



## Reinforcement Learning Process

- 1. Initialize** — Agent starts with no knowledge or policy.
- 2. Explore** — Tries random actions to gather experience and rewards
- 3. Learn** — Updates value estimates or policy based on feedback.
- 4. Improve** — Refines the policy to favor rewarding actions.
- 5. Exploit** — Follows the learned policy to maximize rewards.

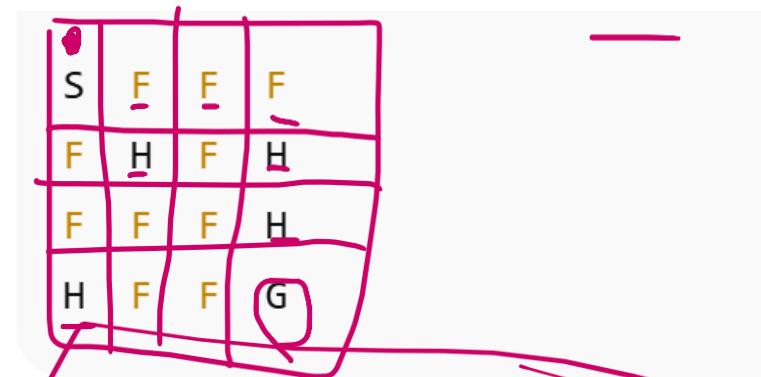
10✓

50✓

100✓

# Frozen Lake

grid



- S = Start (state 0)
- F = Frozen (safe)
- H = Hole (terminal 0 reward)
- G = Goal (terminal +1 reward)

r;  
L;  
U;  
d

Initial State-Value Grid  $V(s) = 0$

1)  
2)

Initial State Values (Episode 1)

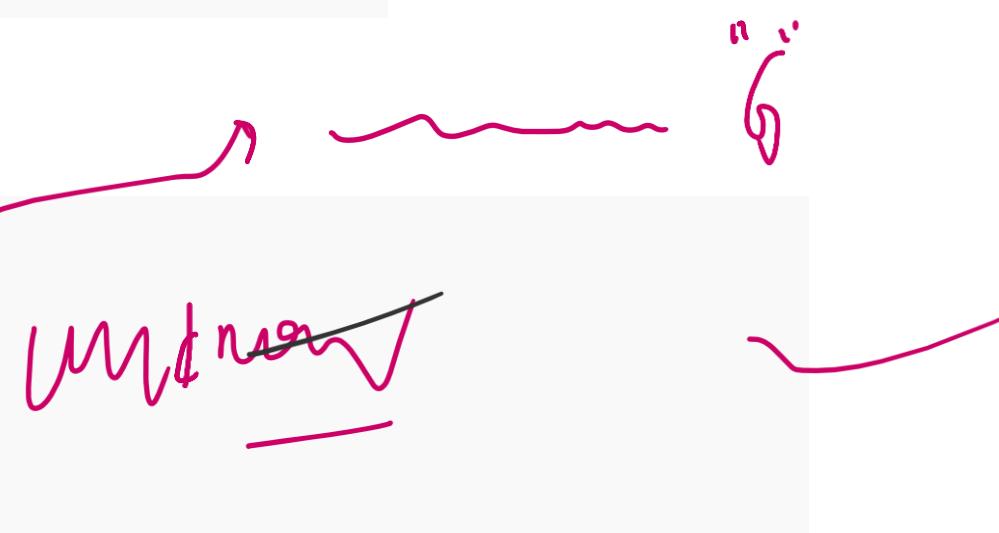
0.00	0.00	0.00	0.00
0.00	H	0.00	H
0.00	0.00	0.00	H
H	0.00	0.00	1.00

Grid - 4x4 - FL

Initial Policy Grid

Initial Policy (Episode 1)

?	?	?	?
?	H	?	H
?	?	?	H
H	?	?	G



- Each "?" means: agent may pick any of {←, ↓, →, ↑} with equal probability.
- Holes (H) and goal (G) are terminal — no actions there.

Safe optimal-looking path

Row\Col	0	1	2	3
0	[S]	[F]	[F]	[F]
1	[F]	[H]	[F]	[H]
2	[F]	[F]	[F]	[H]
3	[H]	[F]	[F]	[G]

Safe optimal-looking path we want the agent to discover (one valid route):

$(0,0) \rightarrow (0,1) \rightarrow (0,2) \rightarrow (1,2) \rightarrow (2,2) \rightarrow (3,2) \rightarrow (3,3)$

Episode 1

0.00	0.00	0.00	0.00
0.00	H	0.00	H
0.00	0.00	0.00	H
H	0.00	0.00	1.00

State-Value

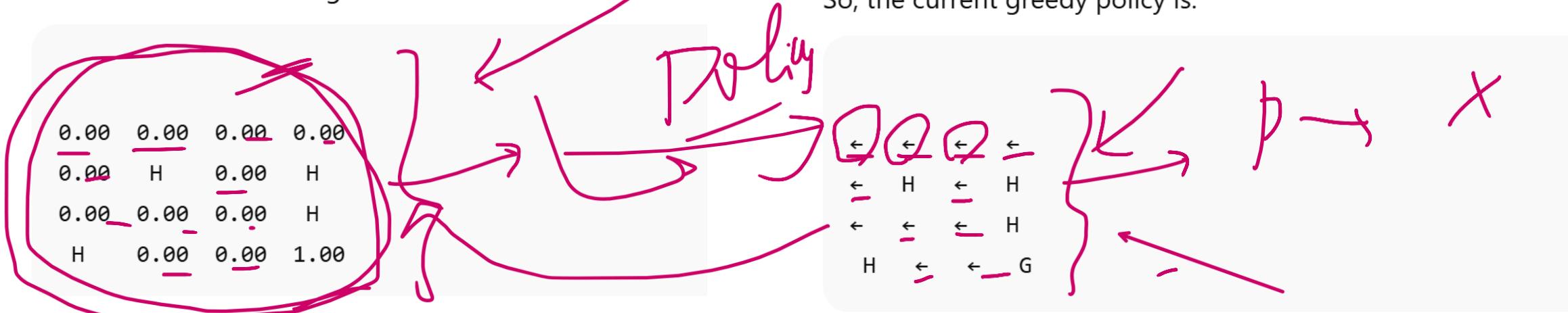
Row\Col	0	1	2	3
0	[S]	[F]	[F]	[F]
1	[F]	[H]	[F]	[H]
2	[F]	[F]	[F]	[H]
3	[H]	[F]	[F]	[G]

Path Taken



No positive reward was received (agent fell into a hole), so values remain unchanged.

All Q-values  $\approx 0 \rightarrow \text{argmax ties} \rightarrow$  arbitrary ( $\leftarrow$  chosen).  
So, the current greedy policy is:



## Episode 2

0.00	0.00	0.00	0.00
0.00	H	0.00	H
0.00	0.00	0.00	H
H	0.00	0.00	1.00

Path Taken

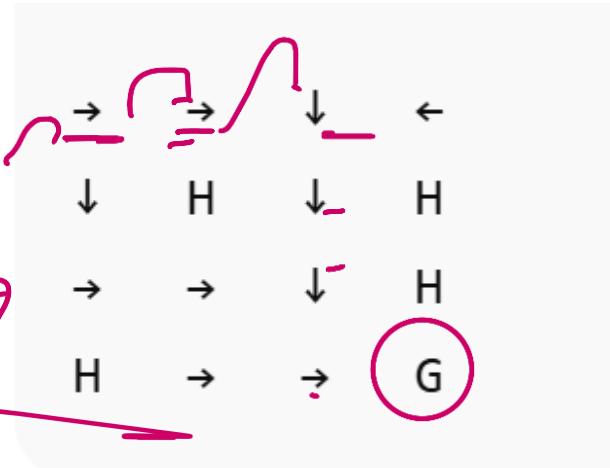
(0,0) → (0,1) → (0,2) → (1,2) → (2,2) → (3,2) → (3,3)  (goal)

Updated State-Value Grid (after episode)

0.30	0.45	0.55	0.10
0.25	H	0.65	H
0.20	0.40	0.80	H
H	0.55	0.90	1.00

Row\Col	0	1	2	3
0	[S]	[F]	[F]	[F]
1	[F]	[H]	[F]	[H]
2	[F]	[F]	[F]	[H]
3	[H]	[F]	[F]	[G]

Policy Grid (after episode)

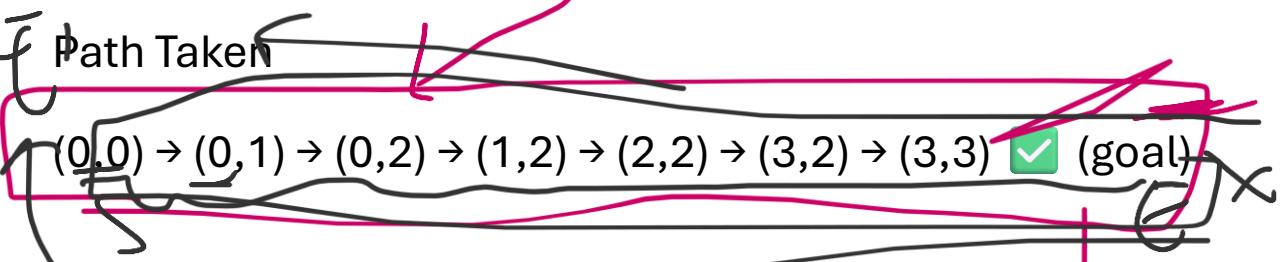


Values increase closer to the goal because those states lead reliably to the +1 reward.

