# REINFORCEMENT LEARNING

JALEY DHOLAKIYA

# RECAP : DEEP LEARNING
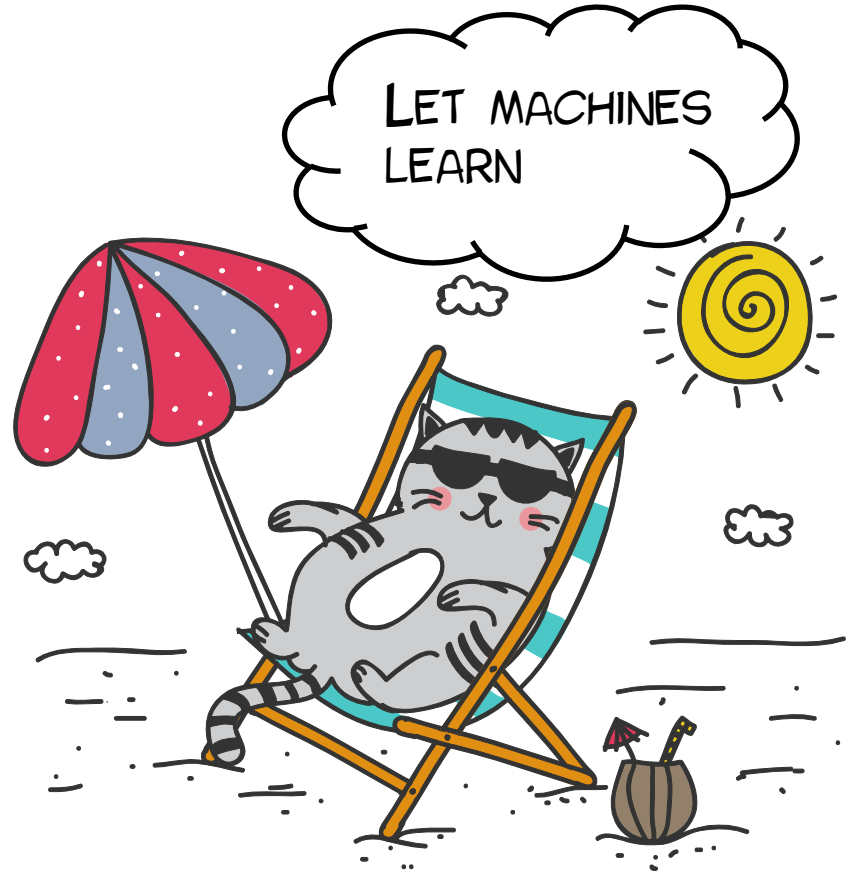


INPUT

MODEL

BISCUITE
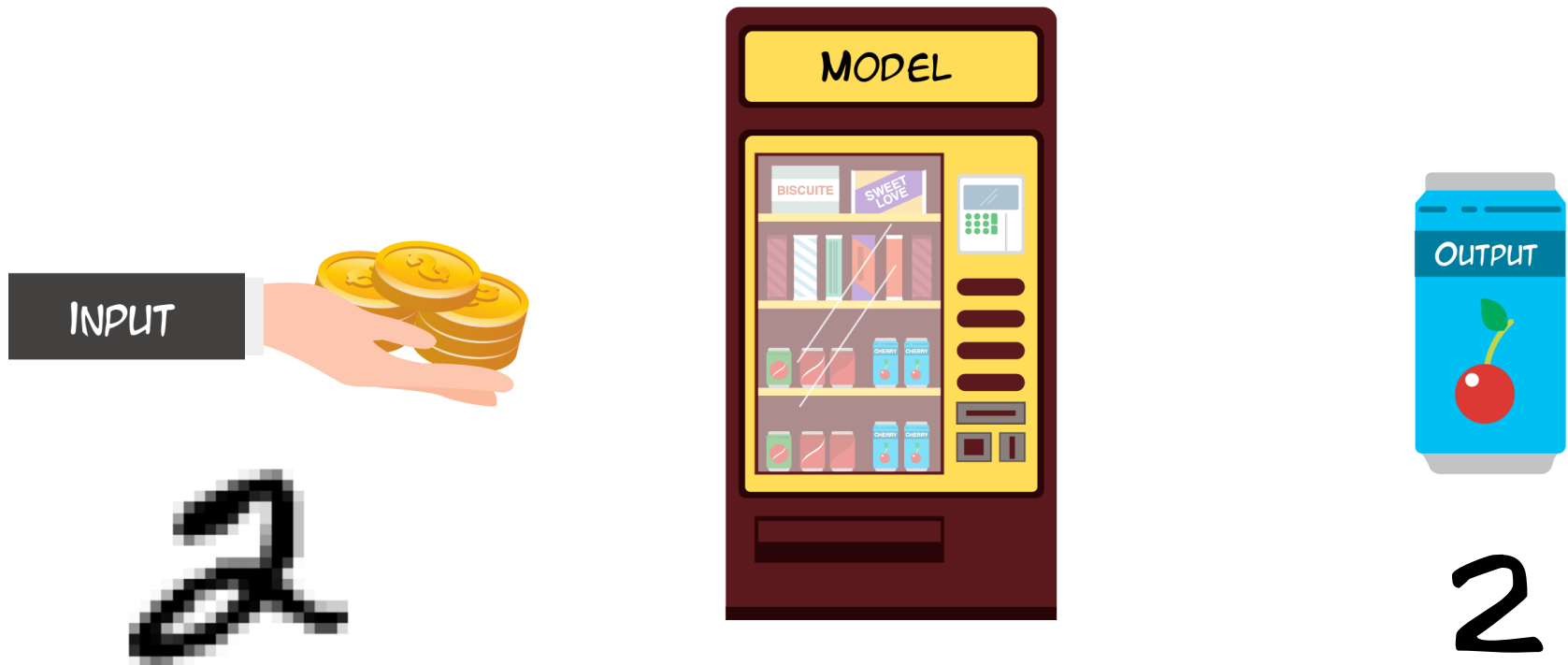SWEET LOVE
CHERRY CHERRY
CHERRY CHERRY

OUTPUT

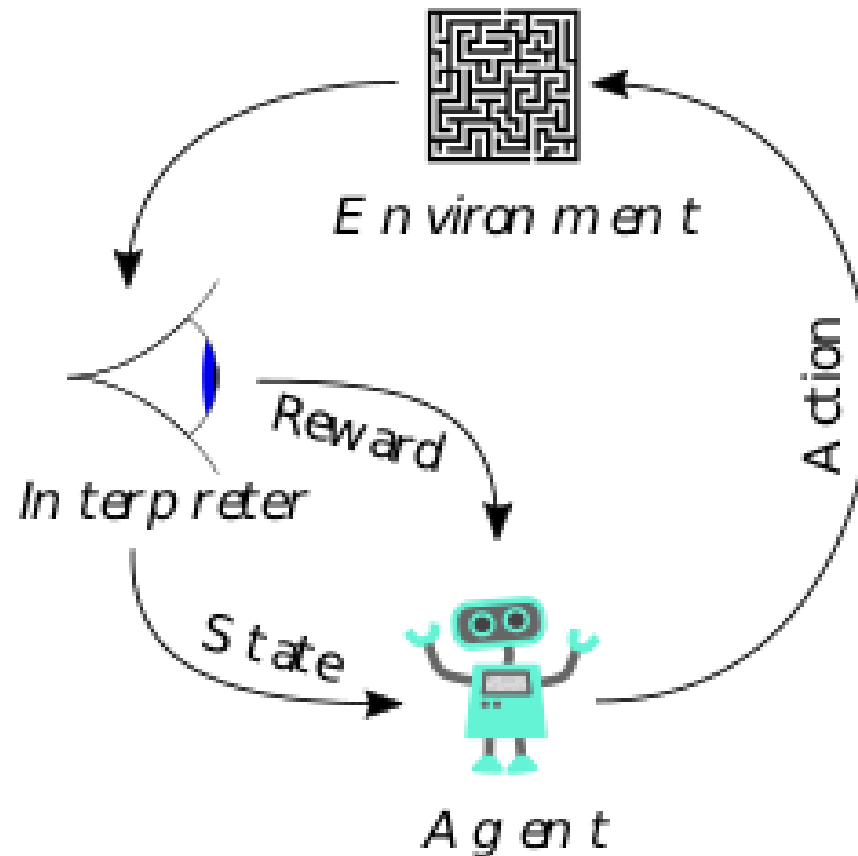$$y \approx f_w(x)$$

# MNIST Digit Classification



$$y \approx f_w(x)$$

# Challenges in Deep learning

- Its Data Hungry
