [H^2 + (\phi^0)^2 + 2\phi^+ \phi^-] - \\
 & \frac{1}{4} g^2 \frac{1}{c_w} Z_\nu^\alpha Z_\nu^\alpha [H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)\phi^+ \phi^-] - \frac{1}{2} g^2 \frac{2c_w}{2c_w} Z_\nu^\alpha \phi^0 (W_\mu^+ \phi^- + \\
 & W_\mu^- \phi^+) - \frac{1}{2} i g^2 \frac{2c_w}{2c_w} Z_\nu^\alpha H (W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2} g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\
 & W_\mu^- \phi^+) + \frac{1}{2} i g^2 s_w A_\mu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{2c_w}{2c_w} (2c_w^2 - 1) Z_\nu^\alpha A_\mu \phi^+ \phi^- - \\
 & g^2 s_w^2 A_\mu A_\nu \phi^+ \phi^- - e^2 (\gamma \partial + m_\lambda^2) e^\lambda - \bar{\nu}^\lambda \gamma \partial \nu^\lambda - \bar{u}_j^2 (\gamma \partial + m_j^2) u_j^2 - \\
 & d_j^2 (\gamma \partial + m_j^2) d_j^2 + i g s_w A_\mu [-(\bar{e}^\lambda \gamma^\mu e^\lambda) + \frac{2}{3} (\bar{u}_j^2 \gamma^\mu u_j^2) - \frac{1}{3} (\bar{d}_j^2 \gamma^\mu d_j^2)] + \\
 & \frac{i g}{2c_w} Z_\nu^\alpha [(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (e^\lambda \gamma^\mu (1 s_w^2 - 1 - \gamma^5) e^\lambda) + (\bar{u}_j^2 \gamma^\mu (\frac{2}{3} s_w^2 - \\
 & 1 - \gamma^5) u_j^2) + (\bar{d}_j^2 \gamma^\mu (1 - \frac{2}{3} s_w^2 - \gamma^5) d_j^2)] + \frac{i g}{2c_w} W_\mu^+ [(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) e^\lambda) + \\
 & (\bar{e}^\lambda \gamma^\mu (1 + \gamma^5) C_{\lambda\alpha} d_j^2)] + \frac{i g}{2c_w} W_\mu^- [(\bar{e}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{d}_j^2 C_{\lambda\alpha} \gamma^\mu (1 + \\
 & \gamma^5) u_j^2)] + \frac{i g}{2c_w} \frac{2c_w}{2c_w} [-\phi^+ (\bar{\nu}^\lambda (1 - \gamma^5) e^\lambda) + \phi^- (\bar{e}^\lambda (1 + \gamma^5) \nu^\lambda)] -
 \end{aligned}$$

评估函数的复杂性

传统方法



博弈树搜索



模拟所有可能未来局面



评估每个未来局面好坏

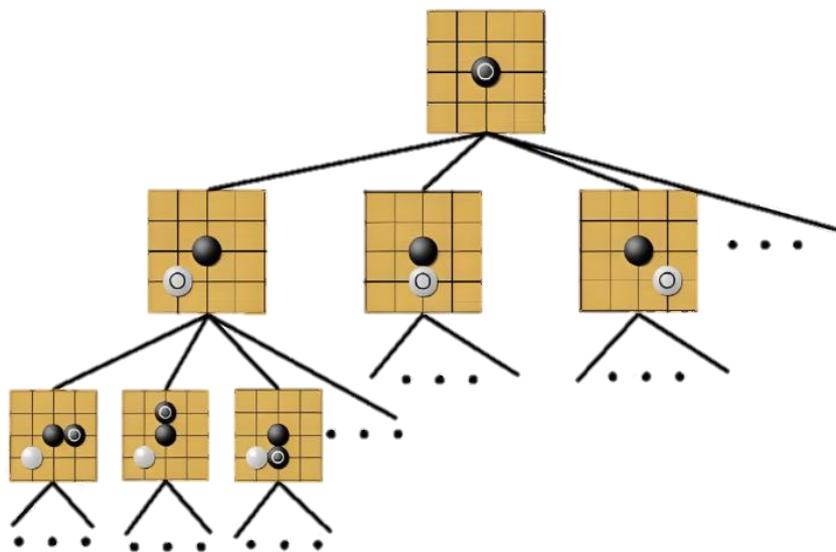


帮助玩家找到**最优策略**

Minimax (或MinMax) 算法



假设双方都是**理性玩家**，并且都试图做出**最优决策**



我方棋步，选择对我最有利 (max)



对方棋步，选择对我最有害 (min)

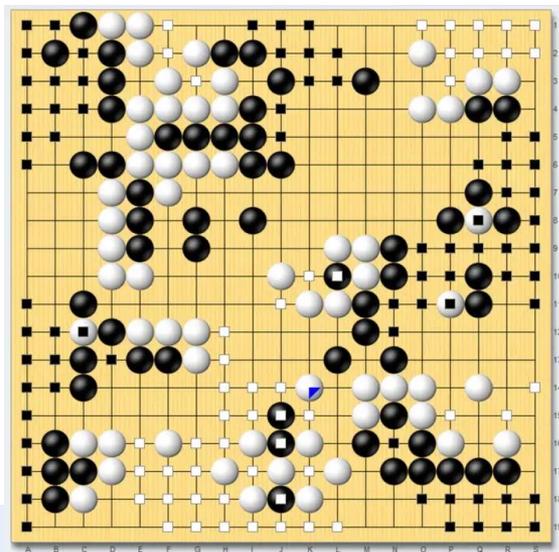


我方棋步，选择对我最有利 (max)



...

