

Deep Learning

A course about theory & practice



UNIVERSITÀ
DI PAVIA

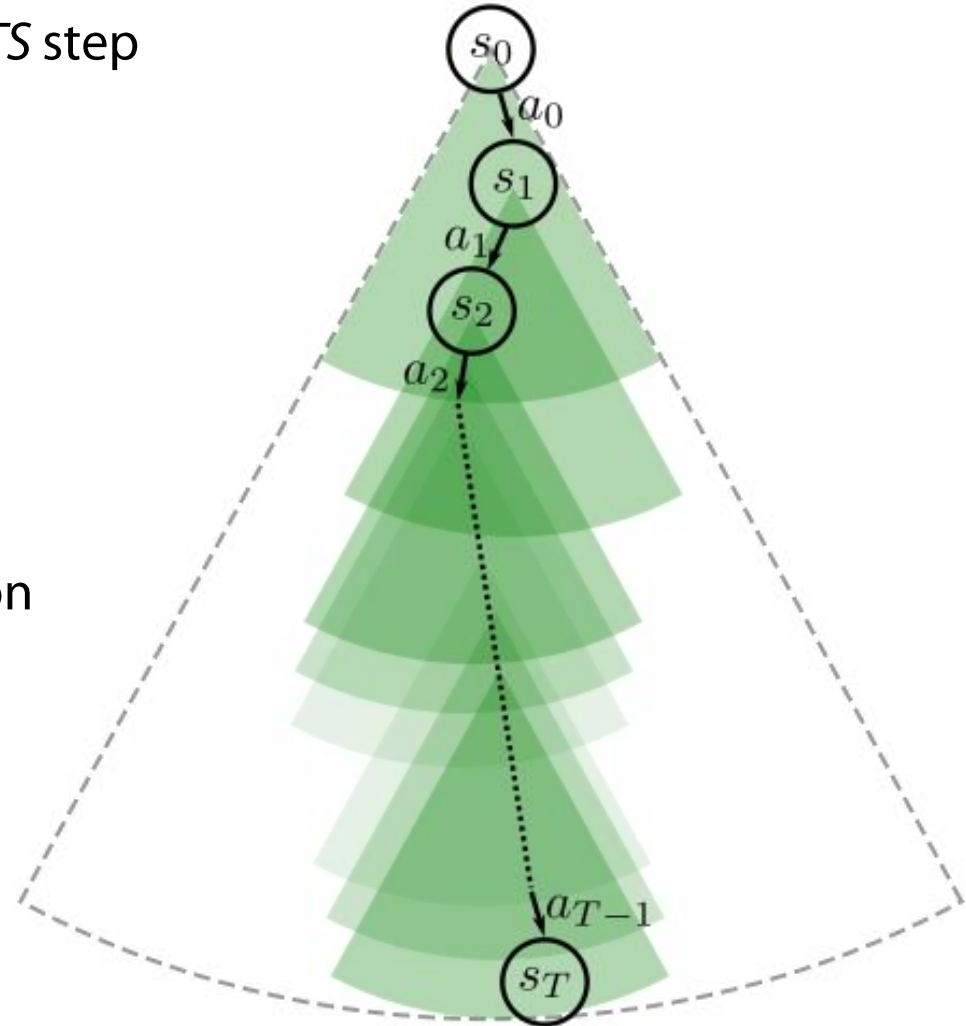
AlphaZero

Marco Piastra

AlphaZero = MCTS + DNN

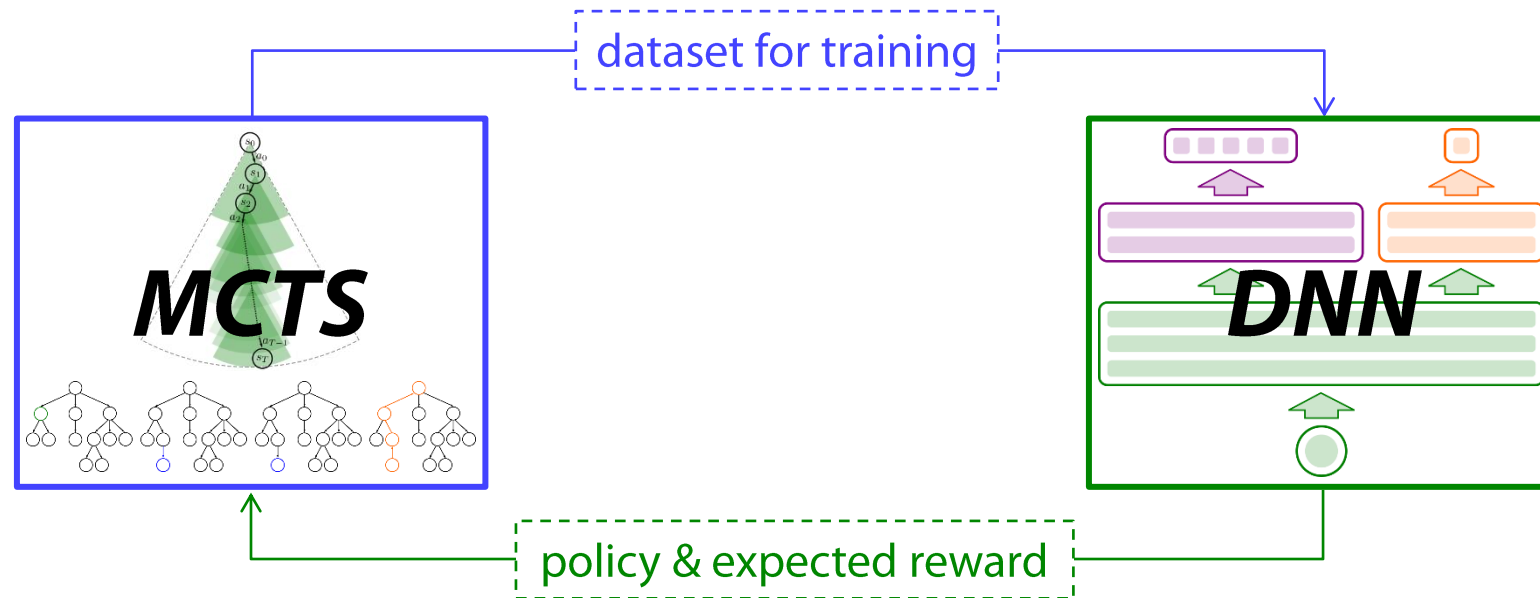
Monte Carlo Tree Search (MCTS) method

- **MCTS** method:
 - memory of past playouts in a single MCTS step
(collected in the tree statistics)
 - knowledge transfer between MCTS steps
(by reusing subtrees already explored)
 - optimal policy only partially defined
(on actually computed states)
 - intrinsically stochastic policy optimization
(the same initial state
can give rise to different optimizations)
 - What about knowledge transfer
between MCTS episodes?
transferring the entire MCTS tree
would rapidly cause its explosive growth...