- No Active participation

# REINFORCEMENT LEARNING

# RL vs DL

- Active Engagement
- Get rid of Manual Data Collection

# Components of RL



AGENT        POLICY = ACTION|STATE

STATE    OBSERVATION    ACTION    REWARD

# ENVIRONMENT
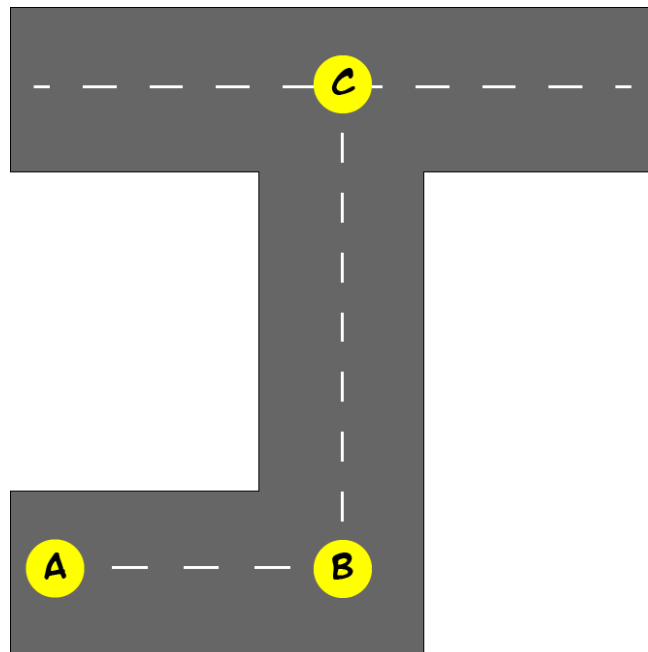
- EPISODIC VS NON-EPISODIC
- FIXED VS DYNAMIC
- SINGLE VS MULTIPLAYER