评估局面



地盘控制

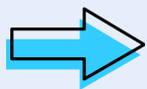
气数

目数

边角控制

其他评估

人工设计
评估函数



局限性

局部与全局难平衡

死活等问题的复杂性

状态空间庞大

博弈树搜索的局限

1. 计算资源的限制

不能构建整个决策树，较浅的搜索深度可能导致不能发现最佳策略

2. 规则和局面特征多样

难以构建精确的评估函数，整个搜索过程的价值将大打折扣



第一代围棋AI“手谈”
达到约13层的搜索深度

蒙特卡洛树搜索 (Monte Carlo Tree Search)

蒙特卡洛树搜索算法开启第二代围棋人工智能

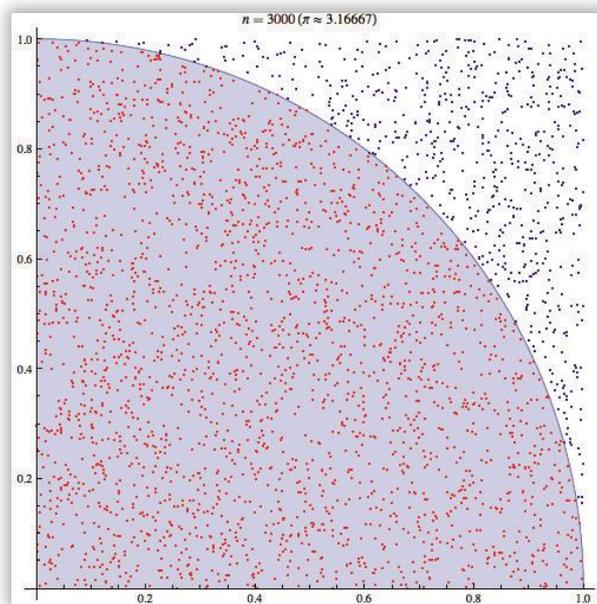


2015年，职业棋手连笑七段在让4子和让5子两场人机对战中，战胜计算机围棋冠军韩国软件“石子旋风”（DolBaram），但在让6子比赛中败下阵来。

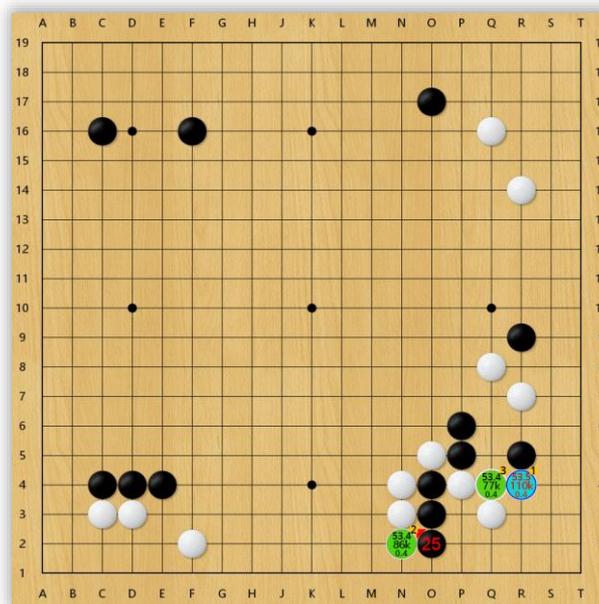
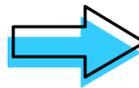
蒙特卡洛方法



一种基于“随机数”的计算方法。



空间内多次落点来近似计算圆周率



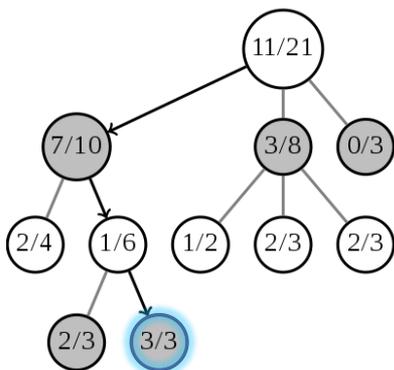
多次模拟棋局后看哪个选点的“获胜概率”最高

蒙特卡洛树搜索 (MCTS)



蒙特卡洛方法和搜索树结构的结合

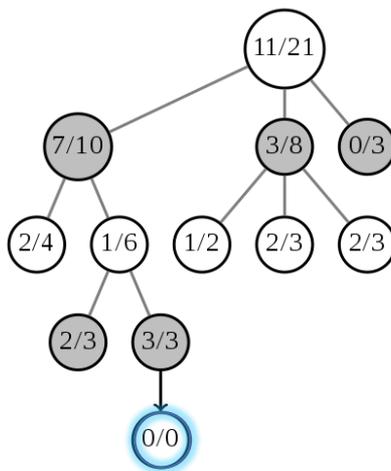
选择 (Selection)



选择一个“最值得看的子节点”

1. 已被证明胜率高
2. 尚未被充分探索

扩展 (Expansion)



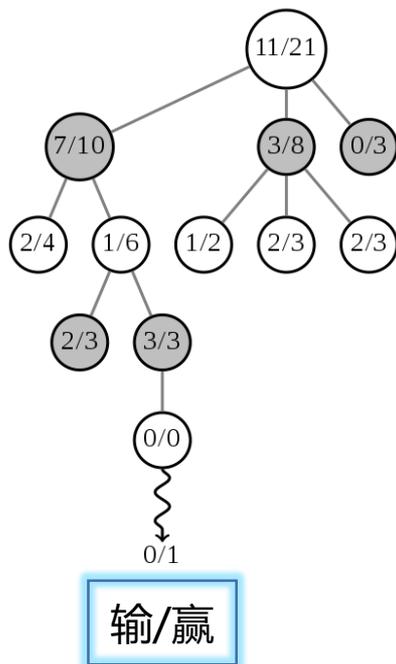
在被选择节点添加一个子节点来扩展搜索树

蒙特卡洛树搜索 (MCTS)



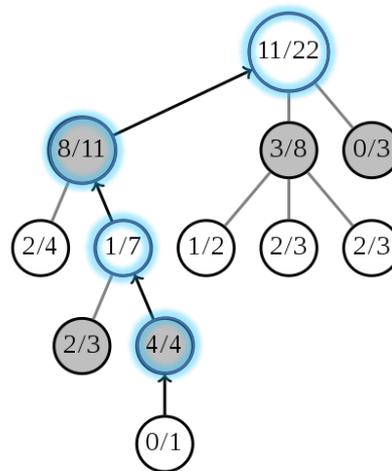
蒙特卡洛方法和搜索树结构的结合

模拟 (Simulation)



随机落子模拟到游戏结束, 评估最终结果

回溯 (Backpropagation)



