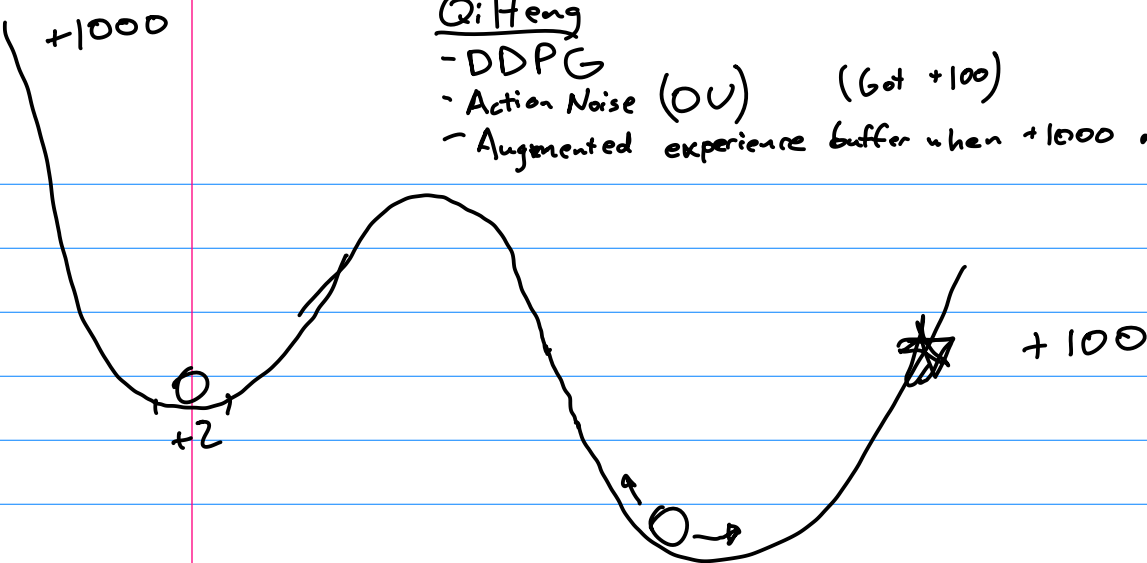


Qi Heng

- DDPG
- Action Noise (0.1) (Got +100)
- Augmented experience buffer when +1000 received



Last Time

The POMDP is PSPACE-Complete

- exact algorithms
(probably) will have exponential complexity

Numerical Approx

Offline
SARSOP

Online

This Time

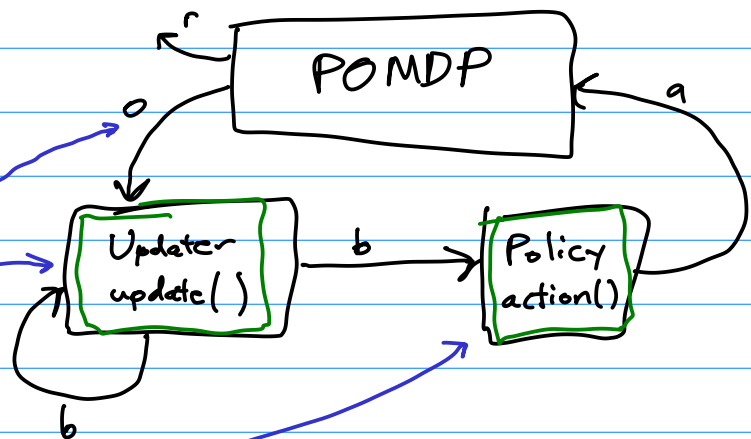
Formulation Approx.

eg. QMDP

First POMDPs.jl for HW 6

```
b = initialize_belief(up, isd)

r_total = 0.0
d = 1.0
while !isterminal(pomdp, s)
    a = action(policy, b)
    s, o, r = gen(DDNOut(:sp,:o,:r), pomdp, s, a, rng)
    r_total += d*r
    d *= discount(pomdp)
    b = update(up, b, a, o)
end
```



$$\xrightarrow{\text{QMDP}} \alpha_a^{(k+1)}[s] = R(s, a) + \gamma \sum_{s'} T(s'|s, a) \underbrace{\max_{a'} \alpha_a^{(k)}[s']}_{Q_{\text{MDP}}(s', a')}$$

$$\text{HW2} \rightarrow V_{\text{MDP}}^{(k+1)}(s) = \max_a \left(R(s, a) + \gamma \sum_{s'} T(s'|s, a) V_{\text{MDP}}^{(k)}(s') \right)$$

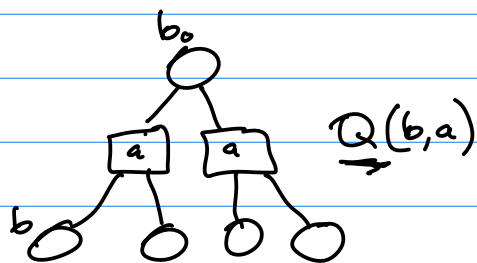
$$V_{\text{MDP}}^*$$

$$\alpha_a[s] = R(s, a) + \gamma \sum_{s'} T(s'|s, a) V_{\text{MDP}}^*(s')$$

$$\pi_{\text{QMDP}}(b) = \underset{a}{\operatorname{argmax}} \underbrace{E[Q_{\text{MDP}}(s, a)]}_{\approx}$$