# VIDEO 1 : COMPONENTS OF RL

# Gym Example

# GOAL OF RL

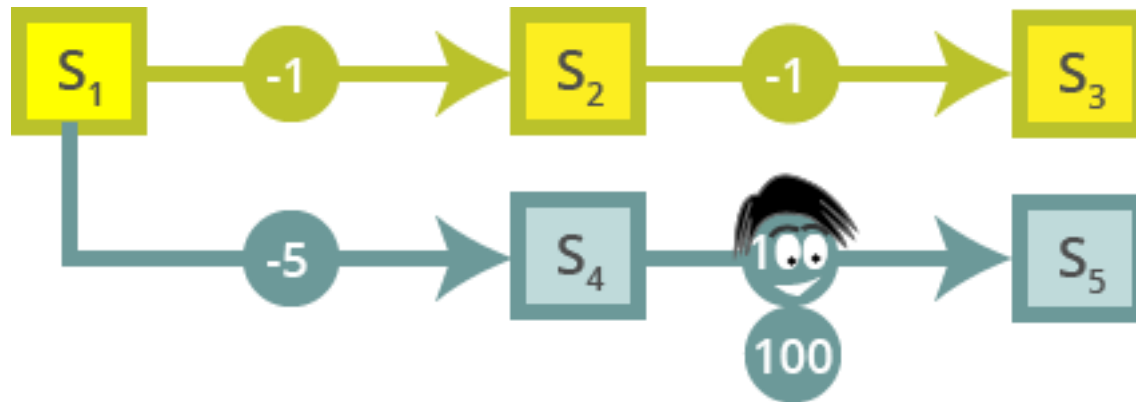| STATE | ACTION |
|-------|--------|
| A | GO STRAIGHT |
| B | TAKE LEFT |
| C | TAKE RIGHT OR LEFT |

# WHICH ACTION IS BEST?

MOST REWARDING ACTION

# IMMEDIATE REWARD?

# IMMEDIATE REWARD?

# DISCOUNTED REWARD

# TYPES OF RL ALGORITHMS
## (BASED ON MODEL)

- MODEL DEPENDENT
    (TOO MUCH MATH/ LESS USEFUL)

- MODEL FREE
    (LESS MATH + MORE FUN)