Episode 3

(Carries over from Episode 2—the agent already knows a good path.)

0.30	0.45	0.55	0.10
0.25	H	0.65	H
0.20	0.40	0.80	H
H	0.55	0.90	1.00



Updated State-Value Grid (after episode)

0.45	0.60	0.70	0.10
0.40	H	0.75	H
0.35	0.55	0.85	H
H	0.70	0.92	1.00

Row\Col	0	1	2	3
0	[S]	[F]	[F]	[F]
1	[F]	[H]	[F]	[H]
2	[F]	[F]	[F]	[H]
3	[H]	[F]	[F]	[G]

$$V_{st_j} = .30 +$$

Policy Grid (after episode)

$\rightarrow$	$\rightarrow$	$\downarrow$	$\leftarrow$
$\downarrow$	H	$\downarrow$	H
$\rightarrow$	$\rightarrow$	$\downarrow$	H
H	$\rightarrow$	$\rightarrow$	G

## Episode 4

0.45	0.60	0.70	0.10
0.40	H	0.75	H
0.35	0.55	0.85	H
H	0.70	0.92	1.00

## Path Taken

Agent mostly exploits but explores a small detour then returns to the known path

~~(0,0) → (0,1) → (1,1) X (hole) (explore)~~

(restart) (0,0) → (0,1) → (0,2) → (1,2) → (2,2) → (3,2) → (3,3) ✓

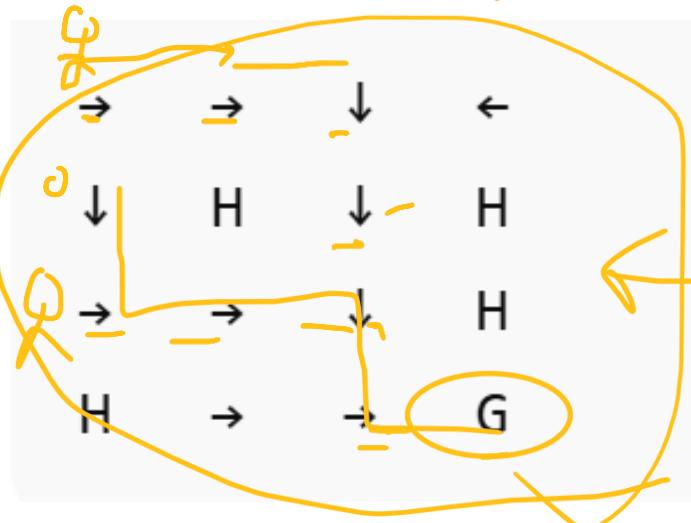
## Updated State-Value Grid (after episode)

0.52	0.66	0.76	0.12
0.48	H	0.80	H
0.40	0.62	0.88	H
H	0.75	0.94	1.00

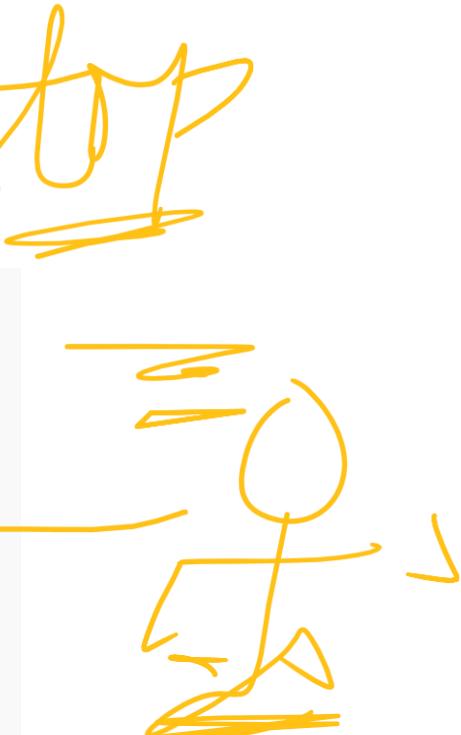
Row \ Col	0	1	2	3
0	[F]	[F]	[F]	[F]
1	[F]	[H]	[H]	[H]
2	[F]	[F]	[F]	[H]
3	[H]	[F]	[F]	[G]



## Policy Grid (after episode)



the successful run reinforces the main path; the hole visit produces no positive reward but teaches the agent to avoid (lowers relative preference).



## Episode 5

0.52	0.66	0.76	0.12
0.48	H	0.80	H
0.40	0.62	0.88	H
H	0.75	0.94	1.00

Row\Col	0	1	2	3
0	[S]	F	F	F
1	F	H	F	H
2	F	F	F	H
3	H	F	F	G

Path Taken

Agent exploits the learned route deterministically  
 $(0,0) \rightarrow (0,1) \rightarrow (0,2) \rightarrow (1,2) \rightarrow (2,2) \rightarrow (3,2) \rightarrow (3,3)$  

Updated State-Value Grid (after episode)

0.60	0.70	0.80	0.13
0.56	H	0.83	H
0.50	0.68	0.90	H
H	0.80	0.96	1.00

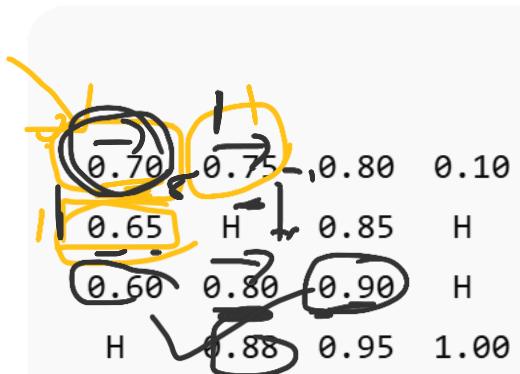
Policy Grid (after episode)

→	→	↓	←
↓	H	↓	H
→	→	↓	H
H	→	→	G

(Policy is now effectively deterministic along the optimal corridor.)

Here's how the **arrows (policy)** are derived from the **state values**  $V(s)$  —the agent looks at the neighboring states' values and points toward the **highest one**.

State-Value Grid ( $V$ )

