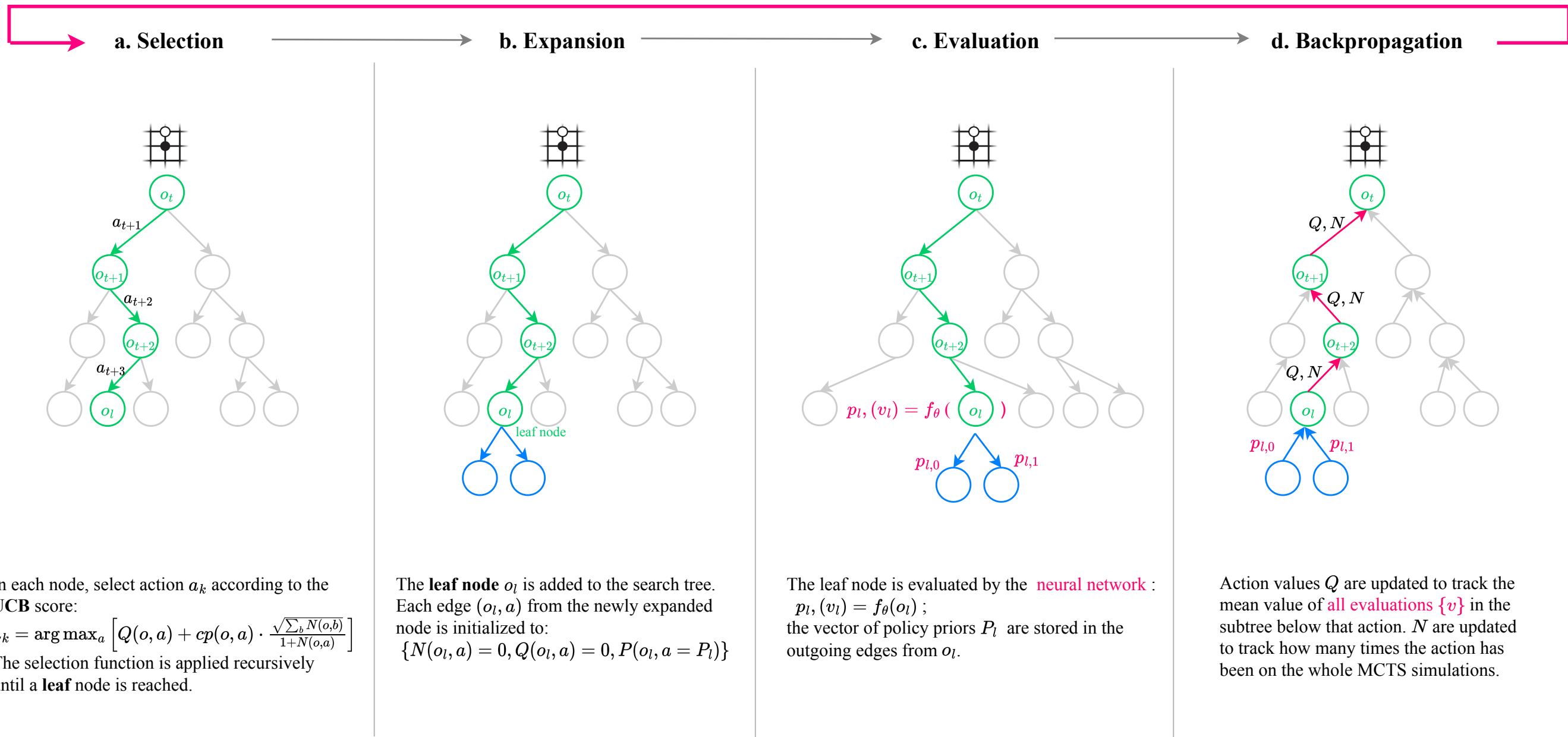
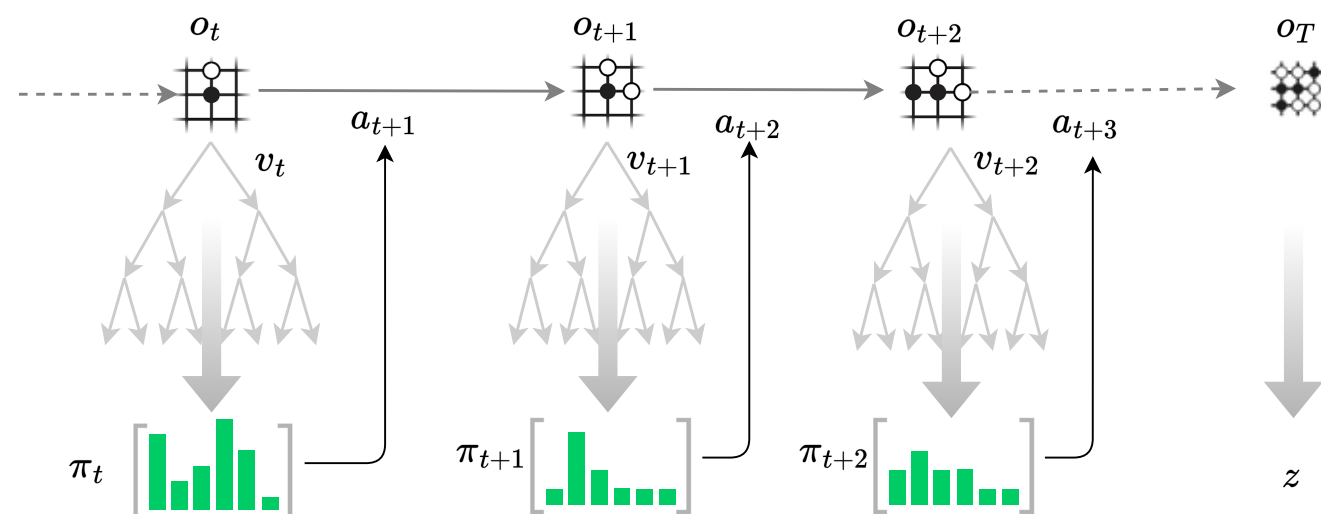


## A. MCTS in AlphaZero



## B. Acting (Self-play)



## C. Training

$$l_t(\theta) = l^v(z_t, \mathbf{v}_t) + l^p(\pi_t, \mathbf{p}_t) + c\|\theta\|^2$$

where,  $z_t$  is the game reward from the **perspective of the current player** and  $\pi_{t+k}$  is the MCTS searched policy at timestep  $t$ .

$\mathbf{v}_t$  and  $\mathbf{p}_t$  is the predicted value and policy from the neural network  $f_\theta$ .

$l^v$  is MSE loss,  $l^p$  is cross-entropy loss.