# Actor - Critic

## Hyperparams:

Initial policy params $\theta$
Policy representation (ANN? Linear?)
Actor step size $\alpha$
Critic step size $\beta$.
Value function representation, $v_w$.
Initial value function weights $w$.

For each episode:
  For each time $t$
    Agent observes $S_t$
    Agent selects action $A_t$ using $\pi_\theta$
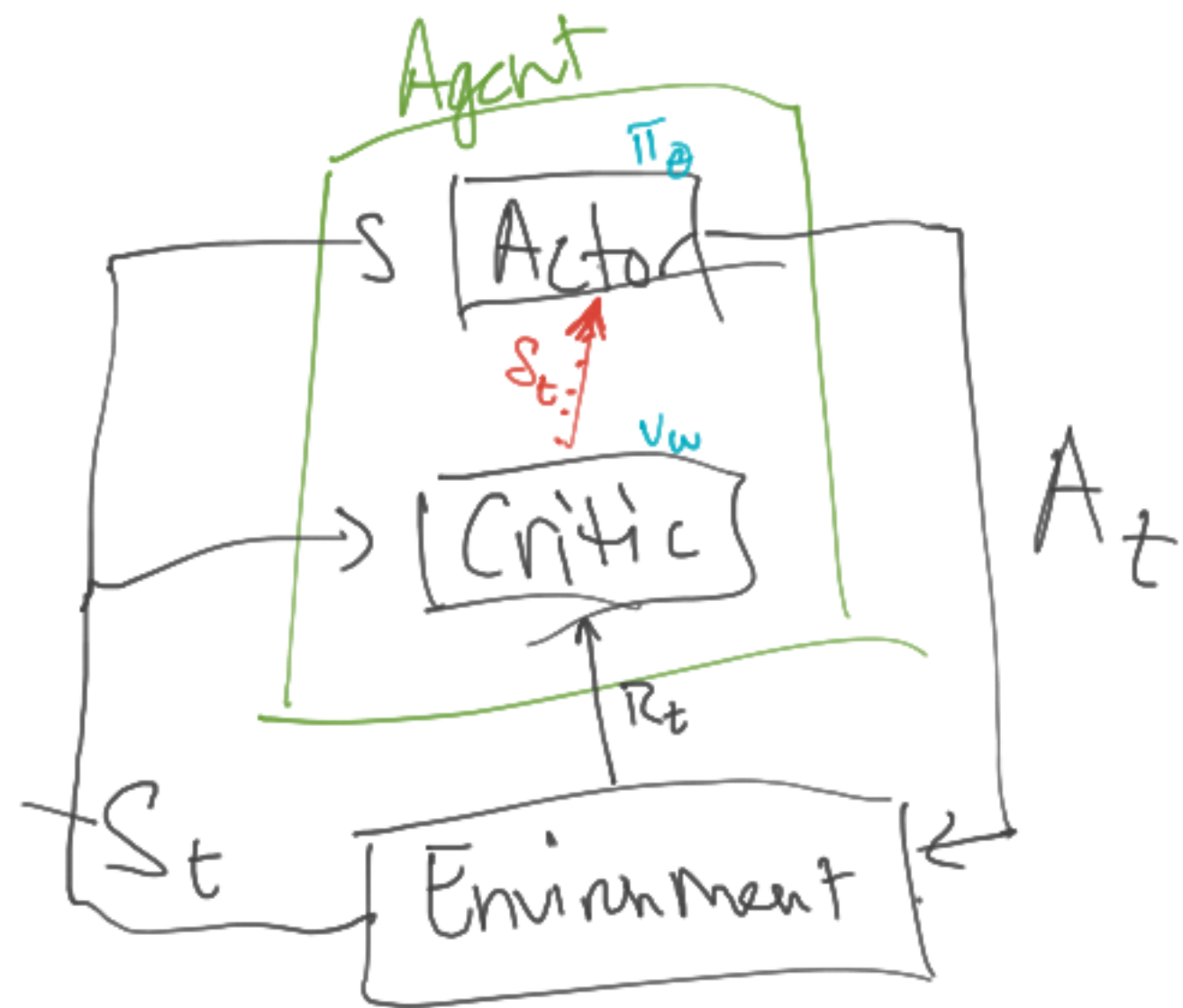    Env responds with $S_{t+1}$ and $R_t$
    $\delta_t = R_t + \gamma v_w(S_{t+1}) - v_w(S_t)$    // TD-error
    $\forall i, \quad \theta_i \leftarrow \theta_i + \alpha \gamma^t \delta_t \dfrac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$   // Actor update
    $\forall j, \quad w_j \leftarrow w_j + \beta \delta_t \dfrac{\partial v_w(S_t)}{\partial w_j}$    // critic update.

Theory says to include In practice it is **bad**. Almost nobody includes this term.

$$v^\pi(S) = \mathbb{E}\left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k} \,\Big|\, S_t = s \,;\, \pi \right]$$



Agent — Actor $\pi_\theta$ — Critic $v_w$ — Environment diagram with $S$, $S_{t-1}$, $R_t$, $S_t$, $A_t$

<u>Psychology</u>

<u>Operant conditioning</u>: learning process through which the strength of a behavior is modified by reward or punishment.