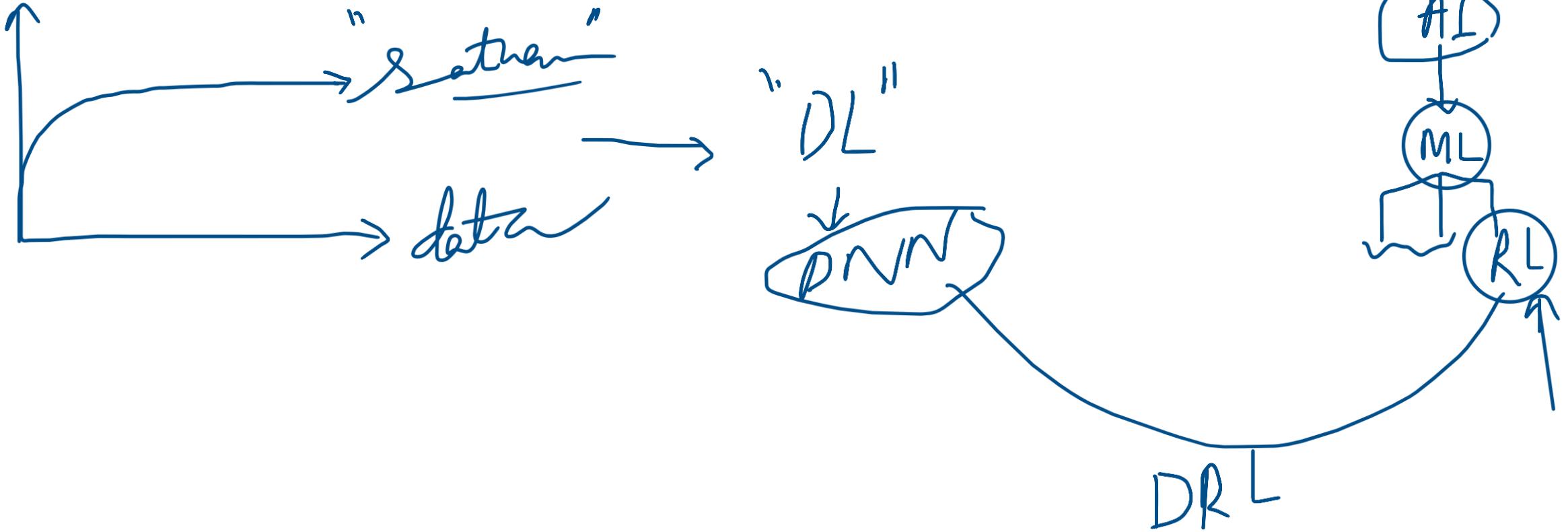


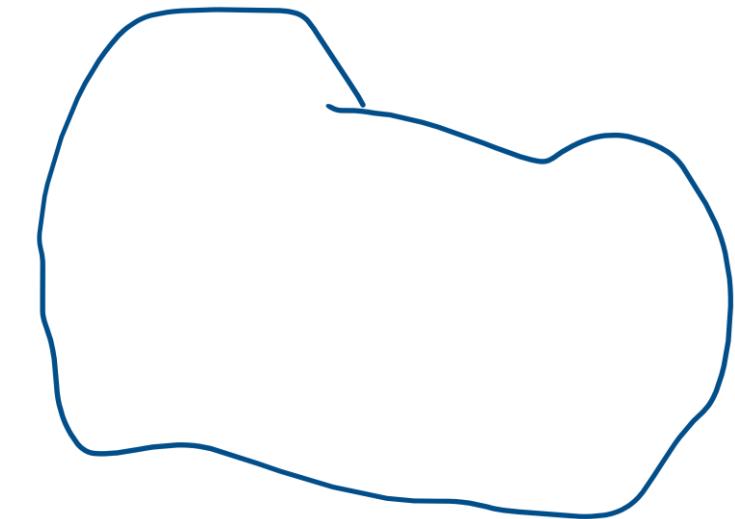
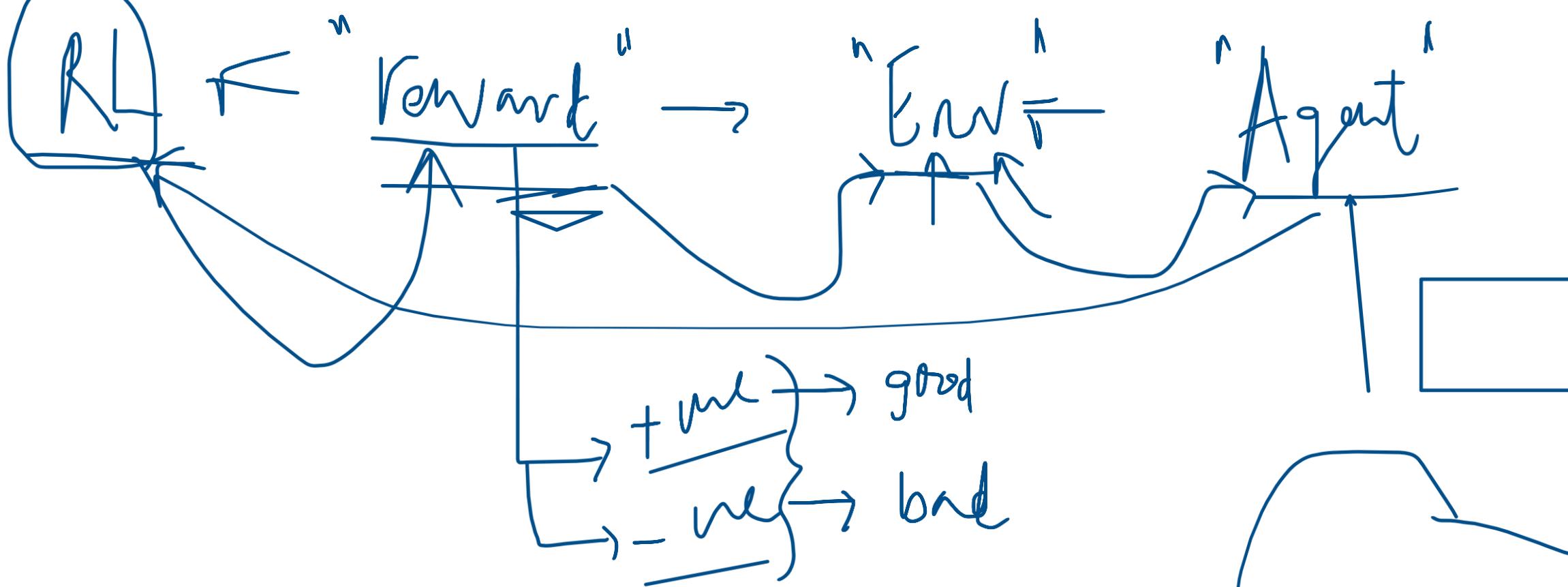
Derived Policy Grid

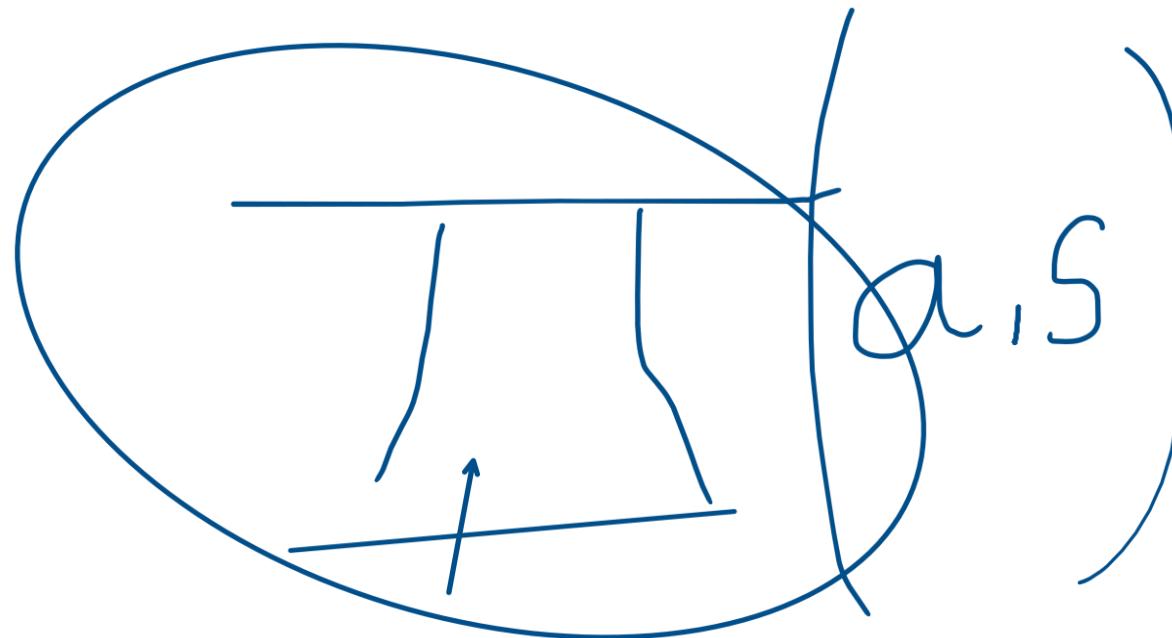


**Interpretation:**

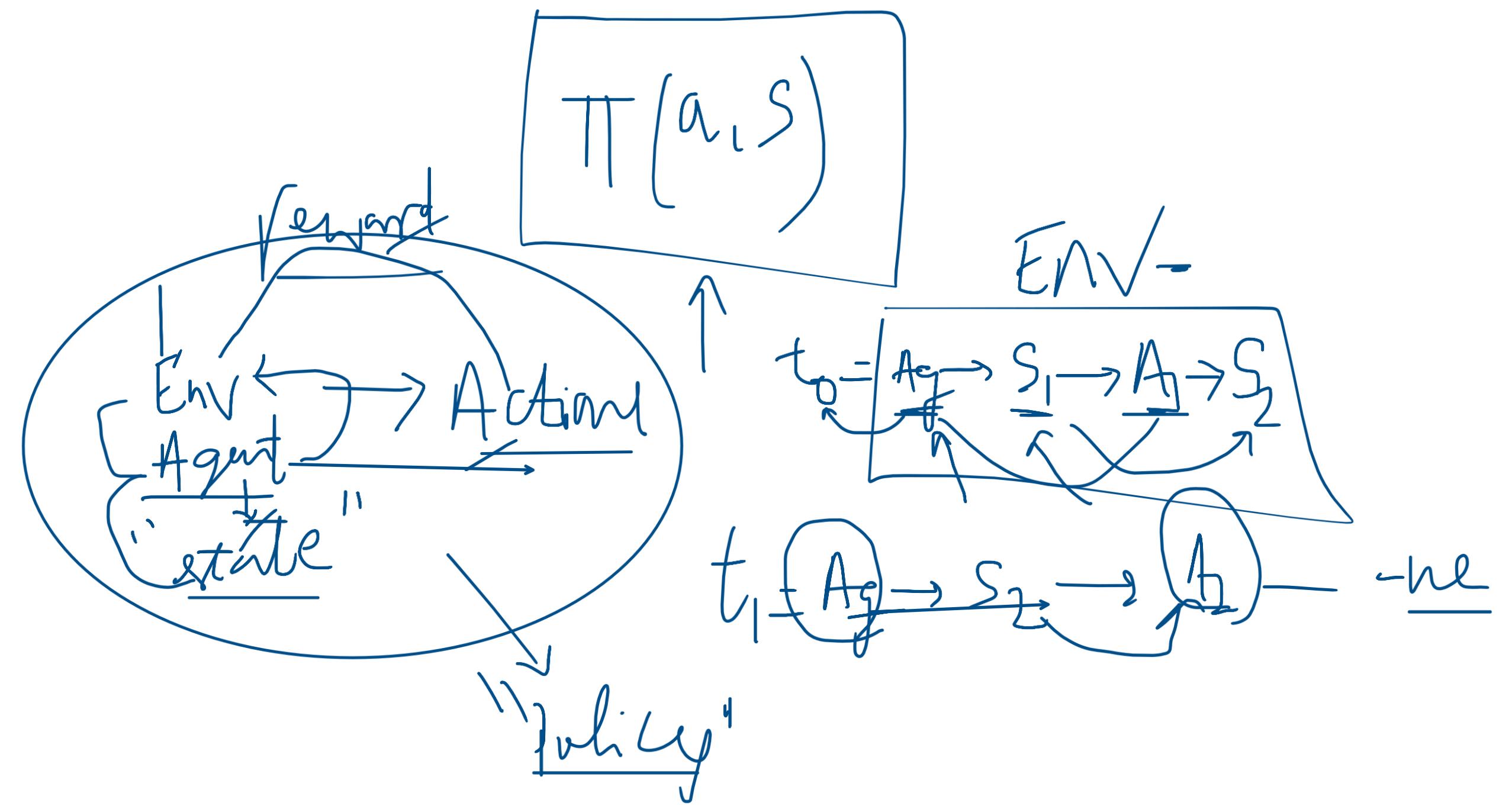
- Each arrow points toward the neighbor with the **highest next-state value**.
- For example:
  - At (0,0) value=0.70 → best neighbor is (0,1)=0.75 → arrow →
  - At (2,2) value=0.90 → best neighbor is (3,2)=0.95 → arrow ↓
- Holes (H) and Goal (G) are terminal — no arrows.