模拟结果沿着路径更新节点的统计信息

- ↓
1. 访问次数
 2. 节点胜率

蒙特卡洛树搜索的优势和局限

优势

适用于庞大状态空间

逐步收敛到最优解

避免复杂的评估函数

局限

模拟阶段的随机走法过于粗糙

效率取决于模拟的次数



第二代AI围棋程序



棋手战鹰
职业二段

传统方法和人类的差别



棋感、大局观，非“机器”所能掌握

蒙特卡洛树搜索法

模拟统计结果驱动

局部累积形成全局

没有积累缺乏创新

局面感知能力较弱



职业棋手

直觉与逻辑驱动

局部与全局平衡

过往与当前棋局

了解定式和棋形

计算机程序想战胜职业棋手，还有很长的路要走



棋手连笑
职业七段

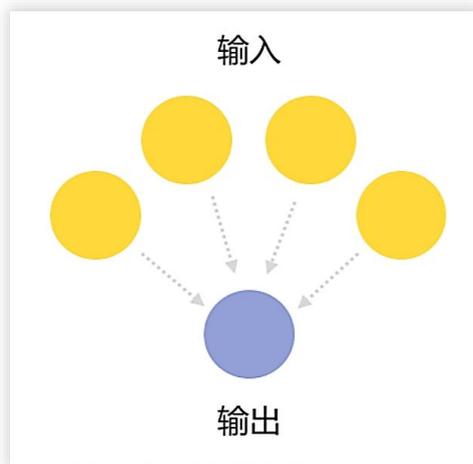
AlphaGo



机器学习

监督学习

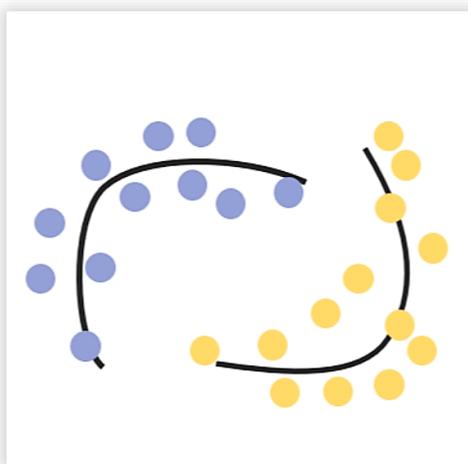
Supervised Learning



有预测目标 Y ，通过输入 X 预测

无监督学习

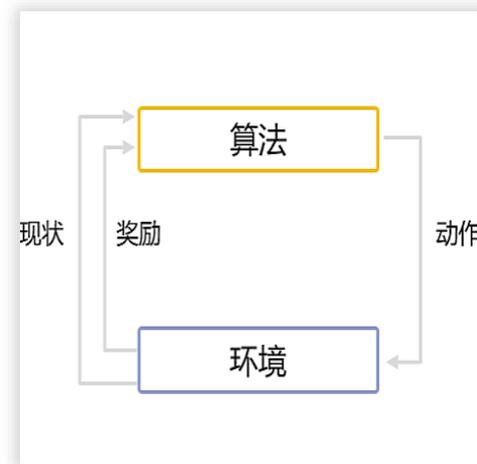
Unsupervised Learning



