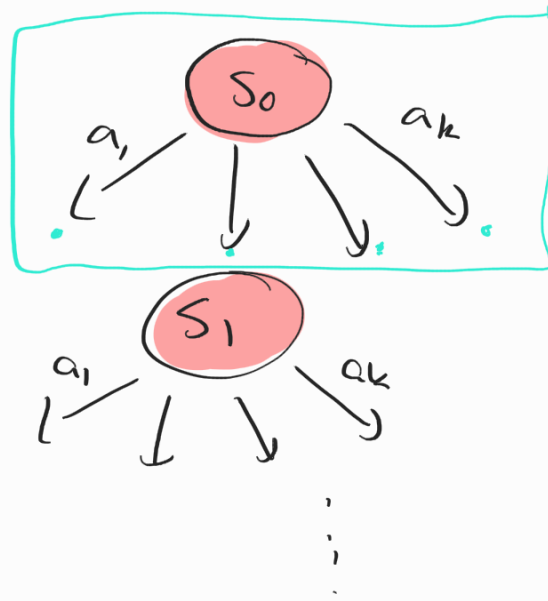


## ② Monte Carlo Tree Search (MCTS)

UCT  $\rightarrow$  extension of UCB

$\hookrightarrow$  Upper Confidence (Bounds) for Trees



$\leftarrow$  a multi-armed bandit question

We are essentially solving nested multi armed bandit problems !

What do we learn ?

	$v$ -values	$\equiv$	state value function
	$q$ -values	$\equiv$	action state value function

Data structure :

Plays :  $S \rightarrow \mathbb{N}$

$s \in S$  Plays(s) how many plays I have simulated that contained s (?)

Values:  $S \rightarrow \mathbb{R}$

values(s) are current estimates for  $val_*(s)$

Play(s) : return an action to play from s.

① Naive Version :

We could simply "brute"  
our value function:

$$\text{play} \quad \arg \max_{a \in A} \underbrace{\sum_{s', r} \Delta(s, a)(s', r) (r + \gamma \text{val}(s'))}_{q(s, a)}$$

② Actual algorithm:

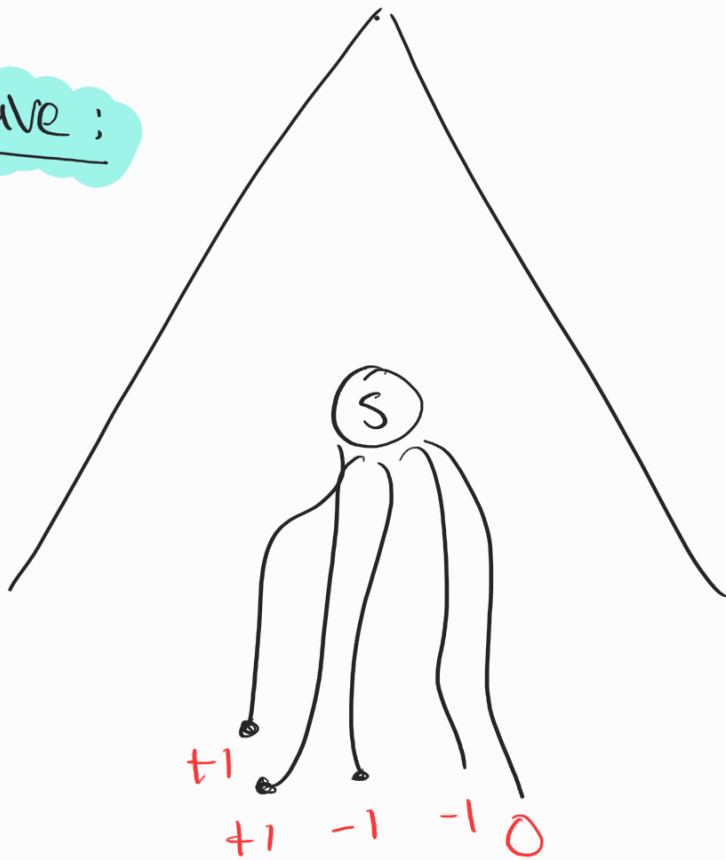
(1) Run a number of simulations from s  
(set a timer)

(2) Now I can trust my value

function, and do as in the naive version above for choosing the action.