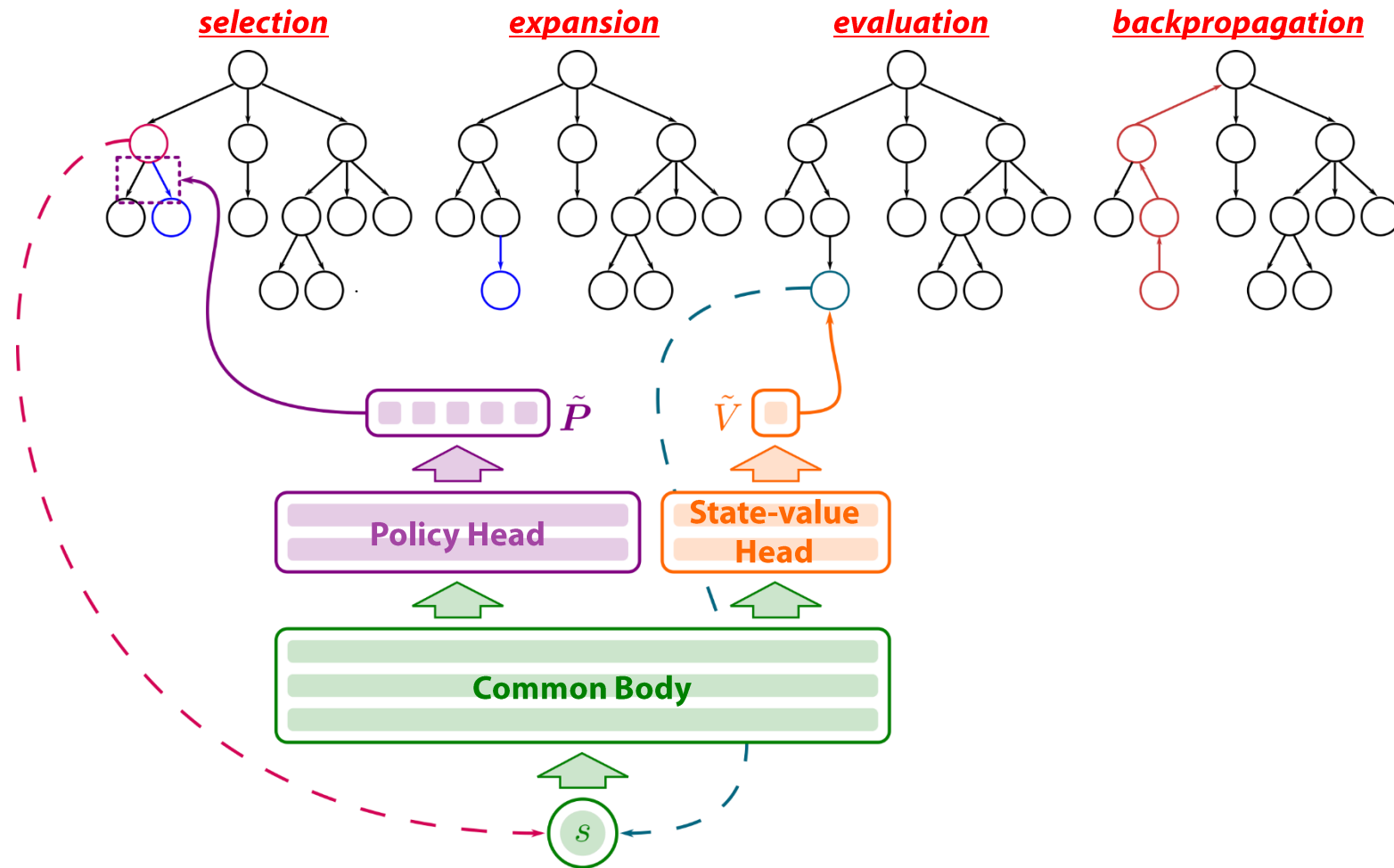
Knowledge transfer between MCTS episodes

- **AlphaZero** [Silver et al. 2017]
 - Monte Carlo Tree Search (MCTS):
improves the policy by focusing on the most promising actions
 - Deep Neural Network (DNN):
learns the improved policy and transfers it between MCTS episodes



AlphaZero

- **AlphaZero = MCTS + DNN**



DNN in AlphaZero

▪ DNN in AlphaZero

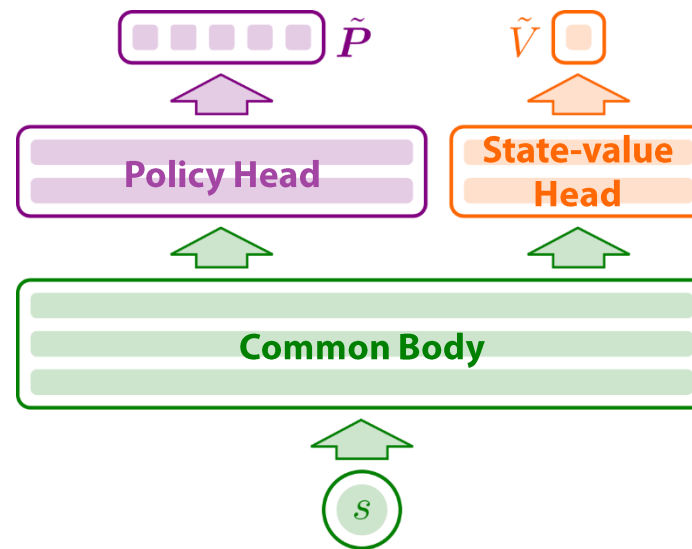
- input: a state s

- output: a probability distribution $\tilde{\mathbf{P}}(s) := [\tilde{P}(a \mid s)]_{a \in \mathcal{A}(S)}$