只通过输入 X 进行分析
并识别模式

强化学习

Reinforcement Learning



通过环境与奖励循环迭代优化出
最合适的动作

深度卷积神经网络

输入层

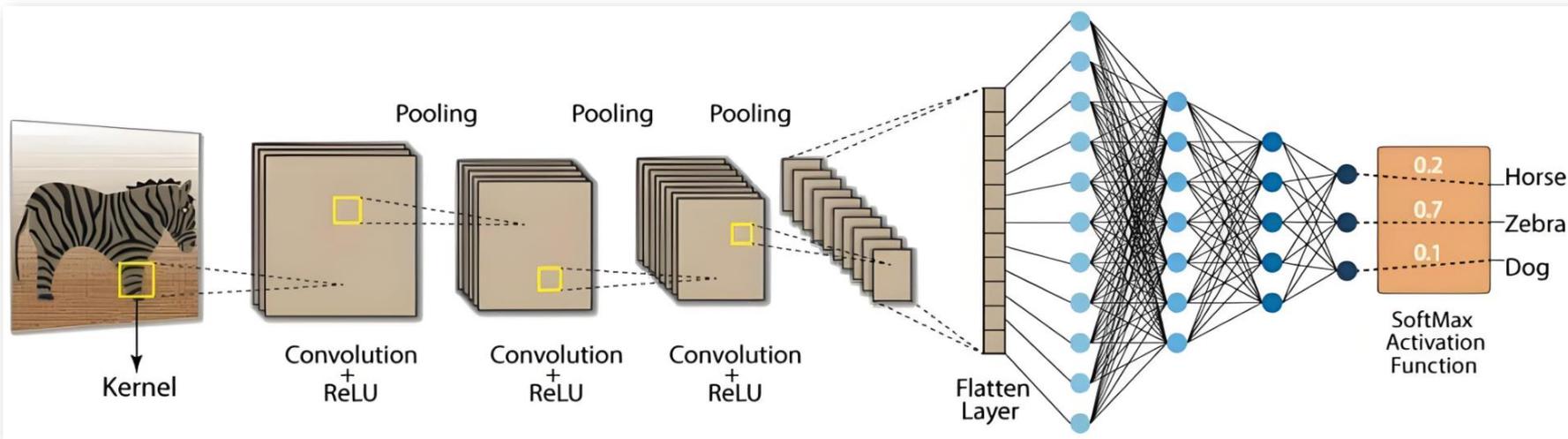
接收多维数组

隐含层

卷积层、汇合层和全连接层

输出层

输出分类标签



AlphaGo



结合蒙特卡洛树搜索 (MCTS) 和深度学习的围棋人工智能

策略网络

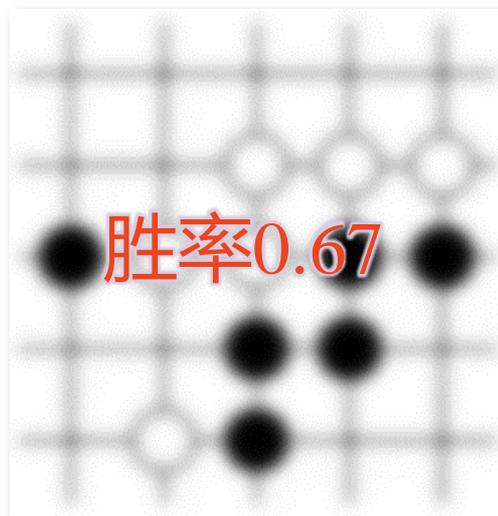
预测下一步最优的棋步

0	0	0.1	0	0
0	0.1	-	-	-
-	-	-	-	-
0	0.5	-	-	0
0	0.1	-	0.2	0

每一个可能的落子点赋予一个概率值，快速筛选出一组有潜力的候选动作

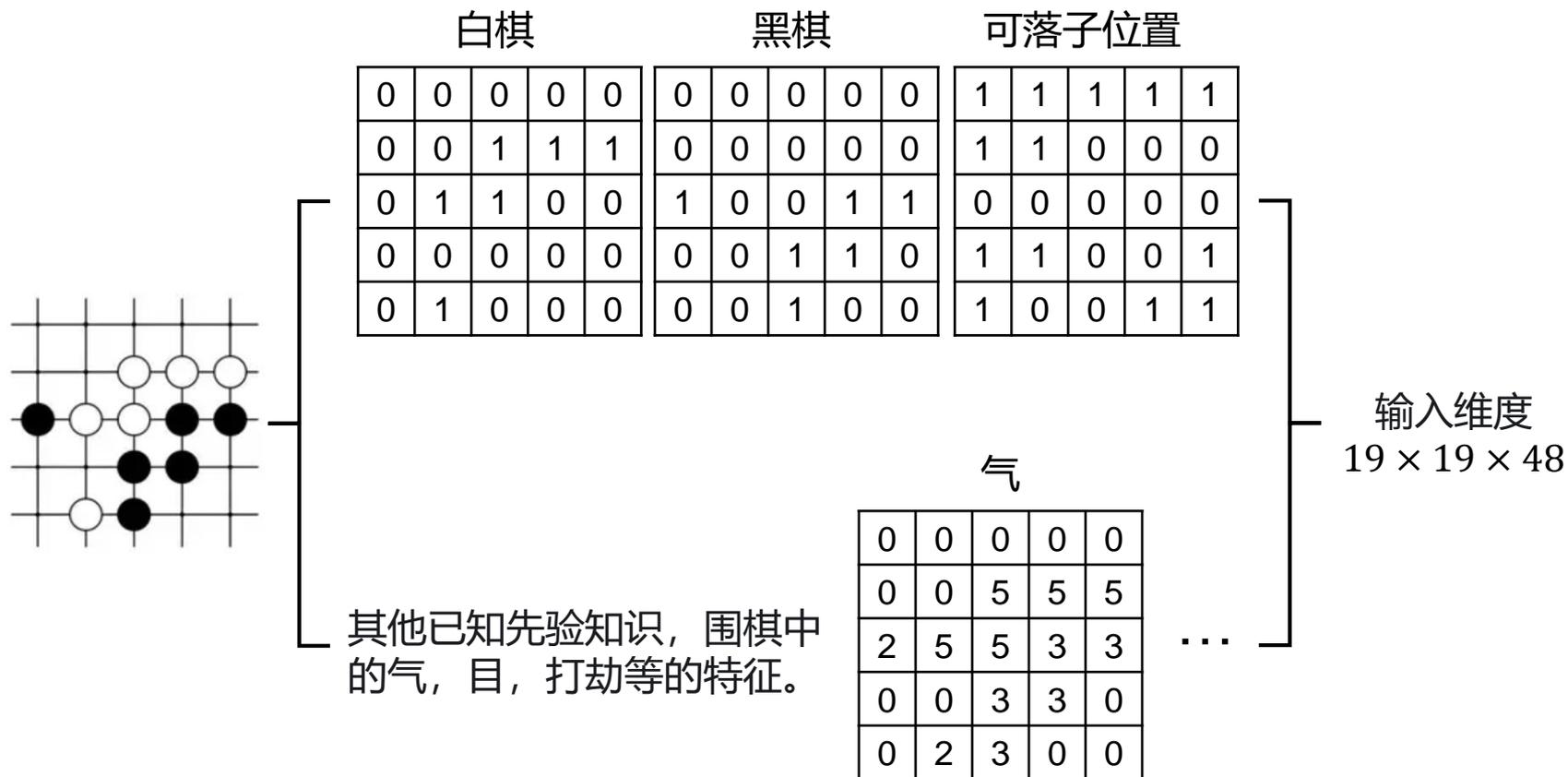
价值网络

评估当前局面的胜负概率



评估整个局面对于当前执棋方的好坏，替代部分模拟对局

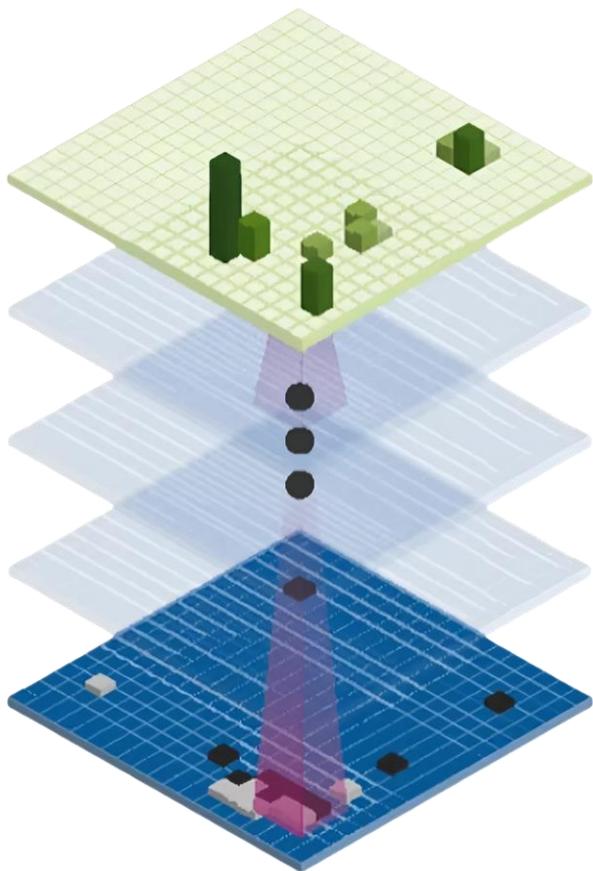
神经网络的输入



策略网络P(Policy Networks)