# TYPES OF RL ALGORITHMS
## (BASED ON POLICY)
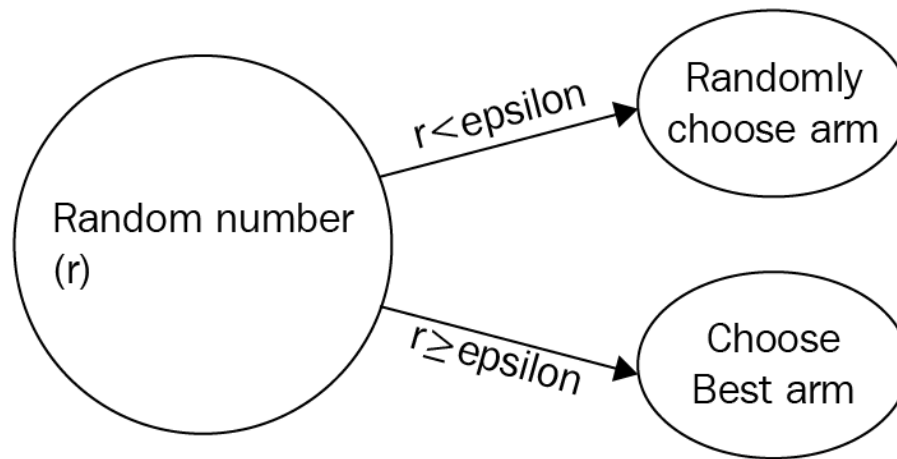
- POLICY EVALUATION

  ESTIMATE THE STATE/ACTION REWARDS
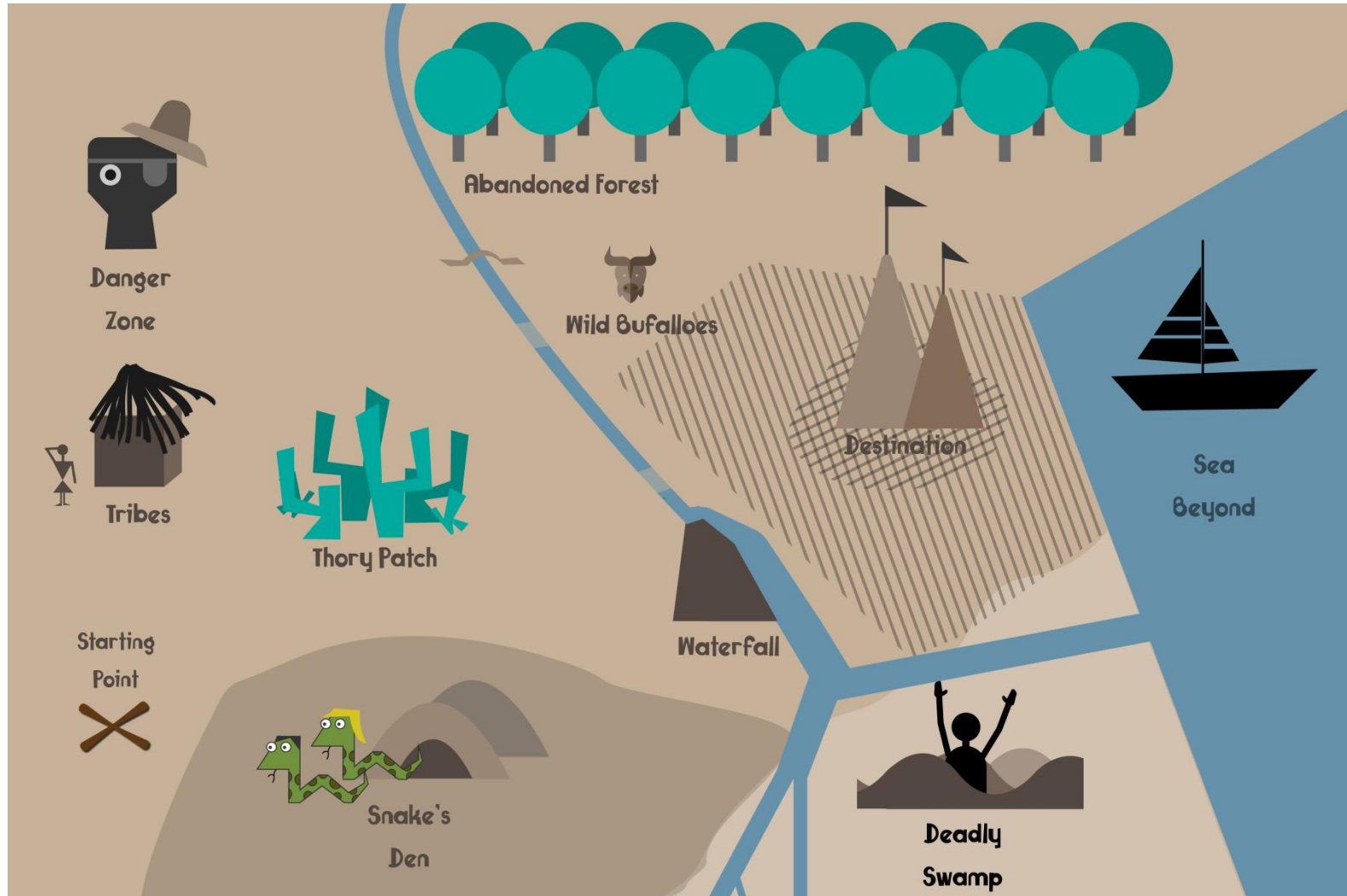  FOR A GIVEN POLICY

- POLICY CONTROL

  TASK IS TO **IMPROVE** A POLICY

# POLICY CONTROL
## EPSILON GREEDY POLICY

# EXPLORE EXPLOIT PROBLEM

# EXPLORE EXPLOIT PROBLEM



STATE VISITATION FREQUENCY

$\pi^{old}(a\,|\,s)$

$\pi^{new}(a\,|\,s)$

# Monte Carlo Evaluation

Farm → Sea → Island → Forest → Dead

Sea → Island → Forest → Mountain → Forest → Dead

Forest → Sea → Island → Sea → Island → Forest → Dead

N = 0   N = 0   N = 0   N = 0   N = 0

# MONTE CARLO EVALUATION

**Algorithm 1:** Monte Carlo Evaluation

1   $N(s) \leftarrow 0, \forall s \in S$
2   $S(s) \leftarrow 0, \forall s \in S$
3   **for** *each episode* **do**
4     **for** *each state(s) $\in$ episode* **do**
5      $G_t \leftarrow 0$
6      $i = 0$
7      **for** $t \leftarrow [t_{curr}, T]$ **do**
8       $G_t \leftarrow G_t + \gamma^i R_t$
9       $i \leftarrow i + 1$
10      $N(s) \leftarrow N(s) + 1$
11      $S(s) \leftarrow S(s) + G_t$
12      $V(s) \leftarrow S(s)/N(s)$

