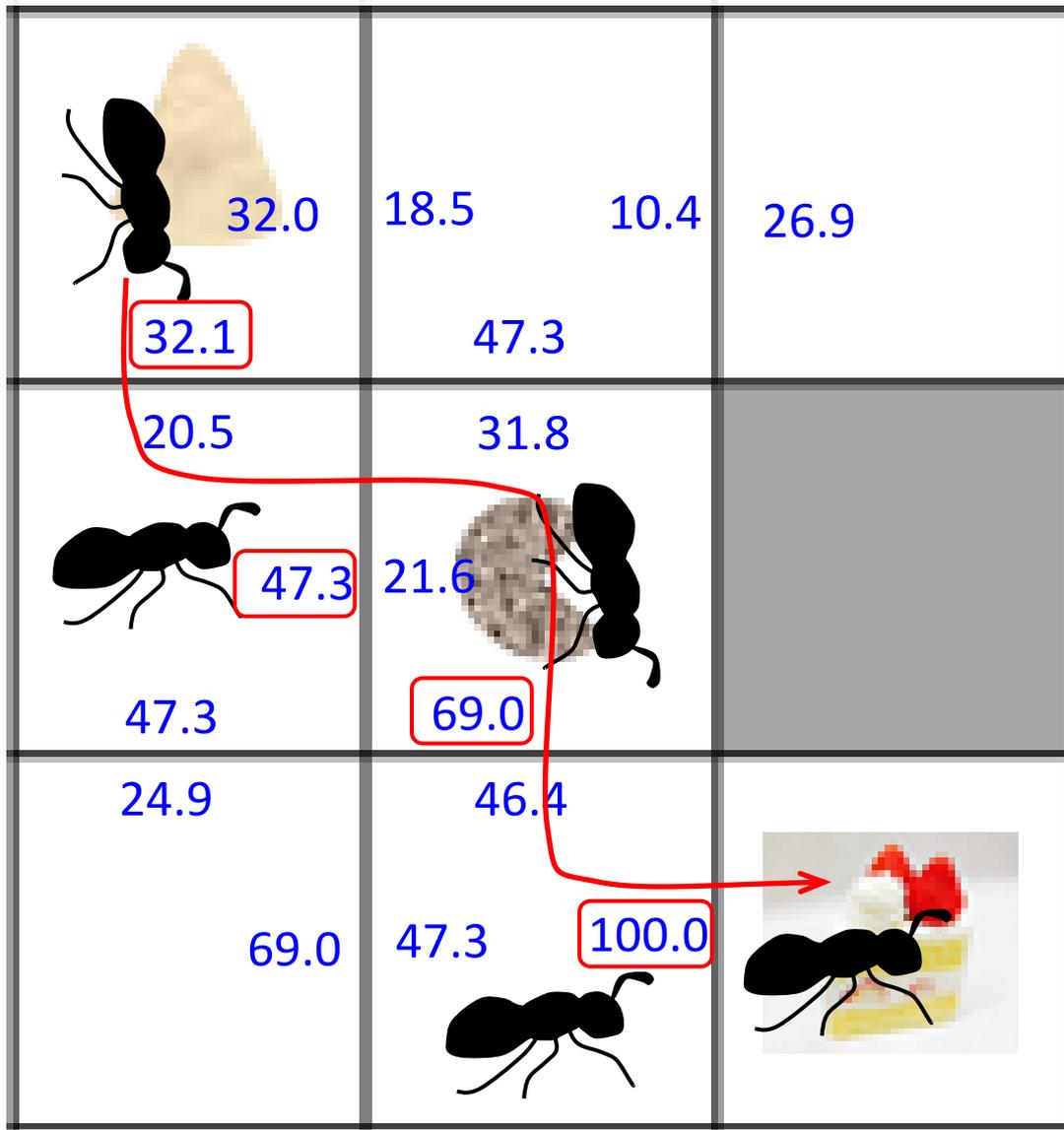


J	K	L	M	N	O	P	Q	R	S	T
step 1 スタート										
現Agent位置										
行	1									
列	1									
状態	1 (s.)									
★										
		終								
状態	初期値	アクション								
	現Q値	右	上	左	下					
	1	3.00	欄外	欄外	1.00					
	2	4.00	欄外	1.00	5.00					
	3	欄外	欄外	9.00	欄外					
	4	2.00	6.00	欄外	5.00					
	5	欄外	2.00	7.00	1.00					
	6	欄外	欄外	欄外	欄外					
	7	8.00	2.00	欄外	欄外					
	8	8.00	1.00	8.00	欄外					
9(到着)	0	0	0	0						
行動決定										
ε	0.98	G乱数	0.175	→	Explore					
		最大Q	E乱数	Act候補						
		Exploit時	3.00		1					
		Explore時		0.800	4					
採用Act	4 (a.)	(下) →								
		次Agent位置								
		行	2							
		列	1							
		次状態	4 (s.)							
Q更新の計算										
r+y MaxQ	3.20									
状態	新Q値	アクション								
		右	上	左	下					
	1	3.00	欄外	欄外	2.10					
	2	4.00	欄外	1.00	5.00					
	3	欄外	欄外	9.00	欄外					
	4	2.00	6.00	欄外	5.00					
	5	欄外	2.00	7.00	1.00					
	6	欄外	欄外	欄外	欄外					
	7	8.00	2.00	欄外	欄外					
	8	8.00	1.00	8.00	欄外					
9(到着)	0	0	0	0						

J	V	W	X	Y	Z	AA	AB	AC	AD	AE	A
step 2											
現Agent位置											
行	2										
列	1										
状態	4 (s.)										
		終									
状態	初期値	アクション									
	現Q値	右	上	左	下						
	1	3.00	欄外	欄外	2.10						
	2	4.00	欄外	1.00	5.00						
	3	欄外	欄外	9.00	欄外						
	4	2.00	6.00	欄外	5.00						
	5	欄外	2.00	7.00	1.00						
	6	欄外	欄外	欄外	欄外						
	7	8.00	2.00	欄外	欄外						
	8	8.00	1.00	8.00	欄外						
9(到着)	0	0	0	0							
行動決定											
ε	0.98	G乱数	0.918	→	Explore						
		最大Q	E乱数	Act候補							
		Exploit時	6.00		2						
		Explore時		0.800	4						
採用Act	4 (a.)	(下) →									
		次Agent位置									
		行	3								
		列	1								
		次状態	7 (s.)								
Q更新の計算											
r+y MaxQ	4.60										
状態	新Q値	アクション									
		右	上	左	下						
	1	3.00	欄外	欄外	2.10						
	2	4.00	欄外	1.00	5.00						
	3	欄外	欄外	9.00	欄外						
	4	2.00	6.00	欄外	4.80						
	5	欄外	2.00	7.00	1.00						
	6	欄外	欄外	欄外	欄外						
	7	8.00	2.00	欄外	欄外						
	8	8.00	1.00	8.00	欄外						
9(到着)	0	0	0	0							

数字はQ値



ゴールまでのステップ数