预测下一步棋的可能落子位置



0	0	0.1	0	0
0	0.1	-	-	-
-	-	-	-	-
0	0.5	-	-	0
0	0.1	-	0.2	0

每一个合法动作的概率分布 $P(a|s)$



Softmax输出层
非线性函数
13层卷积神经网络

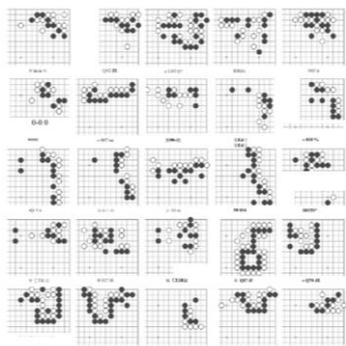


输入

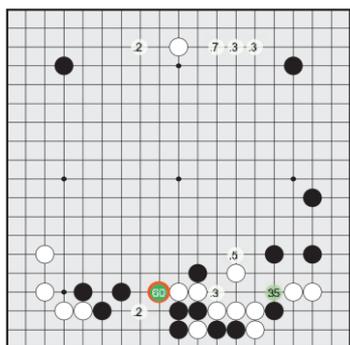
监督学习训练策略网络 P_σ



模仿人类专家下棋



160,000局6至9段人类选手的围棋对弈
30,000,000个构建的 $\langle s, a \rangle$



策略网络输出概率分布 $P_\sigma(a/s)$

0	0	0	0	0
0	0	-	-	-
-	-	-	-	-
0	1	-	-	0
0	0	-	0	0

0	0	0.1	0	0
0	0.1	-	-	-
-	-	-	-	-
0	0.5	-	-	0
0	0.1	-	0.2	0

优化目标

快速走子网络 P_π (Rollout Policy Networks)



用于快速模拟棋局发展，代替原本蒙特卡洛树中模拟阶段的随机落子

轻量化网络设计：

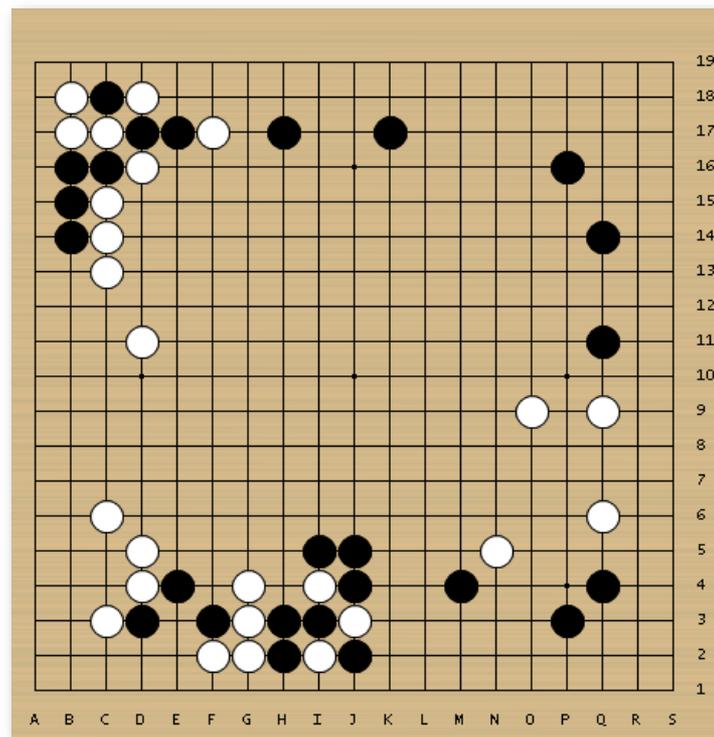
输入仅包含当前黑白棋子的分布，
网络由较少层数的神经网络构成

运算时间：

$$P_\pi \ 2\mu s > P_\sigma \ 3ms$$

落子准确率：

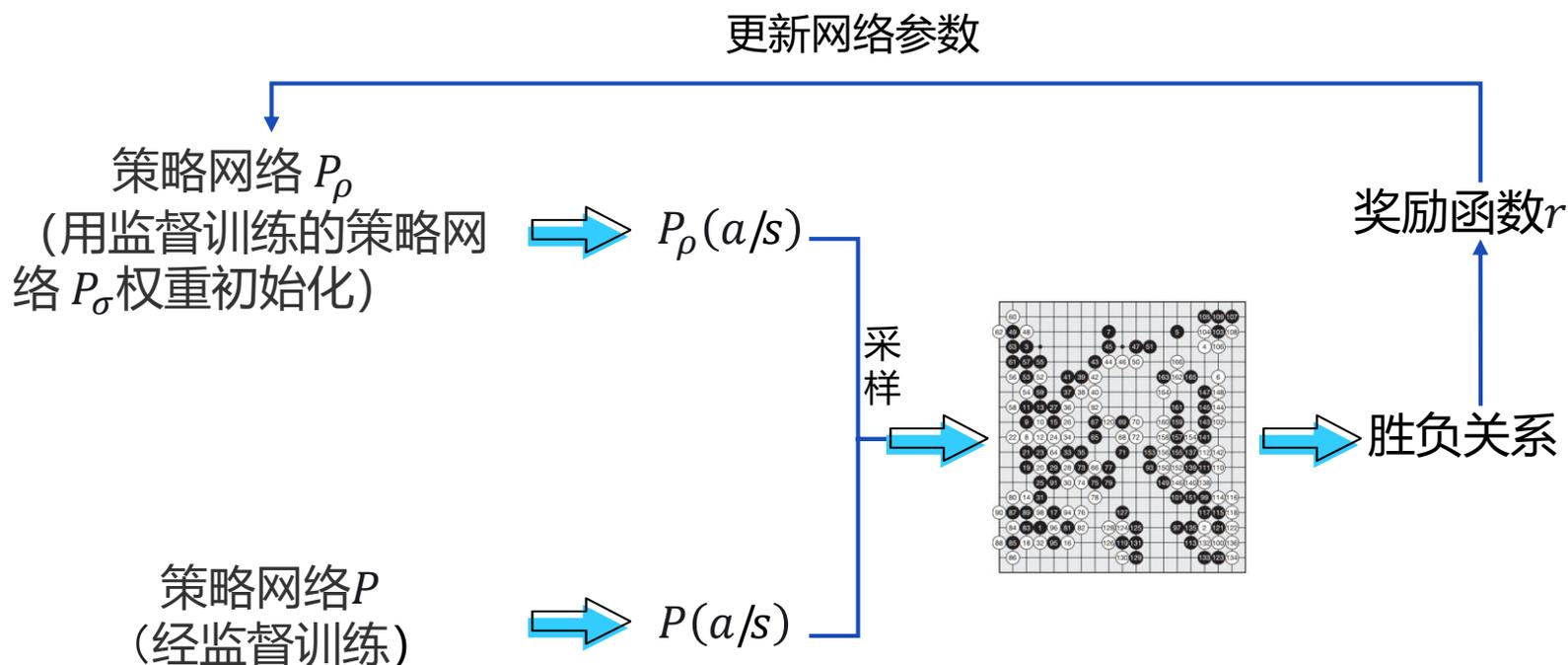
$$P_\sigma \ 55.7\% > P_\pi \ 24.2\% > \text{随机选点}$$