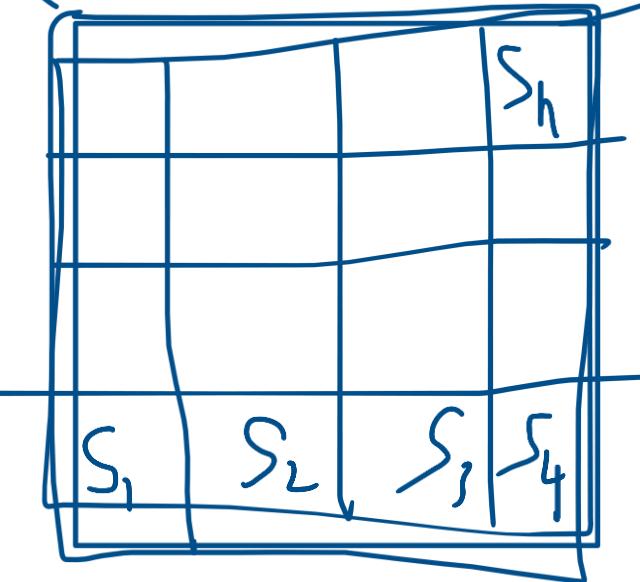




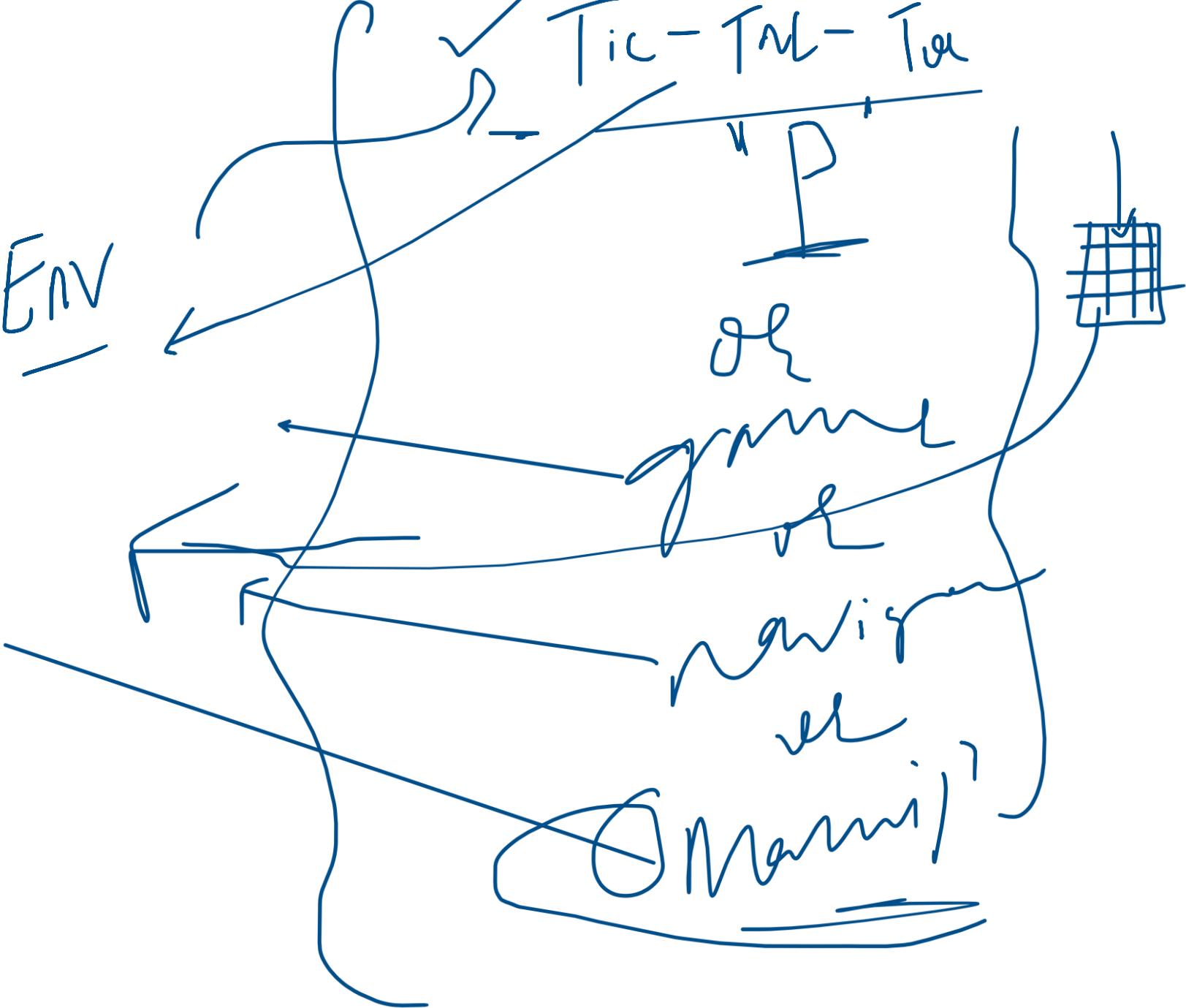
"Prob"  
"π"  
"action"  
"state"  
+ve reward

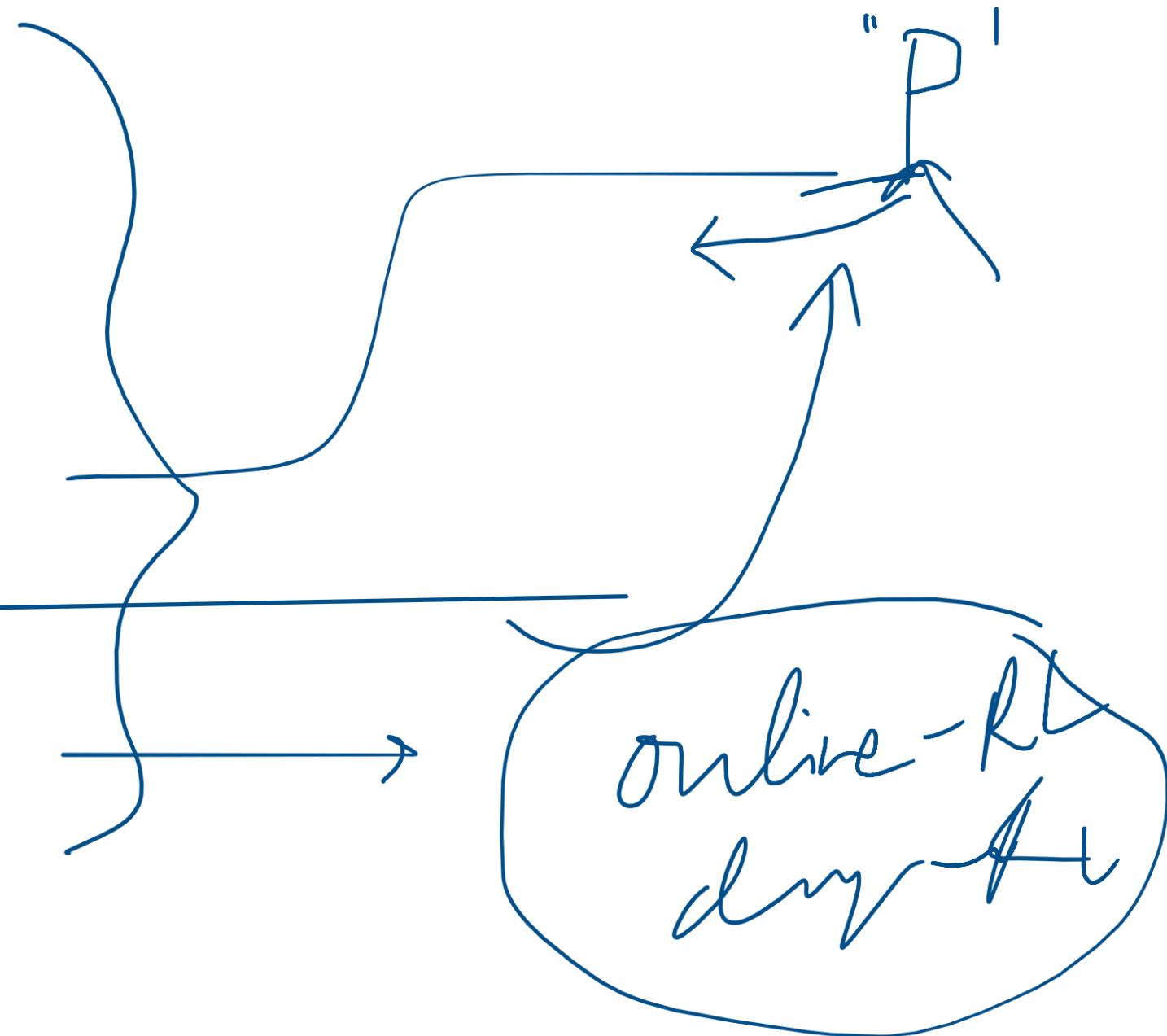
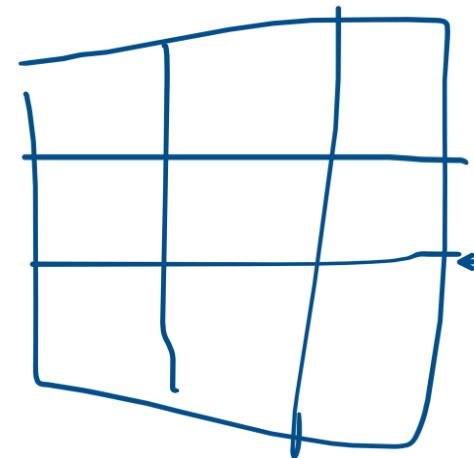
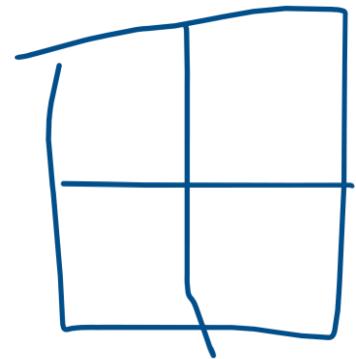