Online POMDP Methods

	Name	S	A	O	T	Z
2007	AEMS <small>Online Planning Algs for POMDPs</small>	D	D	D	E	E
2010	MCTS PO-UCT POMCP	C	D	D	G	G
2013	DESPOT	C	D	D	G	G
2018	POMCPow	C	D (1-2L)	C	G	E
2019	DESPOT- α	C	D	C	G	E
2020	BOMCP	C	C	C	G	E
2021	VOMCPow	C	C	C	G	E



Online
Only tryin to
find a good action
for current
belief

DESPOT
 $Q(b, a)$ only

Offline
Good action for
all beliefs

SARSOP
 α -vectors

AEMS

while time remains

$$b^* = \underset{b \in \text{Fringe}(G)}{\text{argmax}} E(b)$$

expand (b^*)

backup (b^w)

$$E(b) = \gamma^d P(b^d) \hat{E}(b^d)$$

$$\hat{\mathcal{E}}(b) = U(b) - \mathcal{L}(b)$$

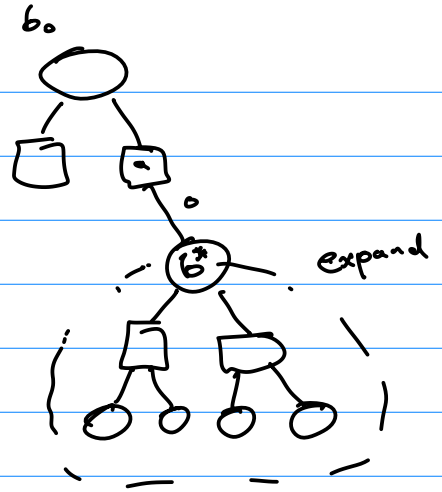
$$P(b^d) = \prod_{i=0}^{d-1} P(o^i | b^i, a^i) P(a^i | b^i)$$

$$P(a|b) = \frac{U(a,b) - L(b)}{U(b) - L(b)}$$

$$P(a|b) = \begin{cases} 1 & \text{if } a = \underset{\substack{\text{argmax} \\ a' \in A}}{a'} \\ 0 & \text{otherwise} \end{cases} \quad U(a, b)$$

AEMS 1

AEMS 2

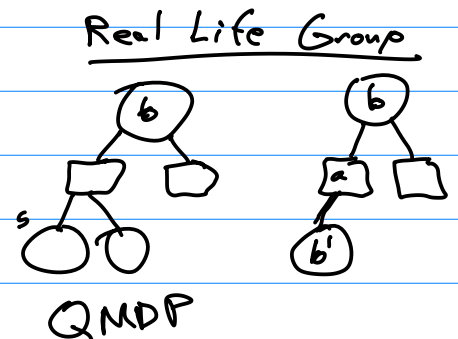
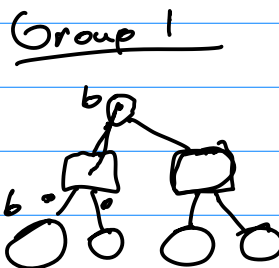
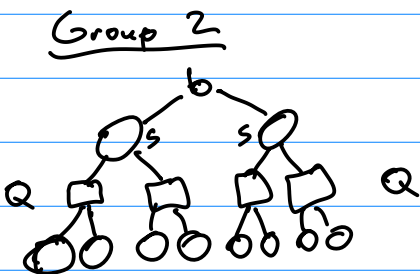
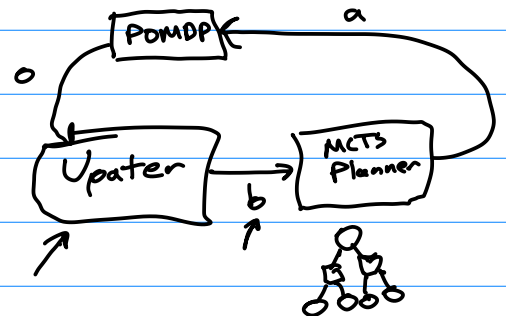


Breakout Rooms

```
function simulate!( $\pi$ ::MonteCarloTreeSearch, s, d= $\pi$ .d)
    if d  $\leq$  0
        return 0.0
    end
     $\mathcal{P}$ , N, Q, c =  $\pi$ . $\mathcal{P}$ ,  $\pi$ .N,  $\pi$ .Q,  $\pi$ .c
     $\mathcal{A}$ , TR,  $\gamma$  =  $\mathcal{P}$ . $\mathcal{A}$ ,  $\mathcal{P}$ .TR,  $\mathcal{P}$ . $\gamma$ 
    if !haskey(N, (s, first( $\mathcal{A}$ )))
        for a in  $\mathcal{A}$ 
            N[(h,a)] = 0
            Q[(h,a)] = 0.0
        end
    end
    return rollout( $\mathcal{P}$ , s,  $\pi$ . $\pi$ , d)
end

a = explore( $\pi$ , s)
s', r = TR(s,a)
q = r +  $\gamma$ *simulate!( $\pi$ , s', d-1)
N[(h,a)] += 1
Q[(h,a)] += (q-Q[(h,a)]) / N[(h,a)]
return q
```