(control)
(searching for a better policy)
- learning a policy.

<u>Classical conditioning</u>: learning procedure in which a biologically potent stimulus (e.g., food) is paired with a previously neutral stimulus (e.g., bell).

(prediction)
- learning a value function.

- Thorndike's Puzzle Boxes.
1898

Sutton & Barto
2nd edition.

RLDM

(performance)
reward.
return

episodes

# Neuroscience

## Dopamine

- Contraction of 3,4 - dihydroxyphenethylamine
- Chemical
- Neurotransmitters.

Cell Body

Axon

Dendrites

chemical "message" (neurotransmitter)

Many different neurotransmitters

synapse

dendrite of post synaptic neuron.

Axon pre-synaptic neuron.

dendrite

Dopamine neuron (Dopaminergic neuron) is a neuron that emits the neurotransmitter dopamine.

Two clusters in mammals:

SNpc and VTA

↓                    ↓
substantia      ventral
nigra            tegmental
pars             area.
compacta

SNpc →dopamine→ Striatum

↳ coordinates motor & action planning, decision making, motivation.
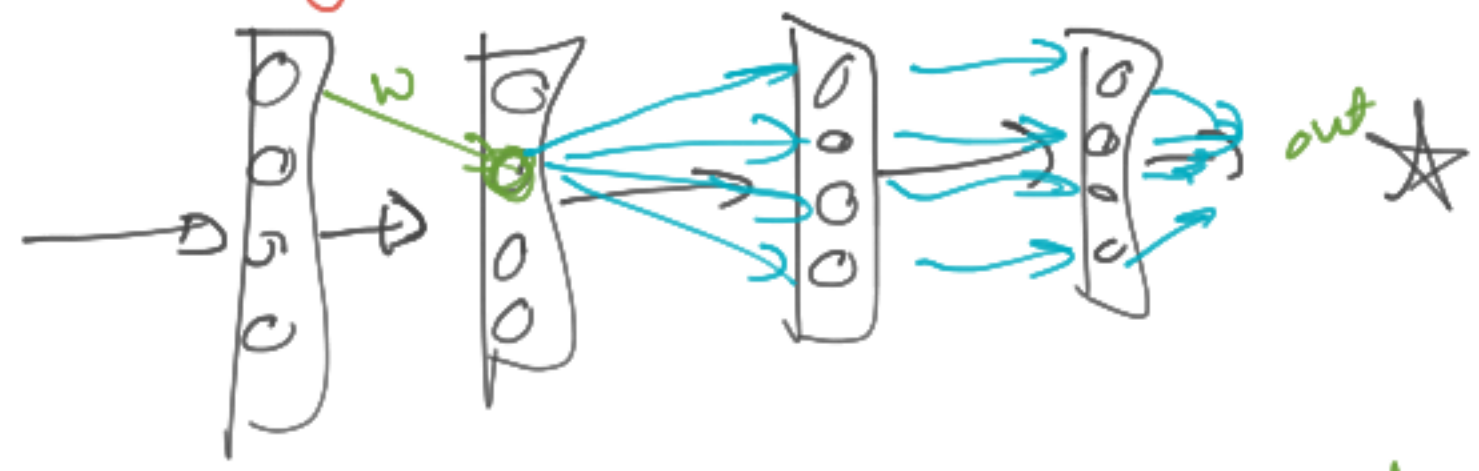
VTA →dopamine→ numerous areas including prefrontal cortex → planning, personality decision making

# Reward Prediction Error Hypothesis for dopamine.

TD error:

Olds & Milner 1954 : Dopamine $\propto$ reward.

---

Brains ✔ do not implement backpropagation.

probabably



$\dfrac{\partial \text{ out}}{\partial w}$

$\Delta_j$ terms propagated backwards through the network.

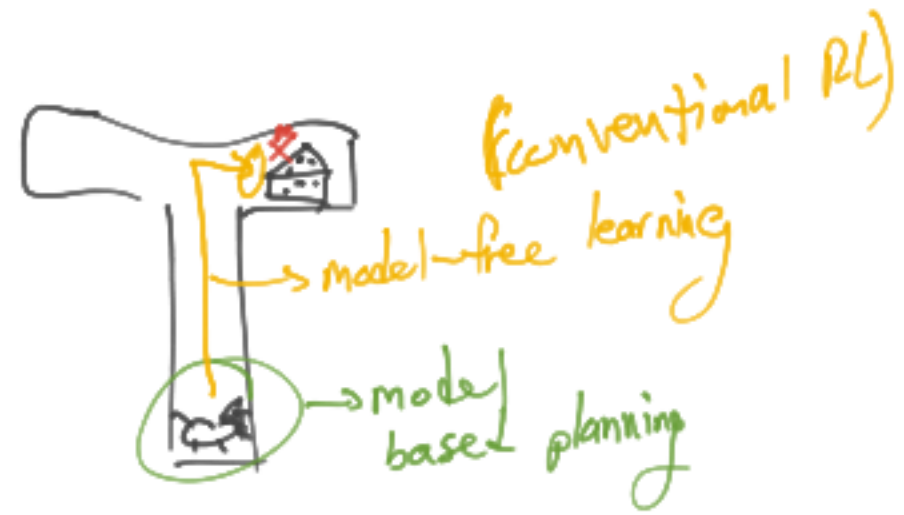→ Information does not seem to pass backwards down the axon.