强化学习训练策略网络 P_ρ



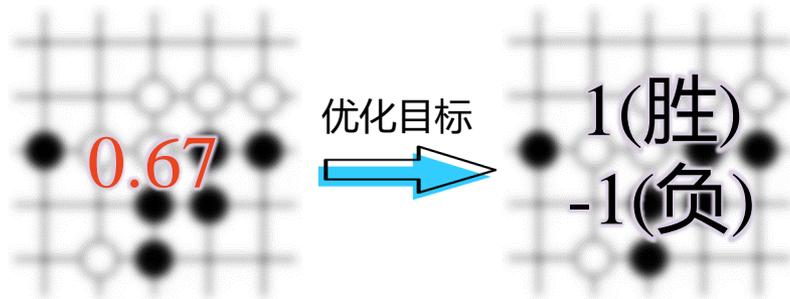
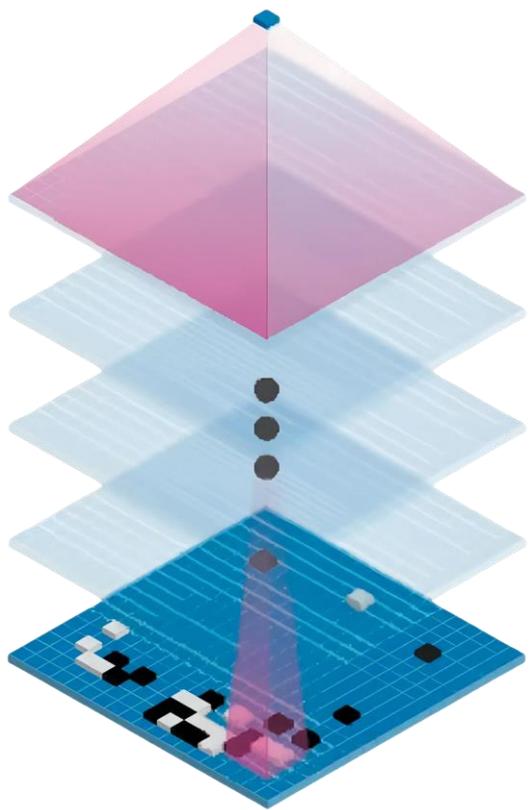
策略网络的“左右互搏”，根据自我对弈的结果来优化决策