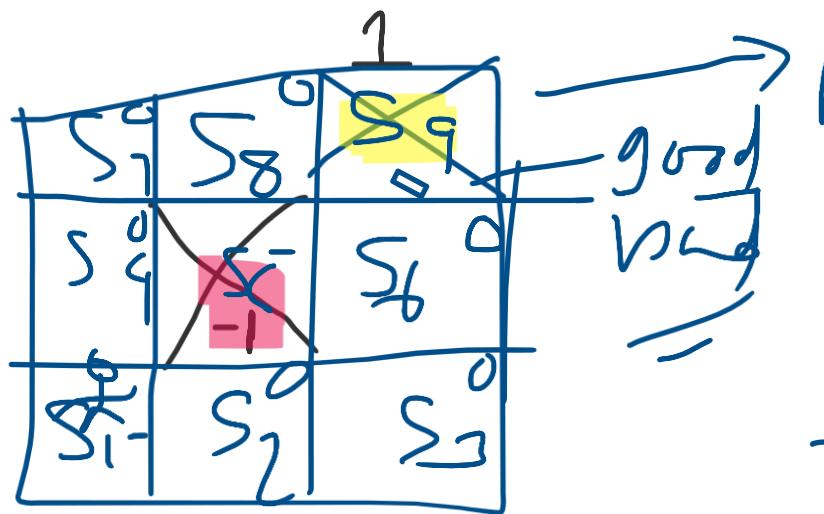
Grid-world



Env







$\xrightarrow{P_t}$

ENV  
Agent  
States { $S_1 \rightarrow S_9$ }

$3 \times 3 = n \times n$

$S_1 \rightarrow R, U, \times, D, \square$   
 $S_2 \rightarrow R, L, U$

Action

$S_5 \rightarrow R, L, U, D$

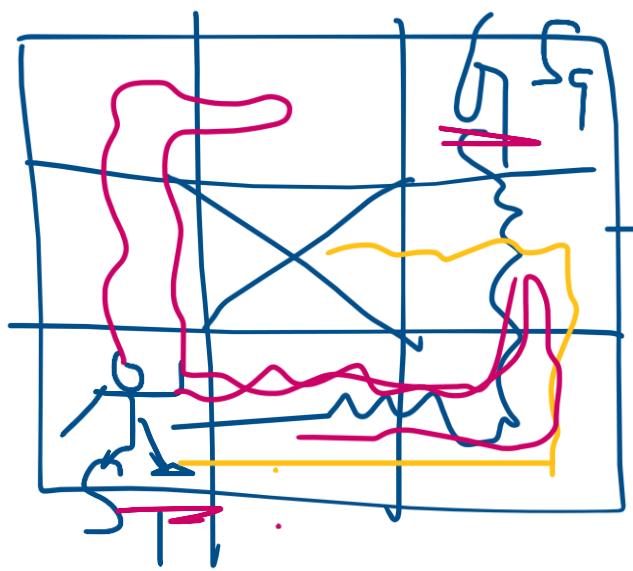
Reward

+wl  $\rightarrow 1$

-wl  $\rightarrow -1$

normal  $\rightarrow 0$





"Initial" ~~Policy~~  
 "Explore" ~~Policy~~

1) G ✓

2) Trap | Bad → X ✓

3) X ↗ "Contract"