Each "neuron" can update using only its input, output, and $S_t$.
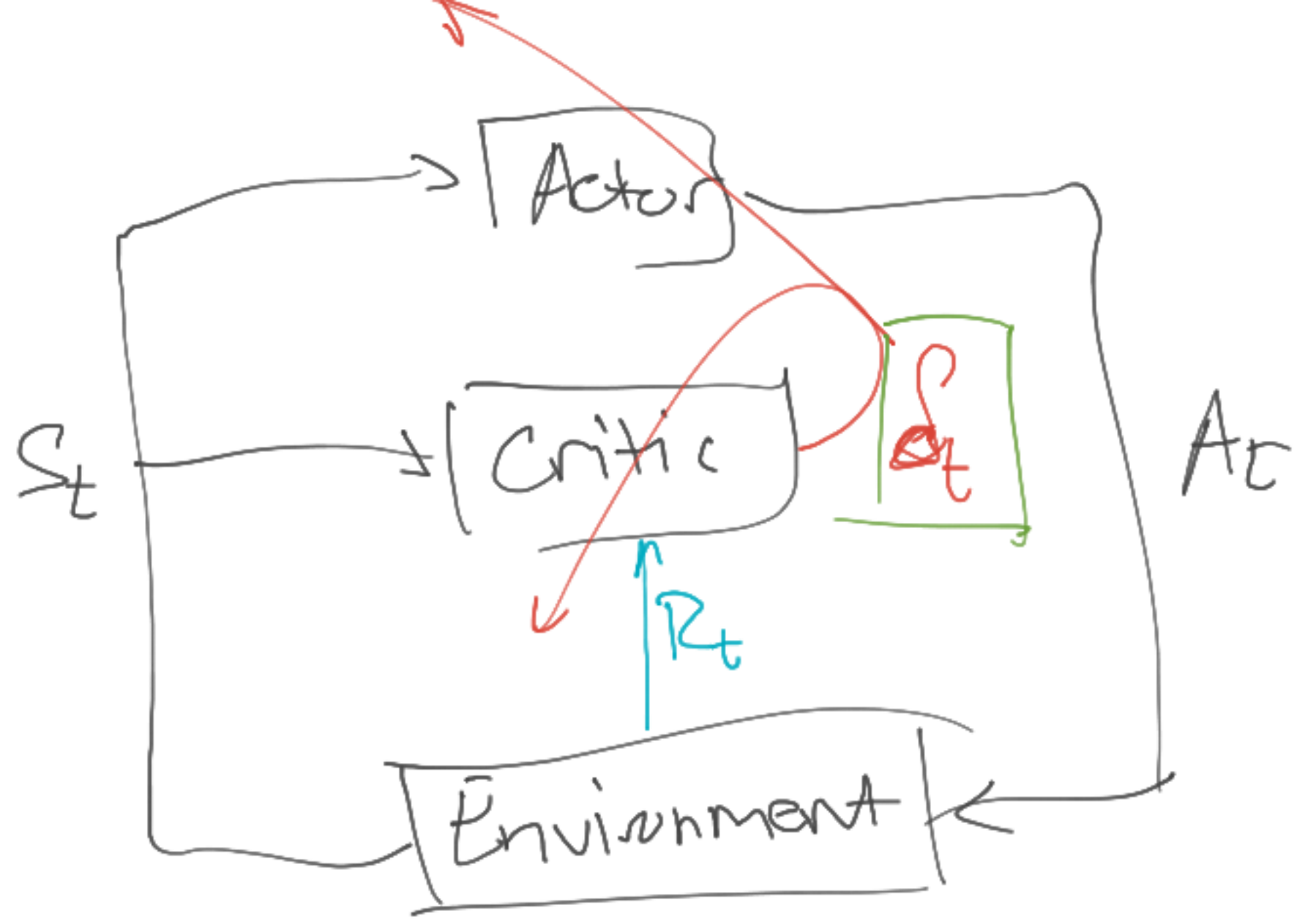
Duplicate of network in reverse.

Coagent networks

Extention of "learning automata"

→ training RL networks without backprop.

# Reward devaluation Studies.



(conventional RL)

→ model-free learning

→ model based planning

→ early learning is model-based planning

→ transitions to a model-free policy over time.

Addiction
Octopi



$S_t$

Actor

Critic

$a_t$

$A_t$

$R_t$

Environment

$\pi^{\text{eat sugar}}(\text{cake in front})$

$\nu^{\text{heroin}}(\text{drugs}) =$