**Algorithm 2:** Monte Carlo Evaluation with incremental updates

1   $N(s) \leftarrow 0, \forall s \in S$
2   $S(s) \leftarrow 0, \forall s \in S$
3   **for** *each episode* **do**
4     **for** *each state(s) $\in$ episode* **do**
5      $G_t \leftarrow 0$
6      $i = 0$
7      **for** $t \leftarrow [t_{curr}, T]$ **do**
8       $G_t \leftarrow G_t + \gamma^i R_t$
9       $i \leftarrow i + 1$
10      $V(s) \leftarrow V(s) + \alpha(G_t - V(s))$

     *or*

11      *or*
12      $N(s) \leftarrow N(s) + 1$
13      $V(s) \leftarrow$
      $V(s) + \frac{1}{N(s)}(G_t - V(s))$

# TEMPORAL DIFFERENCE

$$G_t(s_t) = R_{t+1} + \gamma R_{t+2} + \gamma^{T-2} \cdots R_{T-1}$$

We can write reward for next state as

$$G_{t+1}(s_{t+1}) = R_{t+2} + \gamma R_{t+3} + \gamma^{T-3} \cdots R_{T-1}$$

Therefore we can also write the first equation as

$$G_t(s_t) = R_{t+1} + \gamma \overbrace{\{R_{t+2} + \gamma^{T-3} \cdots R_{T-1}\}}^{G_{t+1}(s_{t+1})}$$
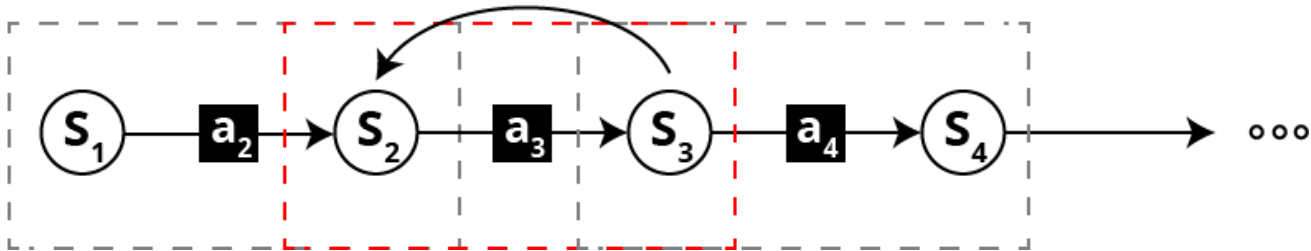
and we all know we want to update value function, closer to discounted reward $(G \leftarrow V)$, therefore our new update for TD Learning becomes

$$V(s_t) = V(s) + \alpha(\hat{G}_t - V(s))$$

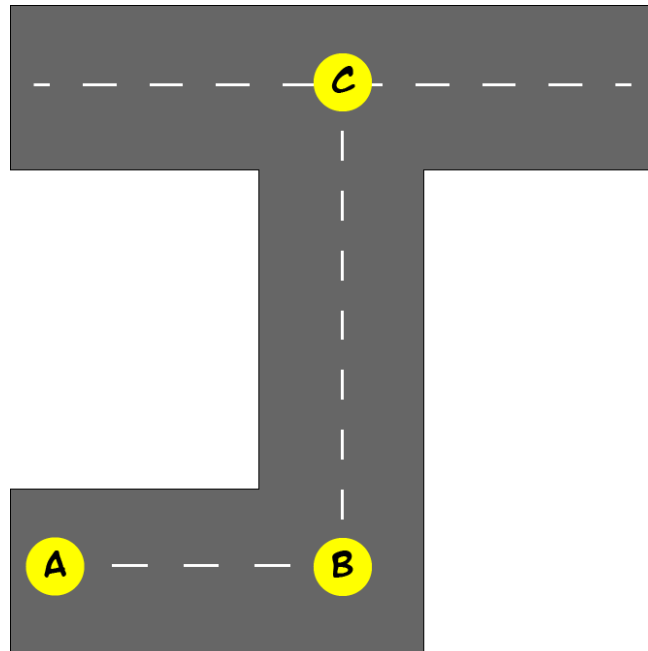$$V(s_t) = V(s_t) + \alpha(R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

# TEMPORAL DIFFERENCE

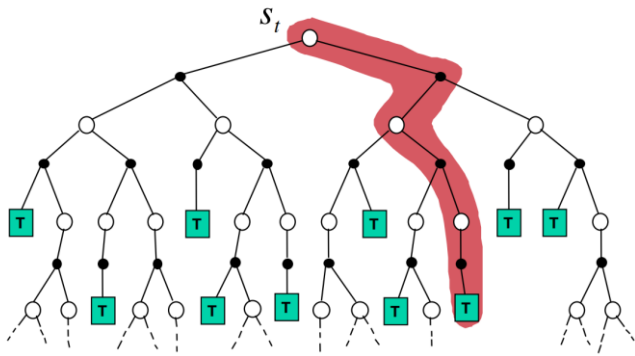UPDATE $V(s_2)$ BASED ON $R_3$ and $V(s_3)$

