



START: Initialize Simulation (RL 2000 Runs)

Run Initialization
set parameter: $K=10000$, $\alpha=0.01$ (UCB constant), $T=1000$ (steps)
Initialize Arrays: $Q = [0, \dots, 0]$, $N = [0, \dots, 0]$ (Q & N are 1D)
Initialize Environment: q -states $N(0, 1)$ (stationary True values)

is $t < K$?
(force initialize phase)

yes
select Action A_{t+1}
(pull each arm once)

no
calculate UCB for all q ?
 $UCB(q) = Q_t(q) + \alpha \sqrt{\ln(t)/N_t(q)}$

select Action $A_{t+1} = \arg\max (UCB)$

Reward & Incremental Update
Get Reward: $R_t \sim N(q\text{-state}[A_t]; 1)$
Update Count: $N_t(A_t) = N_t(A_t) + 1$
Update estimate: $Q_{t+1}(A_t) = Q_t(A_t) + (\alpha / N_t(A_t)) * (R_t - Q_t(A_t))$

Record Metrics: R_t is optimal (A_t is argmax q -state)

is $t < T$?

no
End Run & Aggregate

data



Plot Performance Parameters (Average over runs)
Plot 1: Average Reward vs Steps
Plot 2: % optimal action vs runs



end: Analysis Results