价值网络 V_θ



神经网络学习一个函数 $v(s)$, 直接输出一个该局面赢得比赛的概率



价值网络输出 $v_\theta(s)$

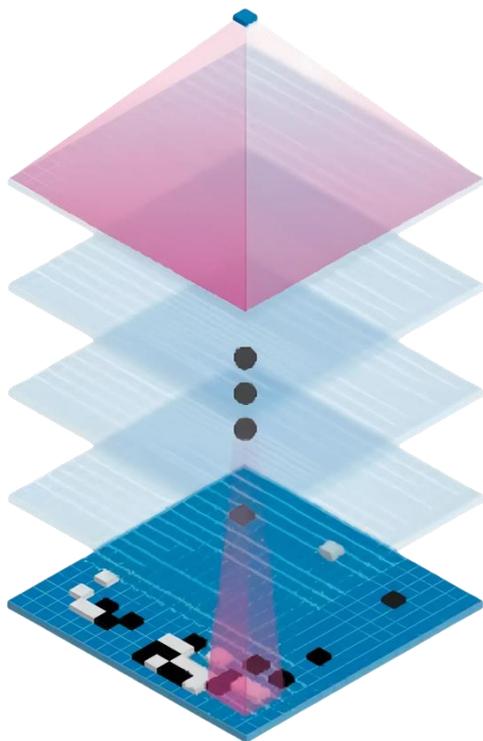
真实的胜负关系 $v^*(s)$

传统方法 — 手工设计复杂的评估函数
— 模拟直到游戏结束再评估

强化学习训练价值网络 V_θ



直接从人类完整棋局中学习价值网络会导致过拟合



棋面胜负预测 $v_\theta(s) \approx$ 本局真实胜负关系 $v^*(g)$



某个棋面下最终赢棋预测值 $v_\theta(s)$



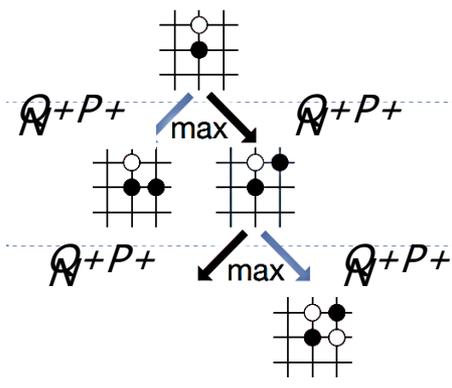
抽样出3000万个棋面 s



策略网络 P_ρ 进行自我对弈生成3000万个棋局 g

AlphaGo中的蒙特卡罗树搜索

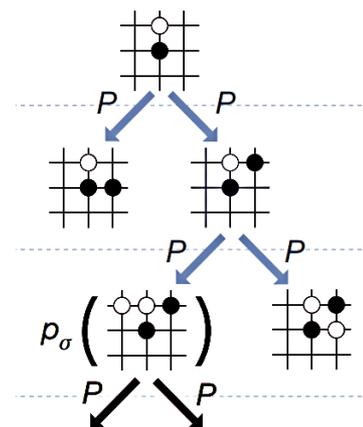
a 选择分支



节点胜率 $Q(s, a)$
 策略网络 $P(s, a)$
 访问次数 $N(s, a)$

} 选择棋面 s 下的最佳动作 a

b 扩展分支

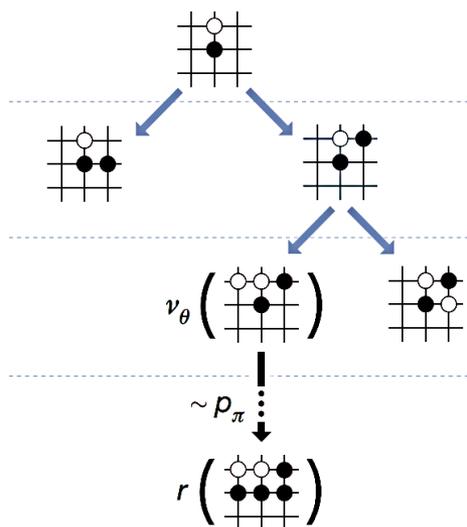


$$P(s, a) = P_{\sigma}(s, a)$$

保存子节点棋面策略网络结果作为先验概率

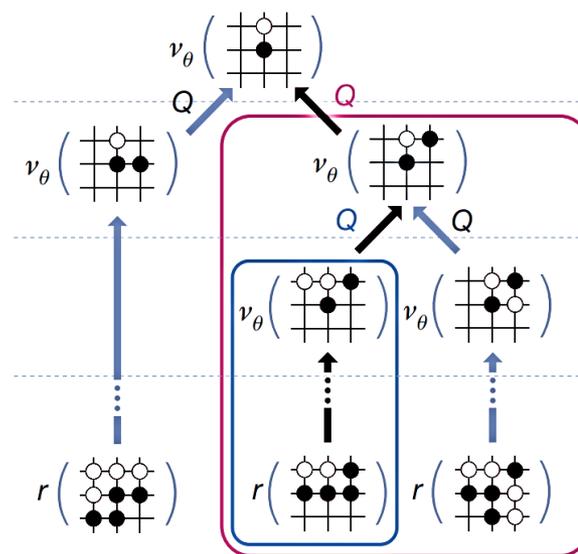
AlphaGo中的蒙特卡罗树搜索

c 结果判断



价值网络输出结果 v_θ
快速走子网络 P_π 模拟 } 胜负评估 r

d 反向传播



反向更新所经过的边的
节点胜率 Q 和访问次数 N

最终决策：「选择访问次数最多的落子节点」

AlphaGo计算需求

传统围棋程序



1 个 线程
1 个 CPU
1 个 GPU

AlphaGo最终版本



40 个 线程
48 个 CPU
8 个 GPU

AlphaGo分布式版本



40 个 线程
1202 个 CPU
176 个 GPU

终局之战

AlphaGo

分布式架构

1202个CPU与176个GPU训练了大约1周



2015年, AlphaGo Fan **5 : 0** 樊
辉二段

AlphaGo Lee

分布式架构

采用48个TPU训练数月



2016, AlphaGo Lee **4 : 1** 世界
冠军李世石九段

终局之战

2017 年 AlphaGo Master

单机版
4 TPU

AlphaGo自我对弈提高棋力
拥有更强大的策略/价值网络
运算量只有AlphaGo Lee的十分之一

3 : 0

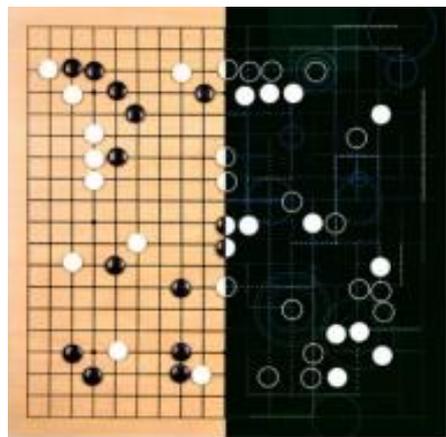


柯洁九段
世界排名第一
“人类最强棋手”

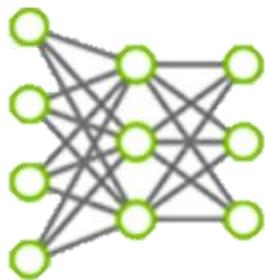
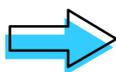
就算阿尔法狗战胜了李世石，但它赢不了我！
(面对AlphaGo Lee时)

AlphaGo Zero

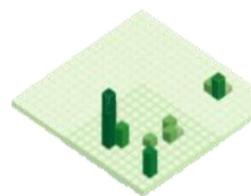
 从“人类专家知识”转向“纯粹的机器学习”