# Reminder : Goal of RL

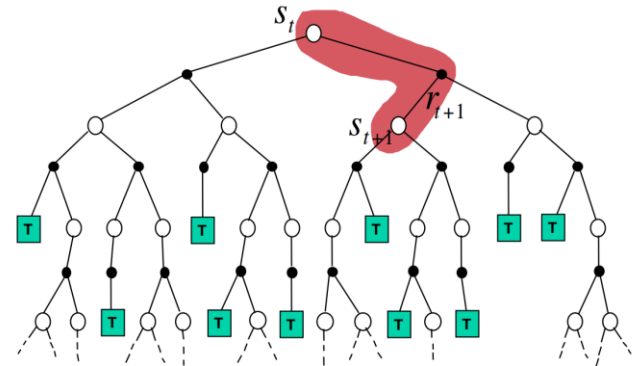| State | Action |
|-------|--------|
| A | Go Straight |
| B | Take Left |
| C | Take Right or Left |

# MC VS TD

## MONTE CARLO

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



$$q(a_t|s_t) = q(a_t|s_t) + a(G_t - q(a_t|s_t))$$

## TEMPORAL DIFFERENCE

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



$$q(a_t|s_t) = q(a_t|s_t) + \alpha((R_t + \gamma q(a_{t+1}|s_{t+1})) - q(a_t|s_t))$$

# TD CONTROL + GREEDY = SARSA
(STATE ACTION REWARD STATE ACTION)

## POLICY EVALUATION

$$q(a_t|s_t) = q(a_t|s_t) + \alpha((\textcolor{red}{R_t + \gamma q(a_{t+1}|s_{t+1})}) - q(a_t|s_t))$$

$$q(a_t|s_t) = \textcolor{red}{(\mathbf{1 - \alpha})}q(a_t|s_t) + \textcolor{red}{\boldsymbol{\alpha}}(R_t + \gamma q(a_{t+1}|s_{t+1}))$$

## POLICY UPDATE

$\epsilon - greedy\ approach$

# SARSA

$$q(a_t|s_t) = \mathbf{(1 - \alpha)}q(a_t|s_t) + \mathbf{\alpha}(R_t + \gamma q(a_{t+1}|s_{t+1}))$$

# Q LEARNING

REQUIRES ONLY SARS PAIRS

$$q(a_t|s_t) = \mathbf{(1 - \alpha)}q(a_t|s_t) + \mathbf{\alpha} * max(R_t + \gamma q(a_*|s_{t+1}))$$

# DEMO TIME

## (GRID-WORLD)

# GOAL OF Q LEARNING

- ESTIMATING ACTION VALUE FUNCTION (Q)

- UPDATE Q FUNCTION VIA MOVING AVG

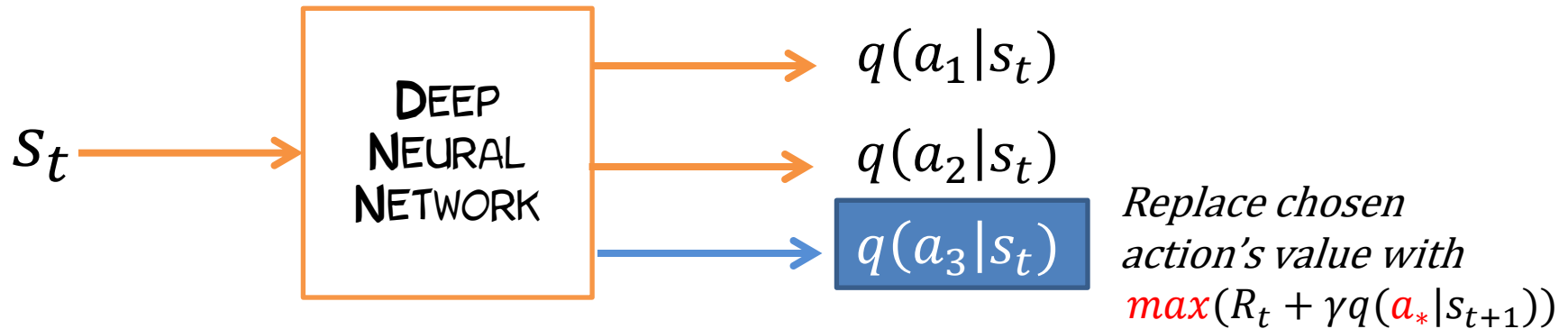$$q(a_t|s_t) = (1 - \alpha)q(a_t|s_t) + \alpha * max(R_t + \gamma q(a_*|s_{t+1}))$$

# Problem with Q Learning

- Can't handle near Infinite states (images)
- Can't handle State as feature vector

# Q LEARNING

$$q(a_t|s_t) = (\mathbf{1 - \alpha})q(a_t|s_t) + \mathbf{\alpha} * max(R_t + \gamma q(a_*|s_{t+1}))$$