stochastic policy (a vector of probabilities)

and a *state-value* $\tilde{V}(s)$

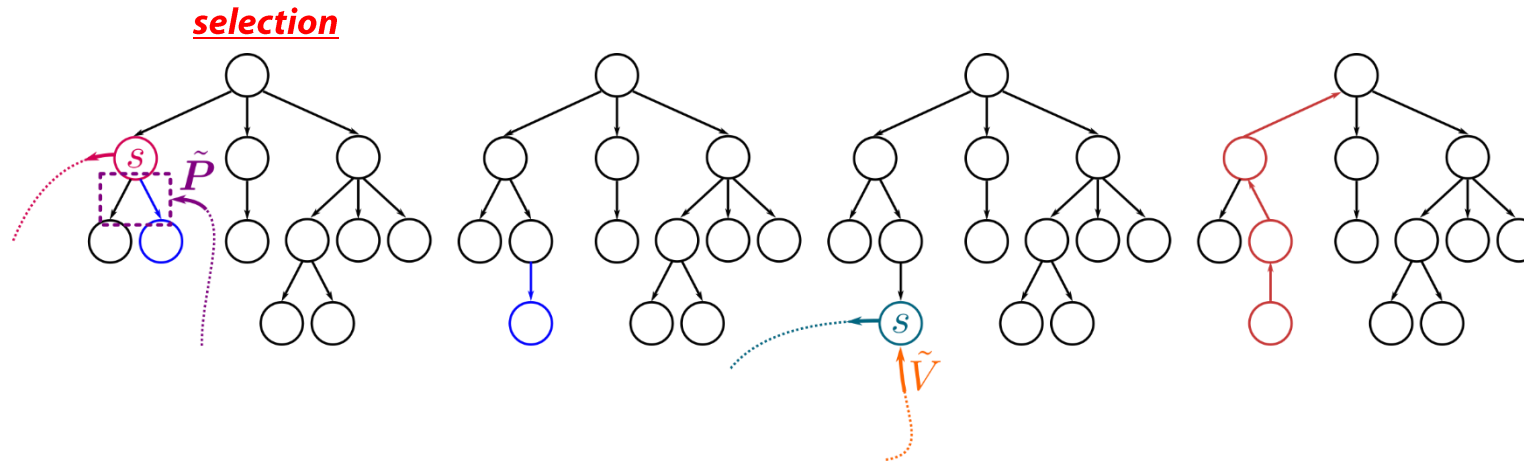
- predicts the expected reward for state s
- acts as an **actor-critic** in the training of **parameters** ϑ of the net



\tilde{V} is compared with the *actual* reward r , which also impacts on training $\tilde{\mathbf{P}}$ by backpropagating through the *Common Body*

MCTS step in AlphaZero

▪ MCTS step in AlphaZero



- selection: UCT policy is replaced with **PUCT** ("Predictor" + UCT)