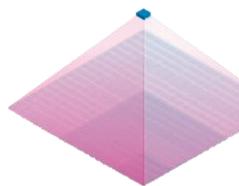
无领域知识输入 s



单一神经网络 f_{θ}



策略 p

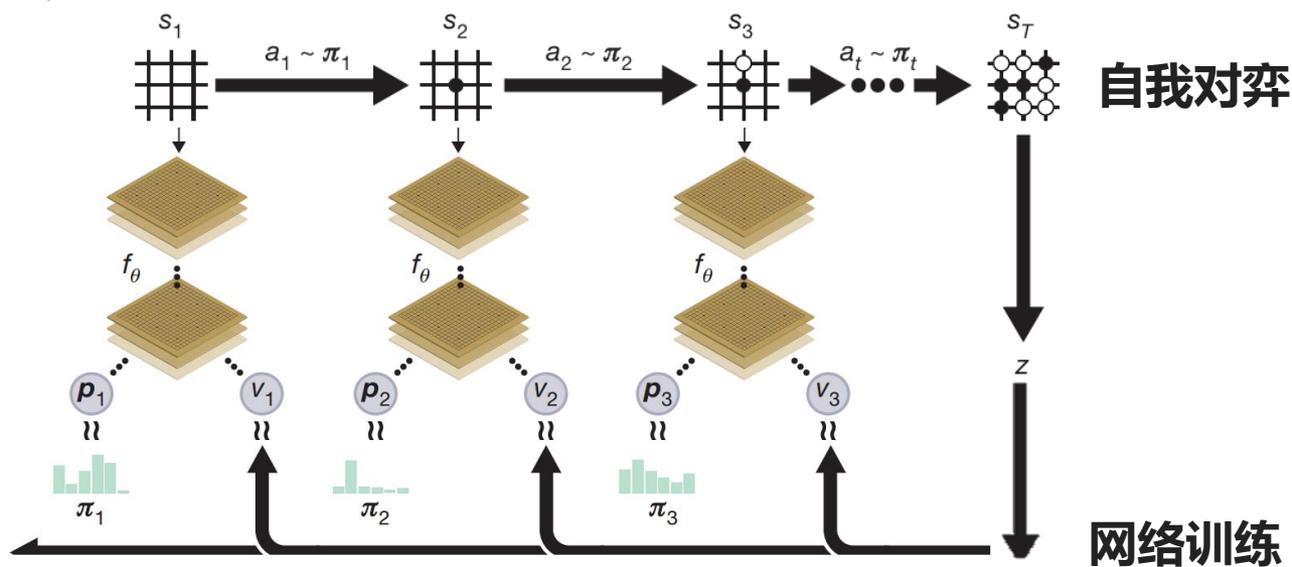


价值 v

AlphaGo Zero——自我对弈



从“人类专家知识”转向“纯粹的机器学习”



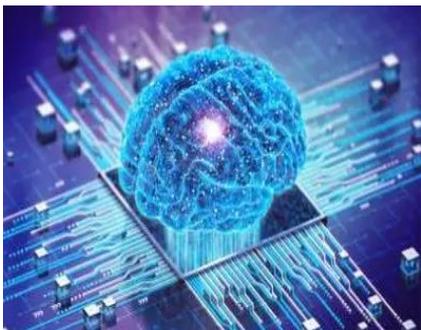
- 优化目标:**
1. 预测落子 \approx 实际落子
 2. 预测获胜概率 \approx 实际对局结果

AlphaGo Zero的里程碑意义

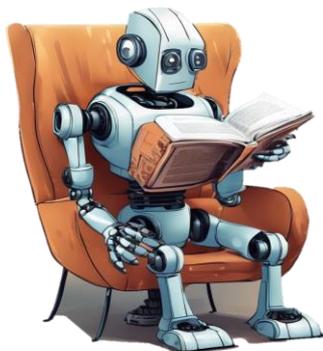
AlphaGo Zero vs. AlphaGo Lee 100 : 0

AlphaGo Zero vs. AlphaGo Master 89 : 11

“一个纯净、纯粹自我学习的AlphaGo是最强的...对于AlphaGo的自我进步来讲...人类太多余了。😏”



证明AI的潜力



AI可以无需人类知识



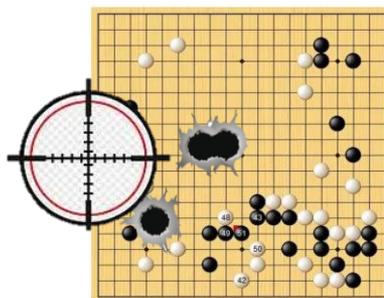
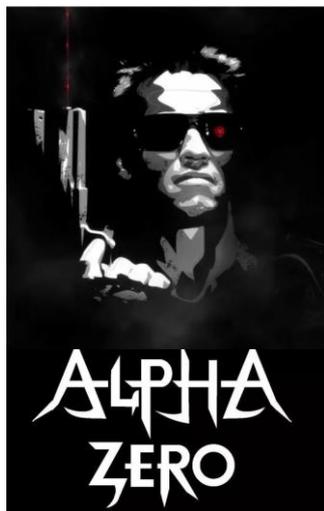
后AlphaGo时代



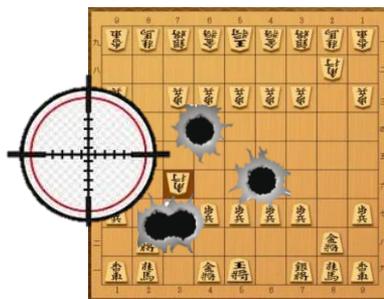
AlphaZero



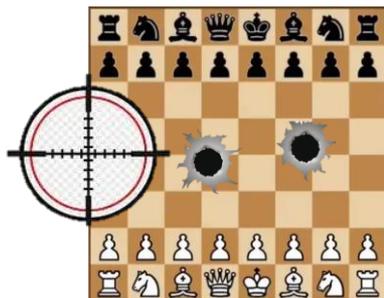
棋类游戏的“终结者”，迈向棋类通用人工智能系统



围棋



将棋



国际象棋

AlphaFold系列



Alpha系列首次应用于科学领域的重大突破

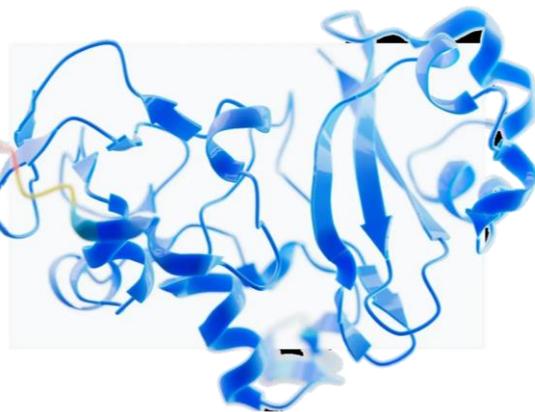
多个氨基酸单元组成的
氨基酸序列

结构预测模块



空间结构 { 氨基酸之间的**夹角分布**
氨基酸之间的**距离分布**

结构生成模块



蛋白质三维结构预测

AlphaStar



应用于即时战略游戏《星际争霸II》



- 实时决策与不完全信息
- 深度强化学习与模仿学习结合
- 多智能体学习

AI新时代的百花齐放



ChatGPT



Stable Diffusion



Sora



具身智能