$$\begin{pmatrix} \tilde{b} = \phi(b) \\ \text{belief representation} \end{pmatrix}$$

$n$ parallel MCTS simulations

$(\mathbf{p}, v) = f_\theta(\tilde{b})$

initial network

$f_\theta = f'_\theta$

**policy evaluation**

**policy improvement**

$f'_\theta = \text{TRAIN}(f_\theta, \mathcal{D})$

$\mathcal{D} = \left\{ \left\{ (b_t, \boldsymbol{\pi}_t, g_t) \right\}_{t=1}^{T} \right\}_{i=1}^{n}$

collected data

SELECTION    EXPANSION    SIMULATION    BACKPROPAGATION

(where $\hat{b} \leftarrow \phi(b)$)

$s \sim \hat{b}$
$s' \sim T(\cdot \mid s, a)$
$o \sim O(\cdot \mid a, s')$

$b' \leftarrow \text{UPDATE}(b, a, o)$

$r + \gamma V_\theta(\tilde{b}')$

Q-value