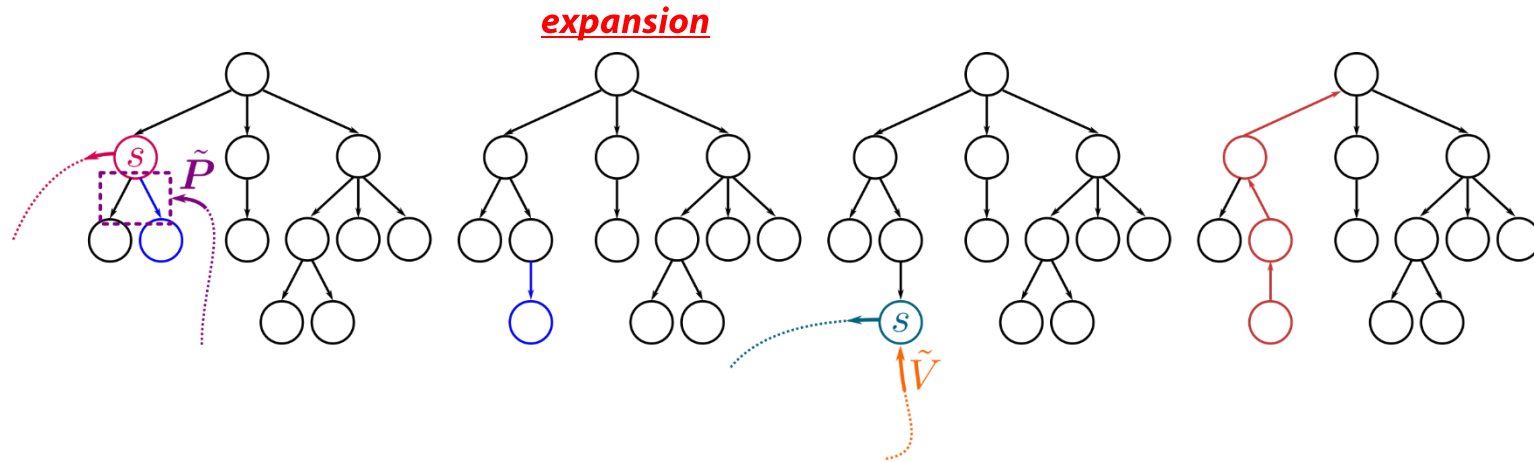
MCTS estimation of $Q(s, a)$ for DNN policy

$$\pi^{\text{PUCT}}(s) := \operatorname{argmax}_a \left\{ \hat{Q}(s, a) + c(s) \tilde{P}(a | s) \frac{\sqrt{N(s)}}{N(s, a) + 1} \right\}$$

exploration rate $c(s) := \log \frac{1 + N(s) + c_{\text{base}}}{c_{\text{base}}} + c_{\text{init}}$
(slowly grows with search time)
avoids division by 0

MCTS step in AlphaZero

▪ MCTS step in AlphaZero

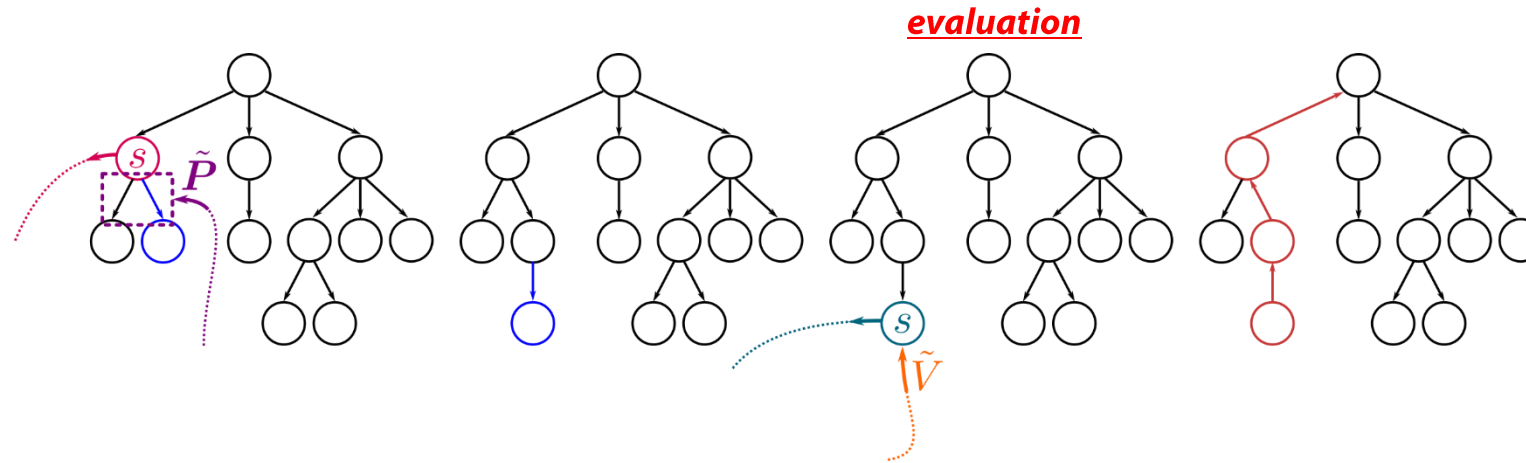


- expansion: initialization of the leaf new node s_L :

$$N(s_L) := 0 \quad \text{and} \quad \forall a \in \mathcal{A}(s_L) \quad N(s_L, a_L) := 0, \quad \hat{Q}(s_L, a_L) := -\infty$$