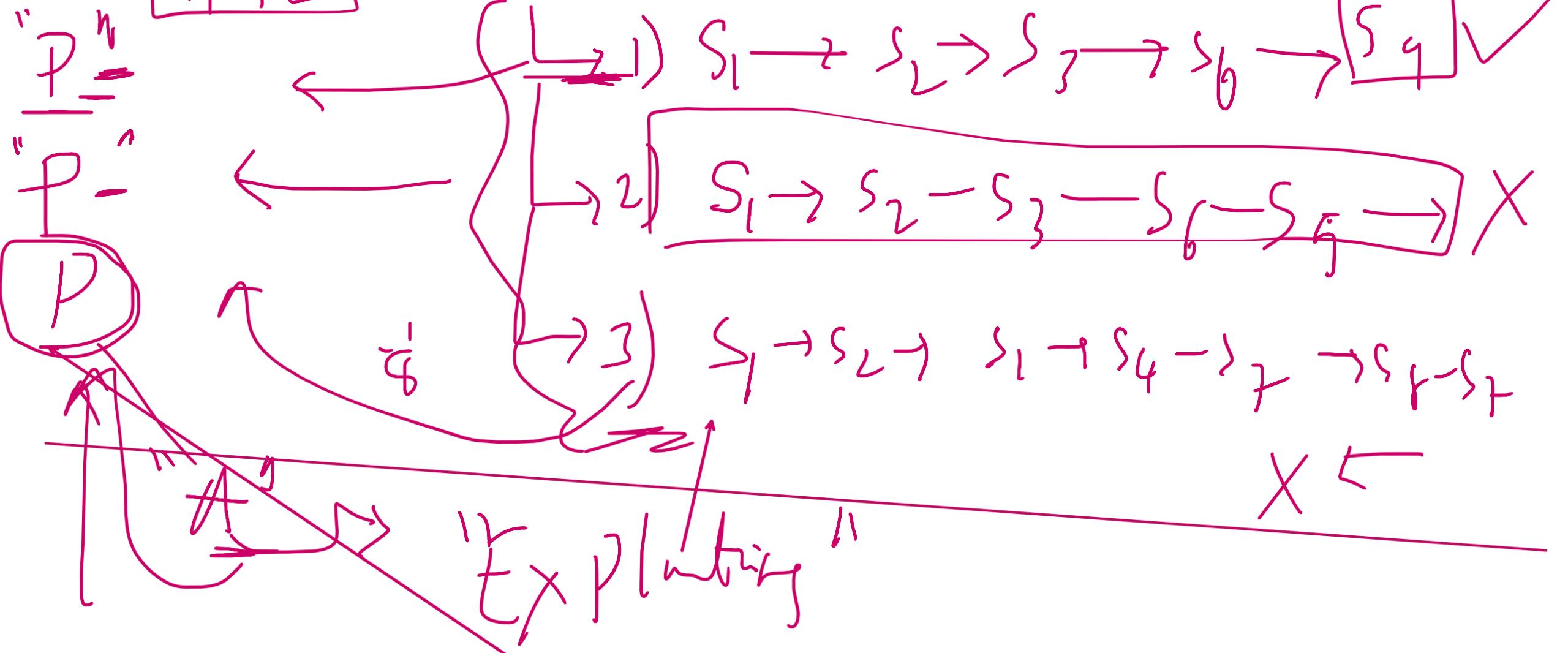
Timer  
~~t~~  
 Stop

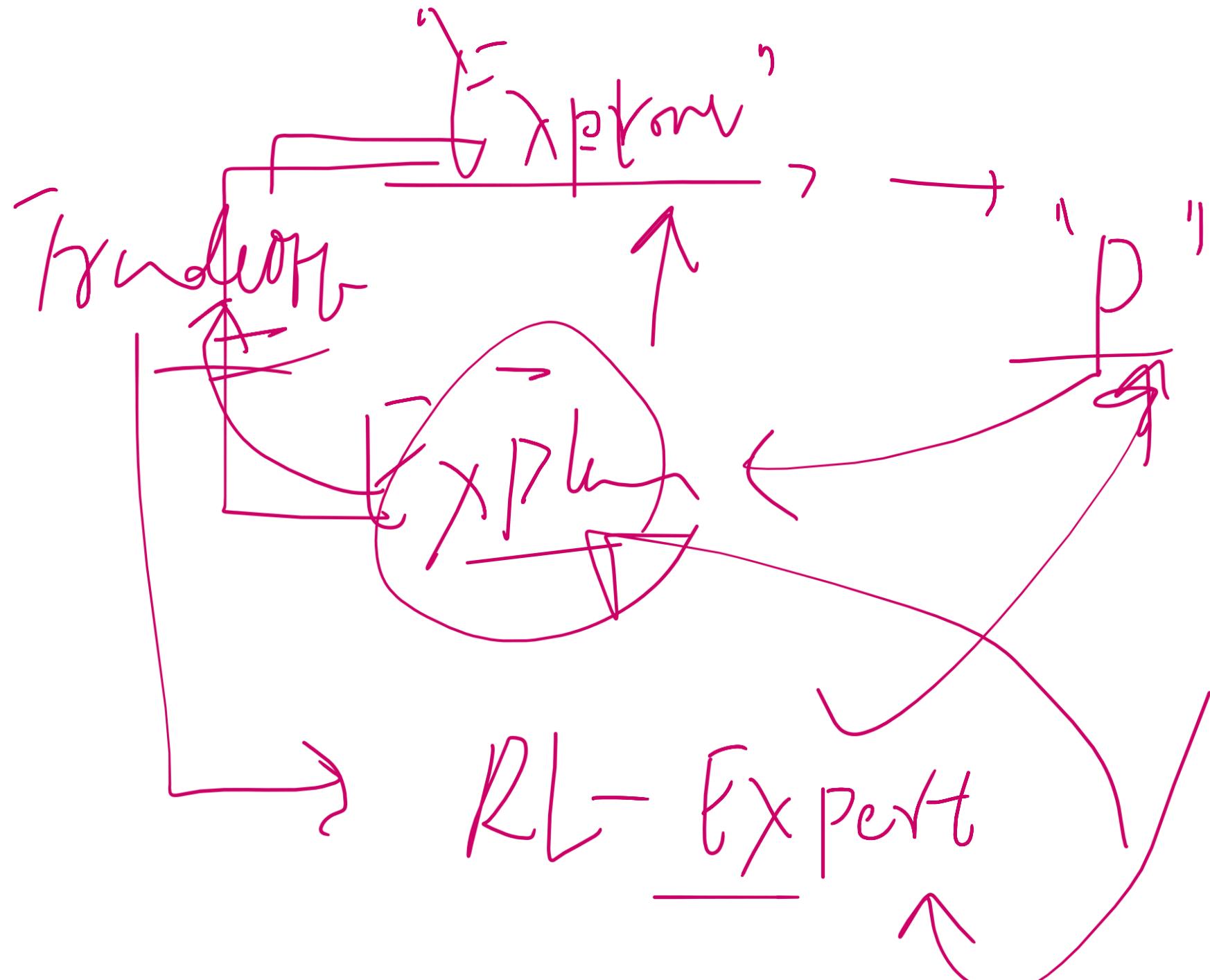


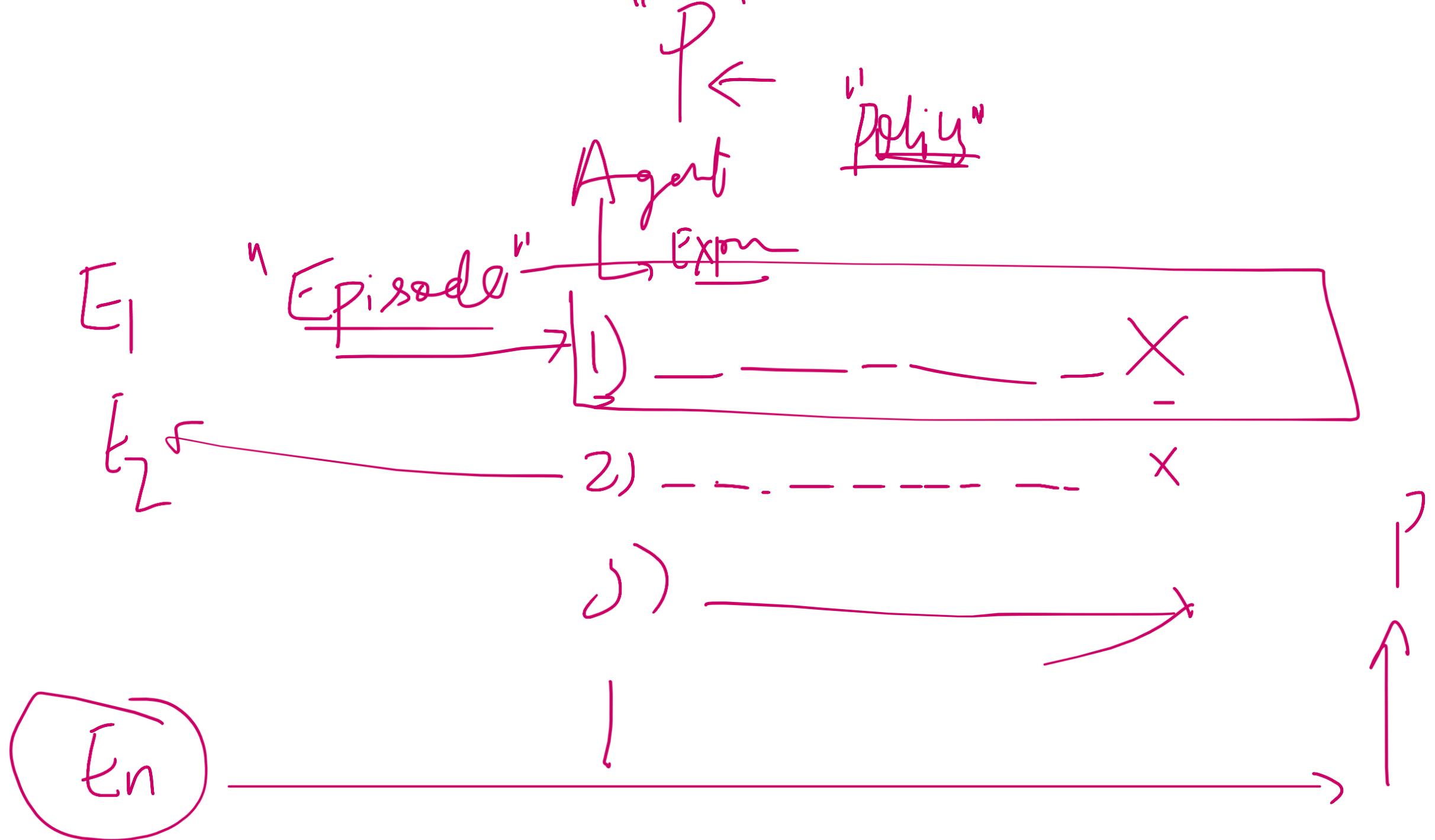
$s_3$	$s_4$	$s_5$	$s_6$	$s_7$
$s_4$	$s_5$	$s_6$		
$s_1$	$s_2$	$s_3$		