# DEEP Q LEARNING

$s_t$ → DEEP NEURAL NETWORK →

$q(a_1|s_t)$

$q(a_2|s_t)$

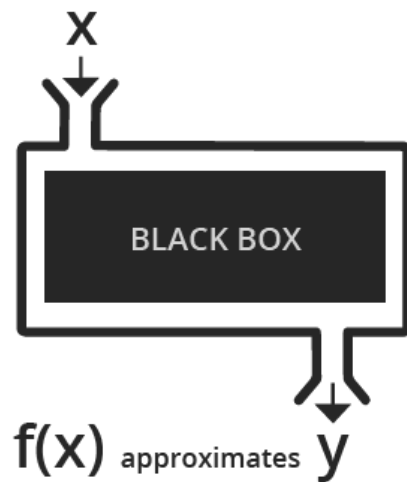$q(a_3|s_t)$

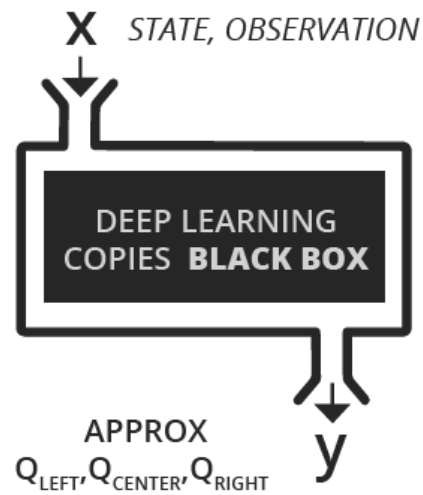*Replace chosen action's value with* $max(R_t + \gamma q(a_*|s_{t+1}))$

Update neural network by setting label for $q(a_t|s_t)$ *as*
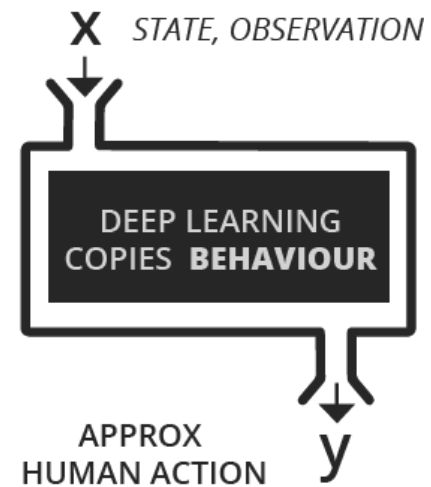$max(R_t + \gamma q(a_*|s_{t+1}))$

$q(a_t|s_t) = neural - network(s_t)$

**FUNCTION APPROXIMATOR**

x

BLACK BOX

$f(x)$ approximates $y$

**DEEP Q LEARNING**

X   *STATE, OBSERVATION*

DEEP LEARNING
COPIES **BLACK BOX**

APPROX
$Q_{LEFT}, Q_{CENTER}, Q_{RIGHT}$   $y$

**BEHAVIORAL CLONING**

X   *STATE, OBSERVATION*

DEEP LEARNING
COPIES **BEHAVIOUR**

APPROX
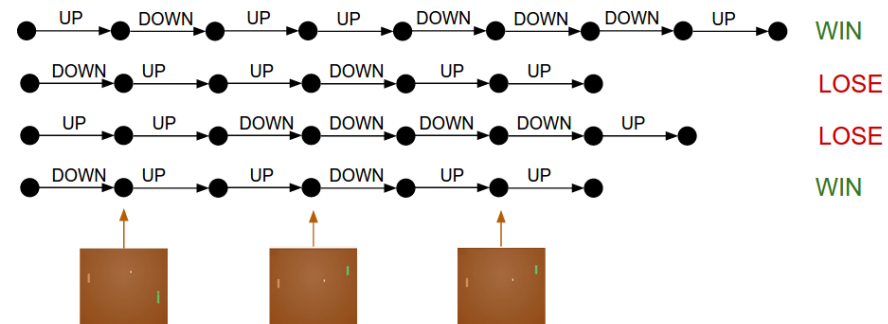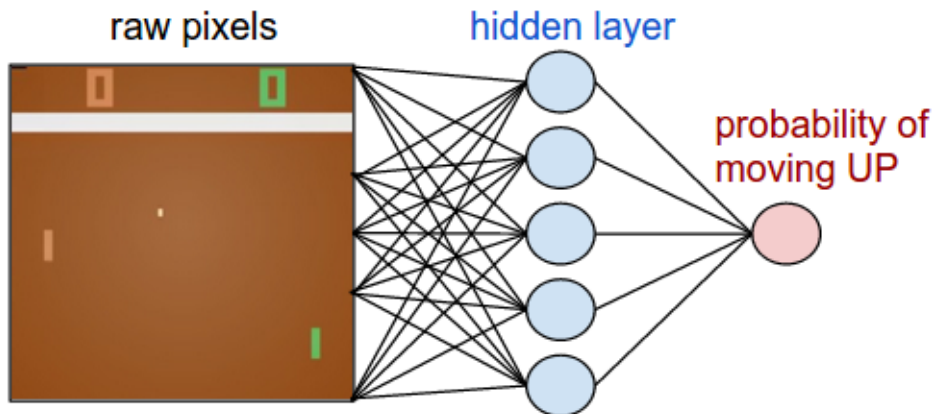HUMAN ACTION   $y$