MCTS step in AlphaZero

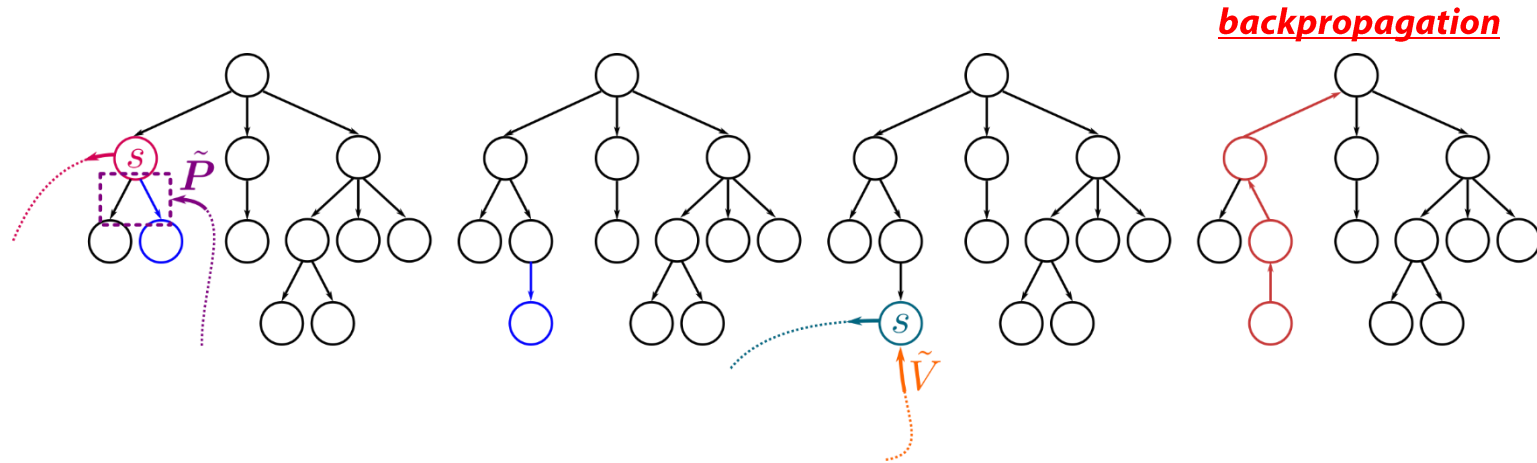
▪ MCTS step in AlphaZero



- expansion: initialization of the leaf new node s_L :
$$N(s_L) := 0 \quad \text{and} \quad \forall a \in \mathcal{A}(s_L) \quad N(s_L, a_L) := 0, \quad \hat{Q}(s_L, a_L) := -\infty$$
- evaluation (in place of simulation): expected reward is $\tilde{V}(s_L)$

MCTS step in AlphaZero

▪ MCTS step in AlphaZero



- expansion: initialization of the leaf new node s_L :

$$N(s_L) := 0 \quad \text{and} \quad \forall a \in \mathcal{A}(s_L) \quad N(s_L, a_L) := 0, \quad \hat{Q}(s_L, a_L) := -\infty$$

- evaluation (in place of simulation): expected reward is $\tilde{V}(s_L)$
- backpropagation: for each state s and action a visited in selection/expansion:

$$\begin{aligned} N(s) &:= N(s) + 1, \\ N(s, a) &:= N(s, a) + 1 \quad \text{and} \quad \hat{Q}(s, a) := \hat{Q}(s, a) + \frac{\tilde{V}(s_L) - \hat{Q}(s, a)}{N(s, a)} \end{aligned}$$