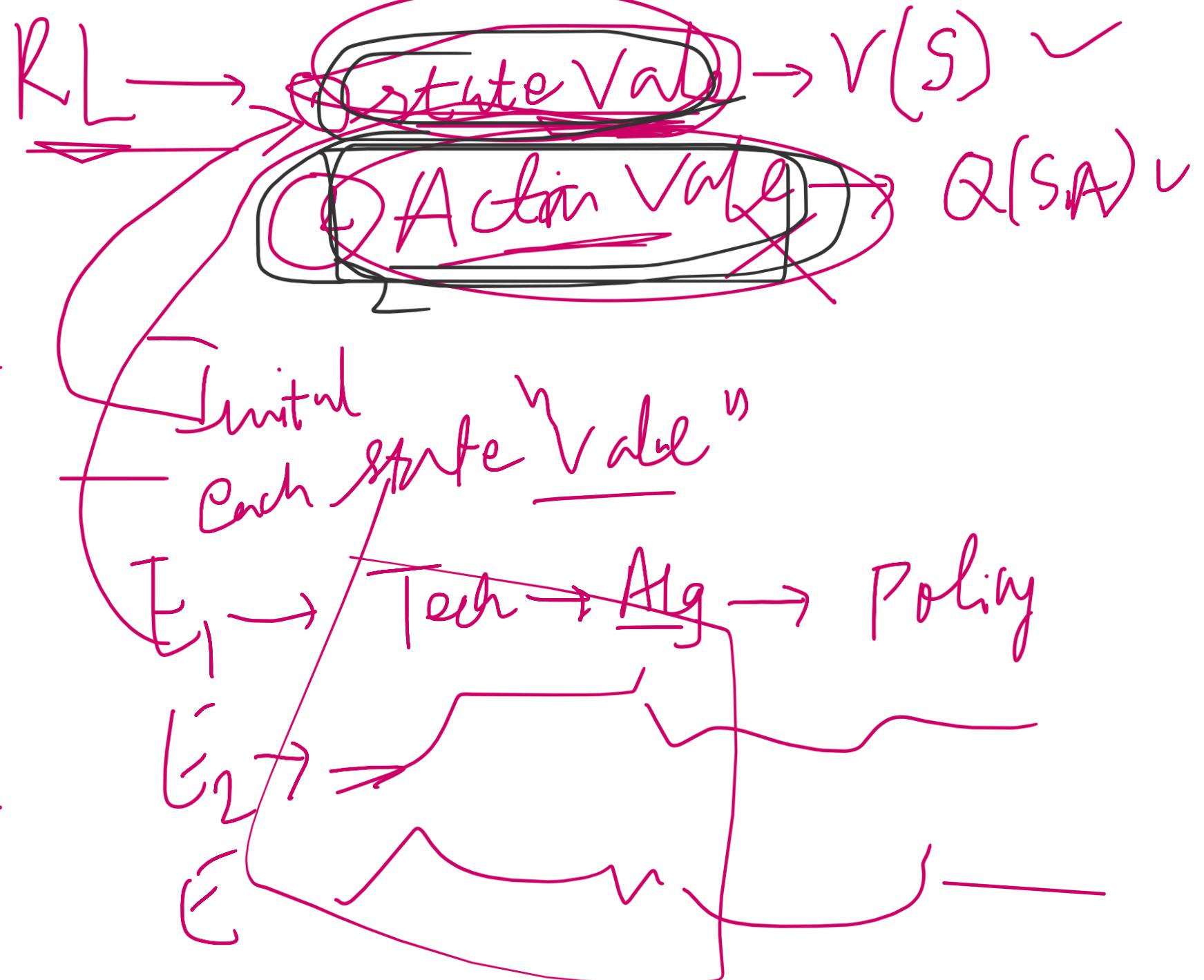
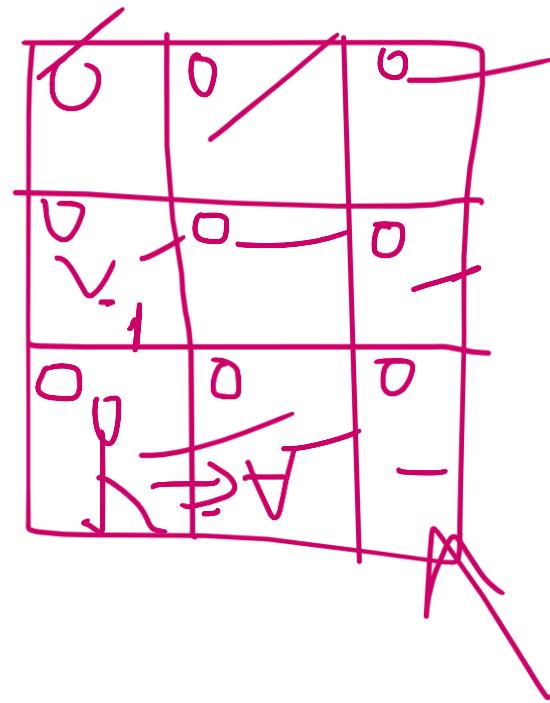
Initial  $\pi \rightarrow$  unknown policy