Simulation (s) :   
 Repeat   
 { construct a trajectory from s   
 , update our statistics using the trajectory

Naive:



val(s) updated using  $(+1, +1, -1, -1, 0)$

Actual MCTS

4 phases:

as long as we have value estimates for all successors

SELECTION : as long as possible, use the statistics (values)

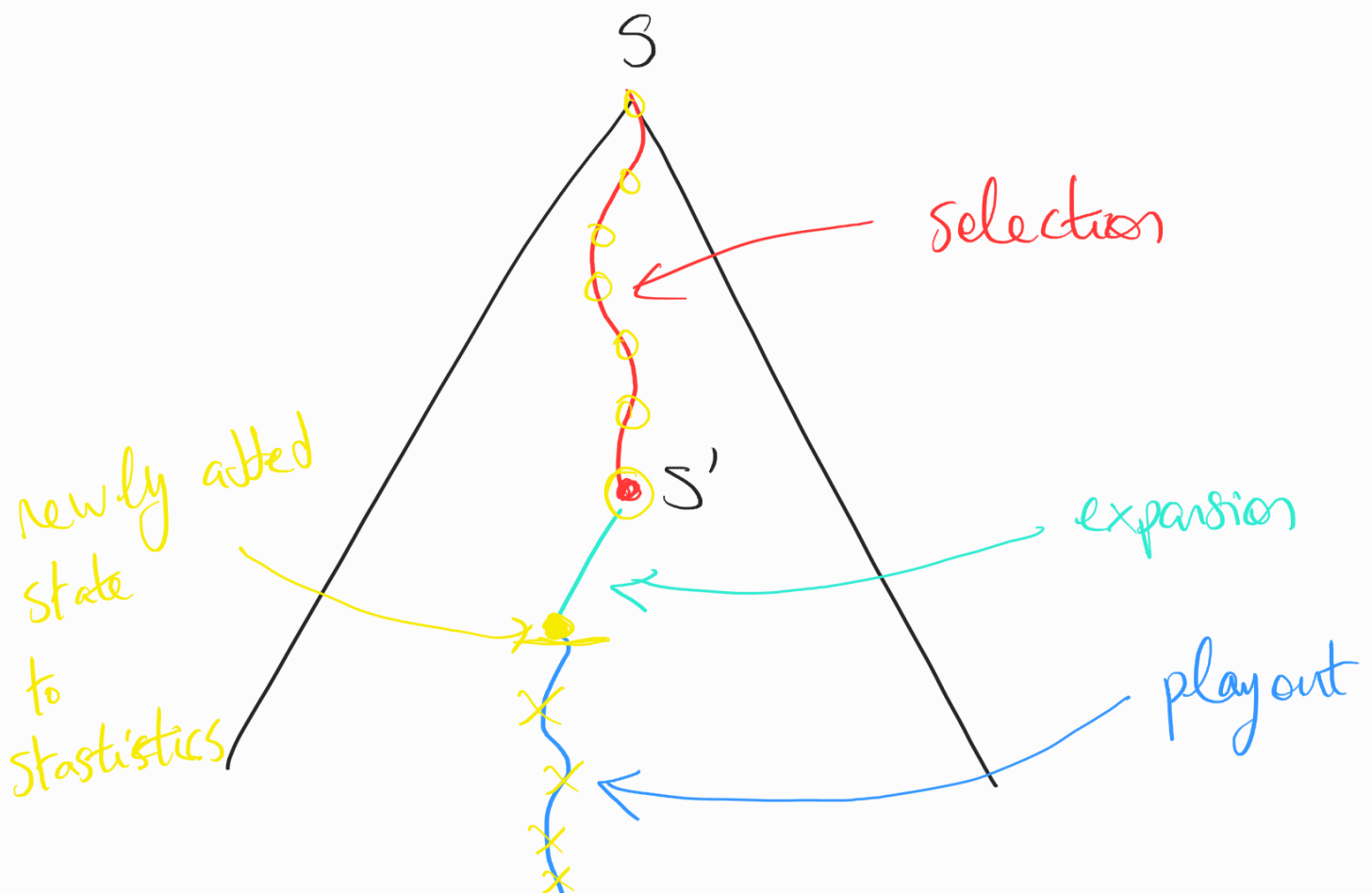
using  $\epsilon$ -greedy or UCB

EXPANSION:  
pick a random successor for which  
we don't have value estimate

PLAYOUT:  
play randomly until trajectory is complete

BACKPROPAGATION

update values of all states visited  
DURING SELECTION AND  
EXPANSION



outcome  
+1

SE

Plays(s)  $t = 1$

Values(s) {  
• average  
• step size  
• temporal difference