# DEMO Time

(CARTPOLE IN OPENAI-GYM)

# POLICY GRADIENT APPROACHES

## CHANGE POLICY RATHER THAN Q VALUES

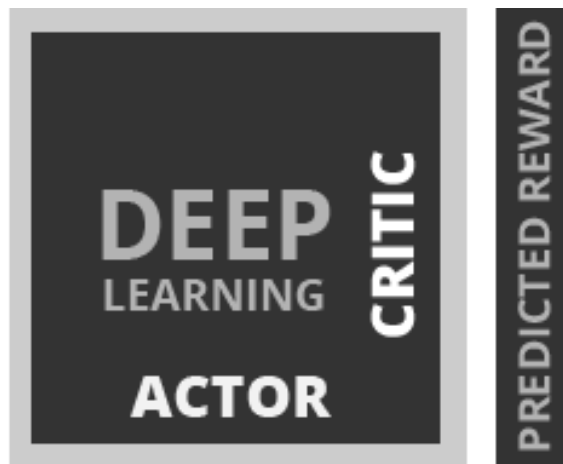POLICY ENTROPY =0.9          POLICY ENTROPY =0.1

# ACTOR CRITIC MODEL

**CRITIC**
*Understands each shot(State),
based on rewards*

**ACTOR**
*Learns to act based on
**advantage factor**
given by critic.*

DEEP LEARNING

CRITIC

ACTOR

$p_1$ $p_2$ $p_3$ $p_4$ $p_5$

PREDICTED REWARD

$-$ EXPECTED REWARD

PREDICTED REWARD

$=$ ADVANTAGE

$R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-t-2} R_T$

ADVANTAGE

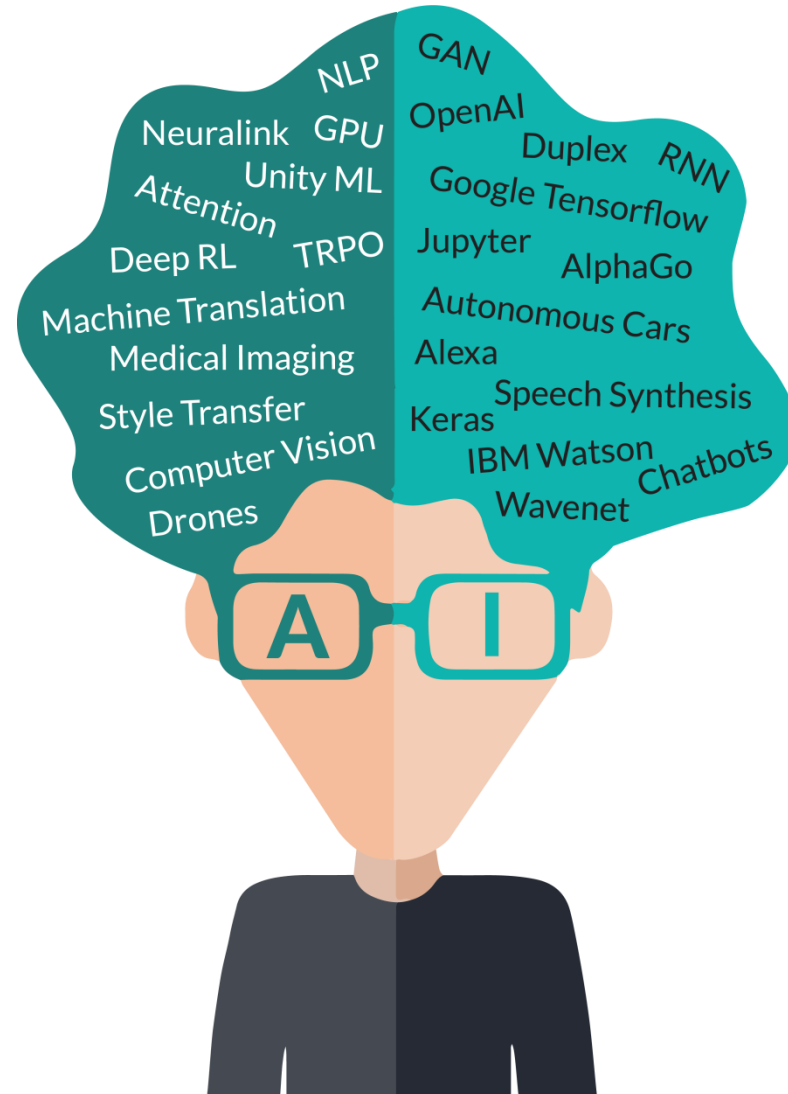HIGH ENTROPY

LOW ENTROPY

DISADVANTAGE

# Other cool stuff

Proximal policy optimization
Monte Carlo Tree Search

# Subscribe

## www.youtube.com/c/crazymuse

# Contact Details

jaley.dholakiya@crazymuse.in
8105207901
**Github link** : www.github.com/crazymuse
**Medium Link** : www.medium.com/crazymuse

*Ping me for any doubt*