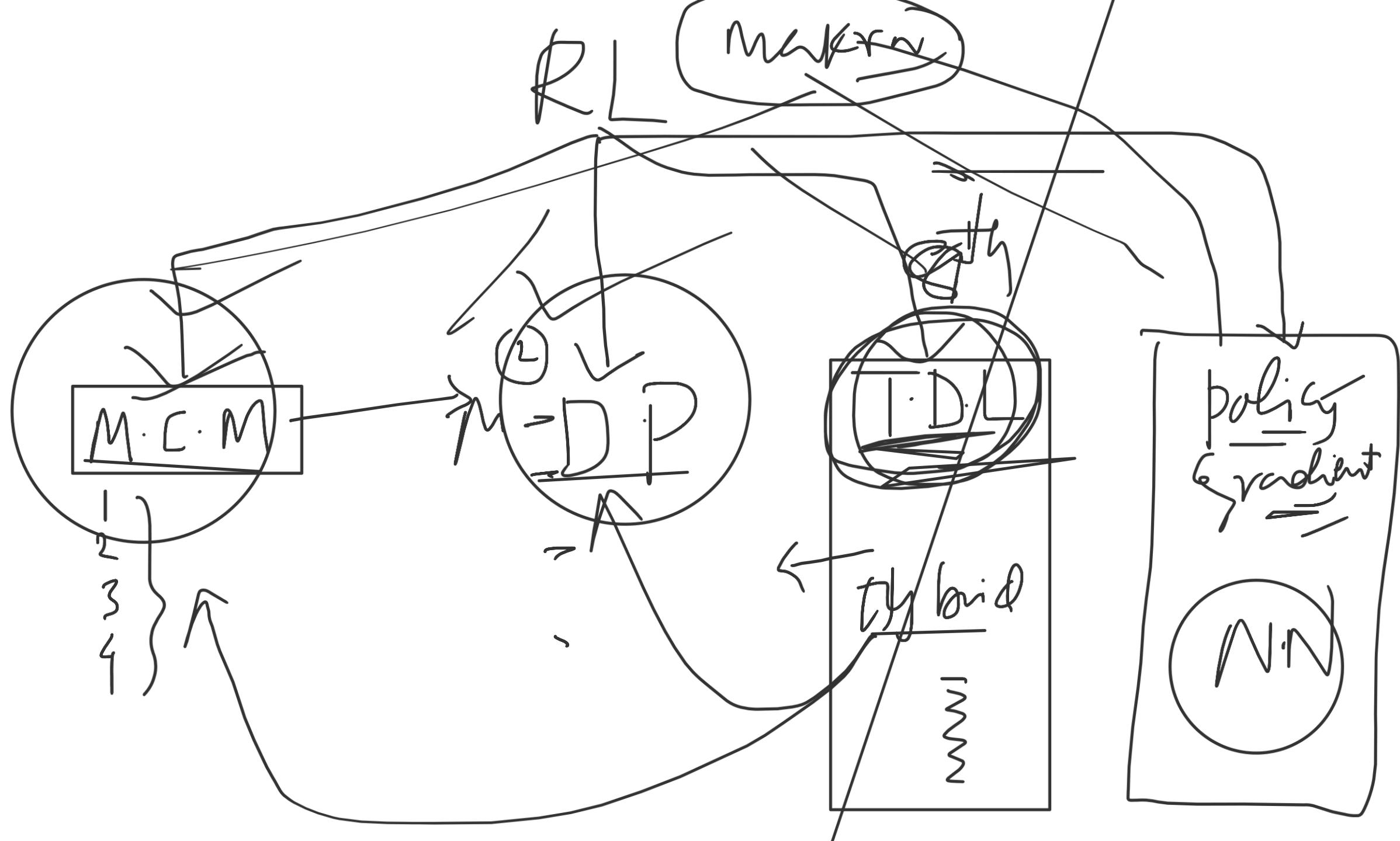
Explore

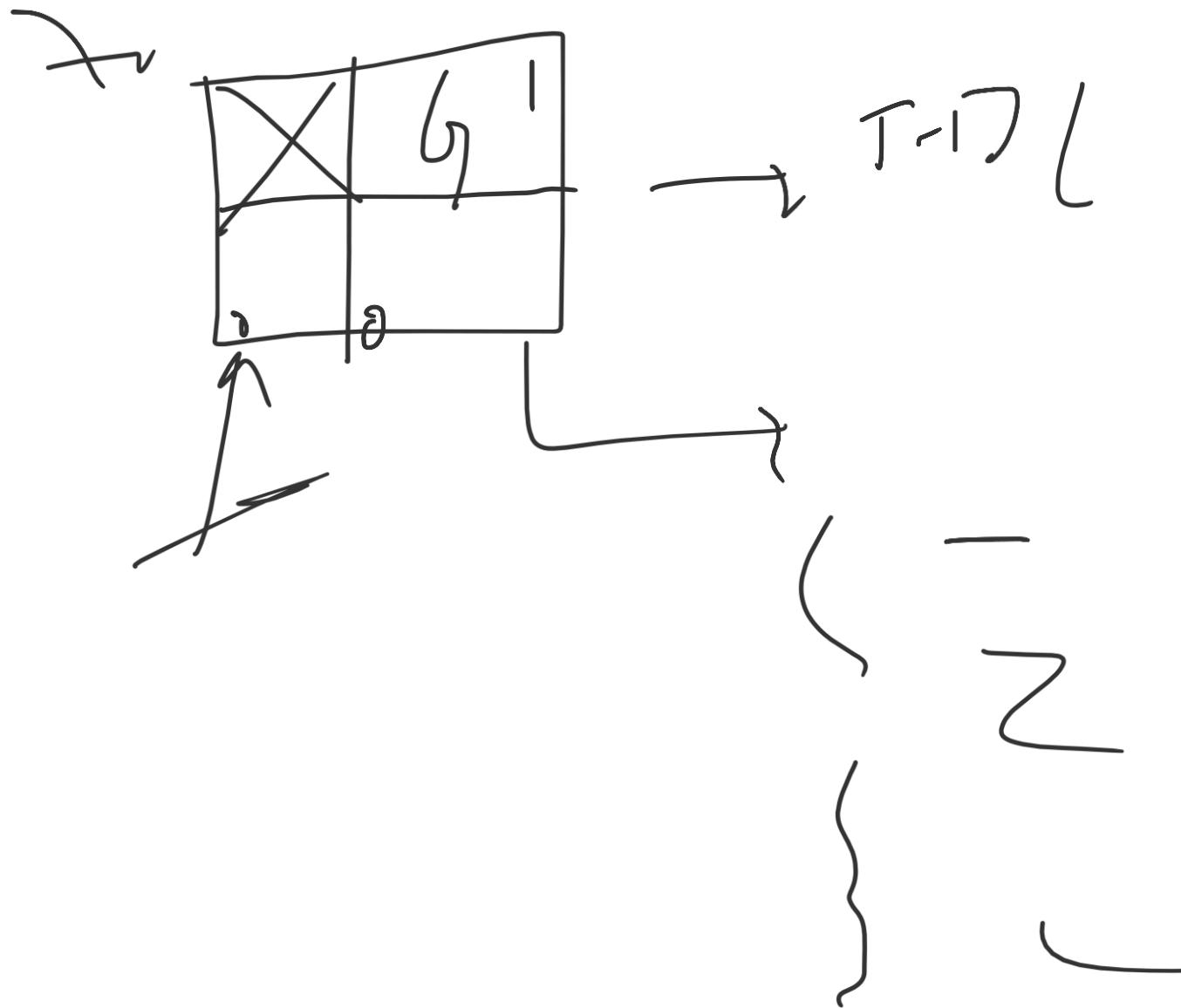


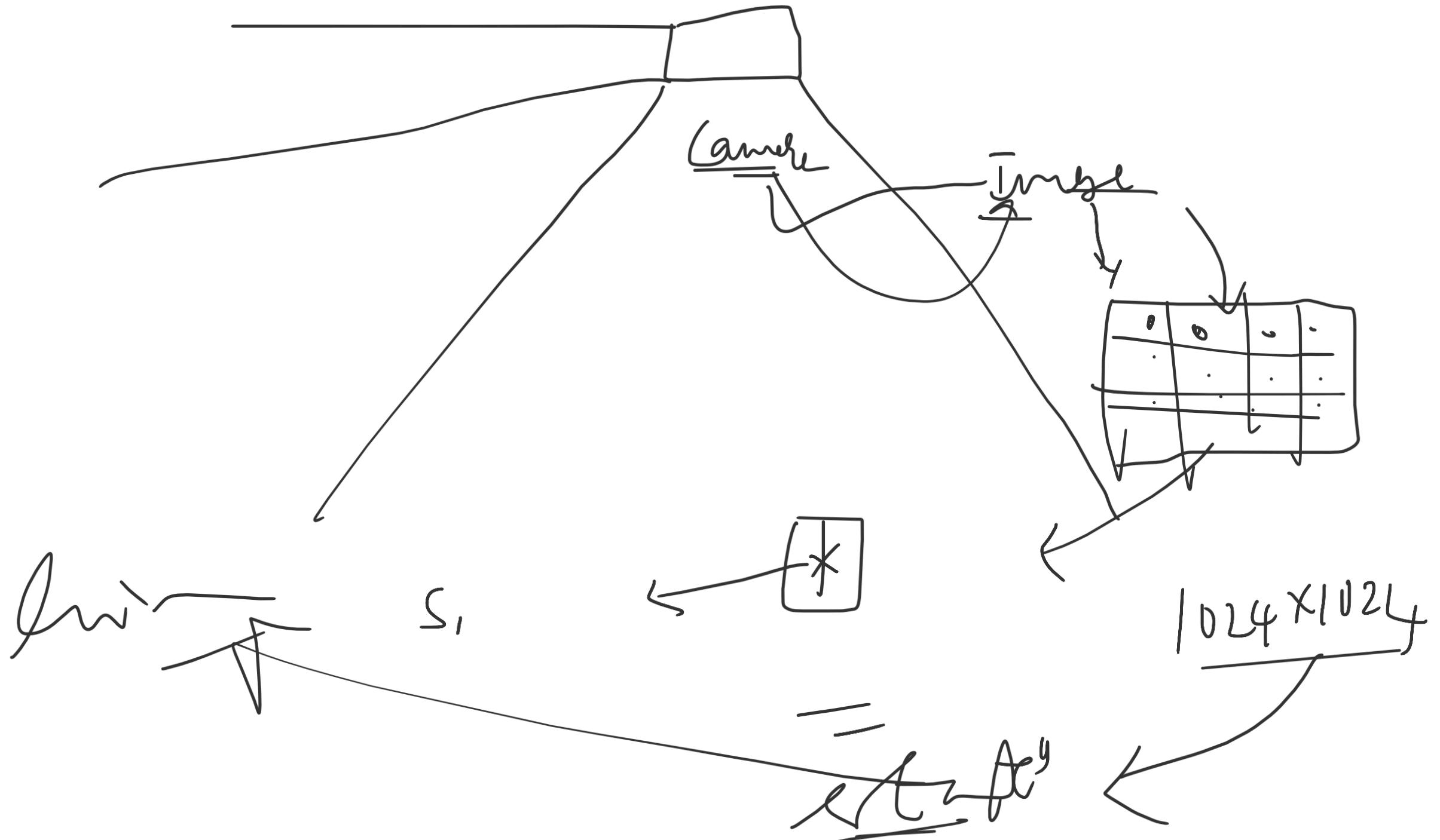




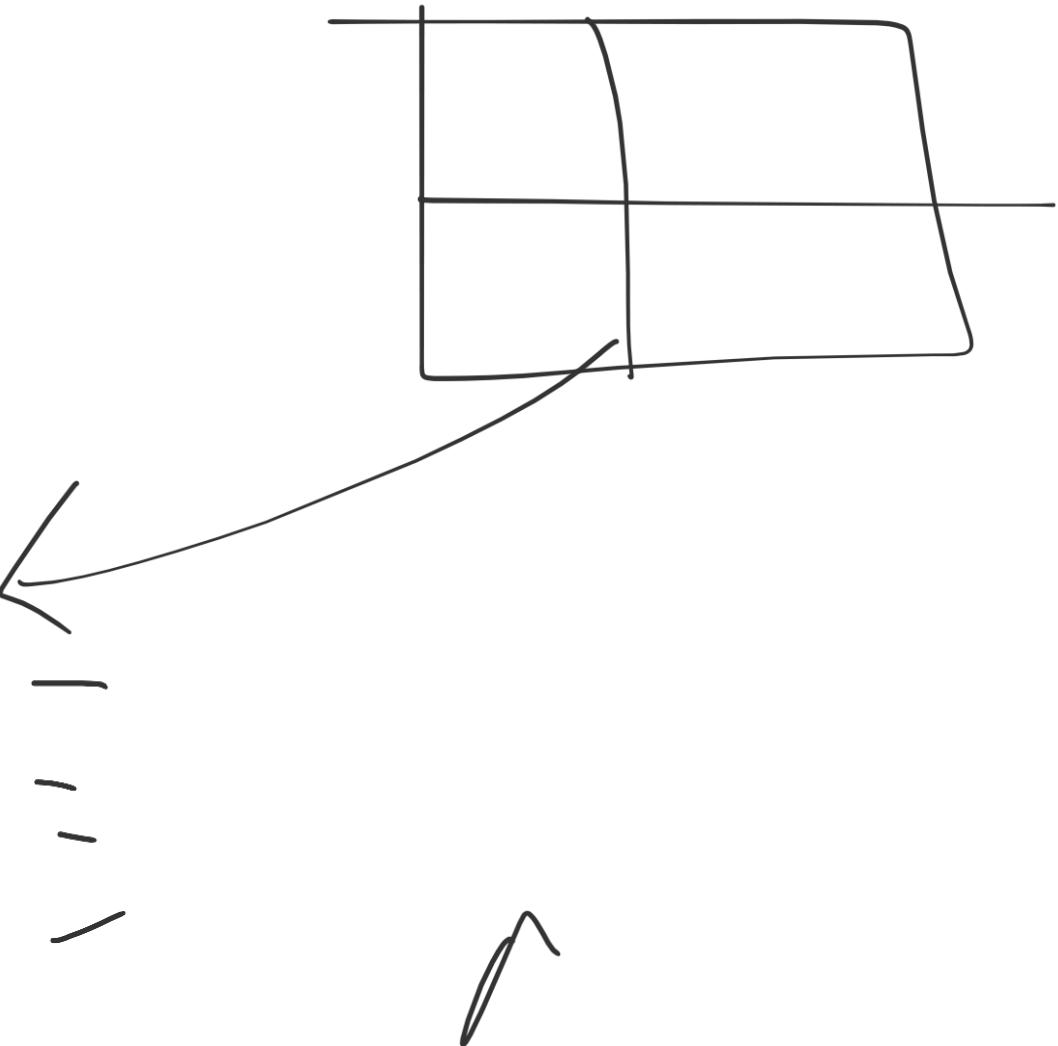












$$V(s) + \min_{\theta} \left( \begin{array}{l} s - \theta \\ | - 0 \end{array} \right)$$