MCTS step in AlphaZero: policies

- Selection policy: **PUCT**

$$\pi^{\text{sel}}(s) := \pi^{\text{PUCT}}(s) := \operatorname{argmax}_a \left\{ \hat{Q}(s, a) + c(s) \tilde{P}(a \mid s) \frac{\sqrt{N(s)}}{N(s, a) + 1} \right\}$$

- Output policy:

$$\pi^{\text{out}}(s) \sim \left[\hat{P}(a \mid s) := \frac{N(s, a)}{N(s)} \right]_{a \in \mathcal{A}(s)}$$

taking frequencies as probabilities
(in place of their argmax as output action)
ensures **exploration**

(the simulation policy does not exist anymore)

DNN training in AlphaZero

- **Data items** from a single MCTS episode:

After an MCTS episode $\mathcal{E} := \langle s_0, a_0, s_1, \dots, a_{T-1}, s_T \rangle$
with actual reward $\hat{V}^{\mathcal{E}} := r(s_T)$:

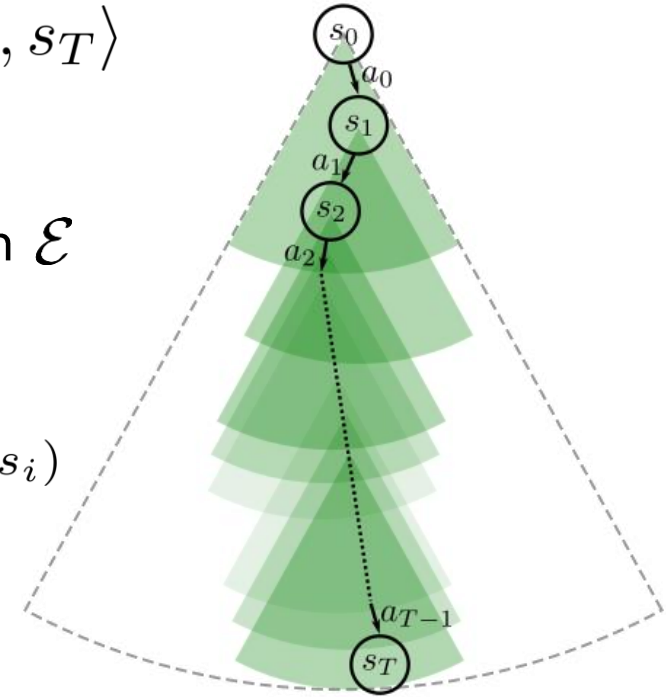
- for each non-terminal state s_i ($i = 0 \dots T - 1$) in \mathcal{E}

$$\hat{\mathbf{P}}(s_i) := \left[\hat{P}(a \mid s_i) := \frac{N(s_i, a)}{N(s_i)} \right]_{a \in \mathcal{A}(s_i)}$$

vector of frequencies

- the **output** of \mathcal{E} is

$$D^{\mathcal{E}} := \left\{ \underbrace{\langle s_i, \hat{\mathbf{P}}(s_i), \hat{V}^{\mathcal{E}} \rangle}_{\text{data item}} \right\}_{i=0 \dots T-1}$$



DNN training in AlphaZero

- **Iteration:**

K times $\left\{ \begin{array}{l} 1) \text{ play one MCTS episode } \mathcal{E}_j \\ 2) \text{ collect data items } D^{\mathcal{E}_j} \end{array} \right.$

3) train the parameters of the DNN by using as **dataset**

$$D := \bigcup_{j=1}^K D^{\mathcal{E}_j}$$

- In the limit of *infinite* iterations:

$$\pi^{\text{DNN}}(s) := \underset{a \in \mathcal{A}(s)}{\operatorname{argmax}} \tilde{P}(a \mid s) \rightarrow \pi^*(s) \quad \forall s$$

AlphaZero

▪ AlphaZero:

- memory of past playouts in a single MCTS step
(collected in the tree statistics)
- knowledge transfer between MCTS steps
(by reusing subtrees already explored)
- knowledge transfer between MCTS episodes
(provided by DNN)
- deterministic policy optimization
with policy defined for all states s :

$$\pi^{\text{DNN}}(s) := \operatorname{argmax}_{a \in \mathcal{A}(s)} \tilde{P}(a \mid s)$$