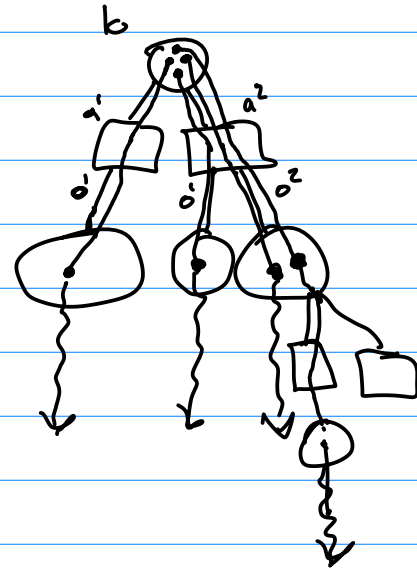
How can we modify MCTS
to work for POMDPs
avoiding explicit Bayesian
Belief Updates?
($b, a_1, o_1, a_2, o_2 \dots o_n$)



Do Work { MCTS
PO-UCT

need $h' \leftarrow G(h, a)$ have $s', o, r \leftarrow G(s, a)$

```
function simulate( $\pi$ ::HistoryMonteCarloTreeSearch, s, h, d)
  if d ≤ 0
    return 0.0
  end
   $\mathcal{P}, N, Q, c = \pi.\mathcal{P}, \pi.N, \pi.Q, \pi.c$ 
   $S, \mathcal{A}, TRO, \gamma = \mathcal{P}.S, \mathcal{P}.\mathcal{A}, \mathcal{P}.TRO, \mathcal{P}.\gamma$ 
  if !haskey(N, (h, first( $\mathcal{A}$ )))
    for a in  $\mathcal{A}$ 
       $N[(h, a)] = 0$ 
       $Q[(h, a)] = 0.0$ 
    end
     $b = [s == s' ? 1.0 : 0.0 \text{ for } s' \text{ in } S]$ 
    return rollout( $\mathcal{P}, b, \pi.\pi, d$ )
  end
  a = explore( $\pi, h$ )
   $s', r, o = TRO(s, a)$ 
   $q = r + \gamma \cdot \text{simulate}(\pi, s', \text{vcut}(h, (a, o)), d-1)$ 
   $N[(h, a)] += 1$ 
   $Q[(h, a)] += (q - Q[(h, a)]) / N[(h, a)]$ 
  return q
end
```

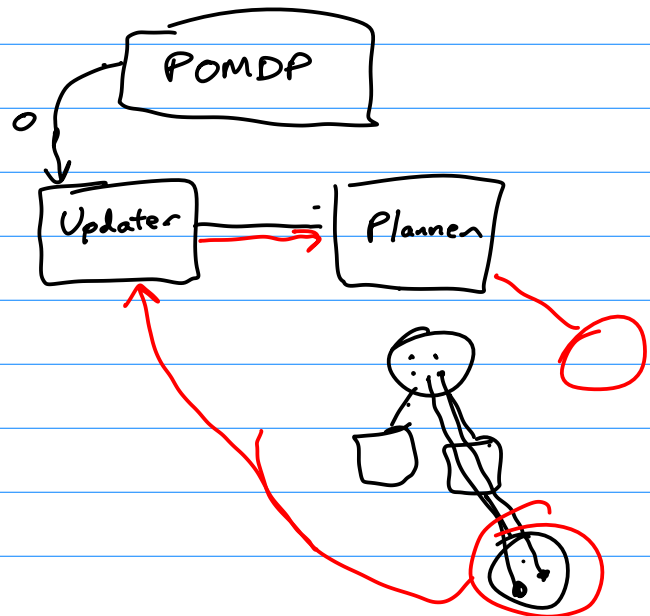


10 actions
10 observations

POMCP

Re-use planning
simulations for interaction
loop belief updates

Does this work
in practice? No



DESPOT

not MCTS; Heuristic Search

K scenarios

ϕ_i generate beforehand

$$s' \leftarrow G(s, a, \phi_i[i])$$

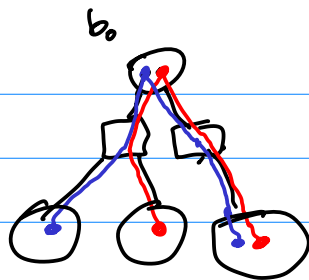
Example $s' = s + a + w$

$$\phi_i[i] = 2$$



$$G(1, 1, \phi_i[i]) = 1 + 1 + 2 = 4$$

$$G(1, 2, \phi_i[i]) = 1 + 2 + 2 = 5$$



weighted, normalized

$$U(b)$$

μ

$$L(b)$$

$$a^* = \arg \max_a \mu(b, a)$$

$$o^* = \arg \max_o E(\tau(b, a, o))$$

↑ weighted excess