# Advanced Topics in AI
## Temporal Difference Value Learning



Instructor: Prof. Dr. techn. Wolfgang Nejdl

Leibniz University Hannover
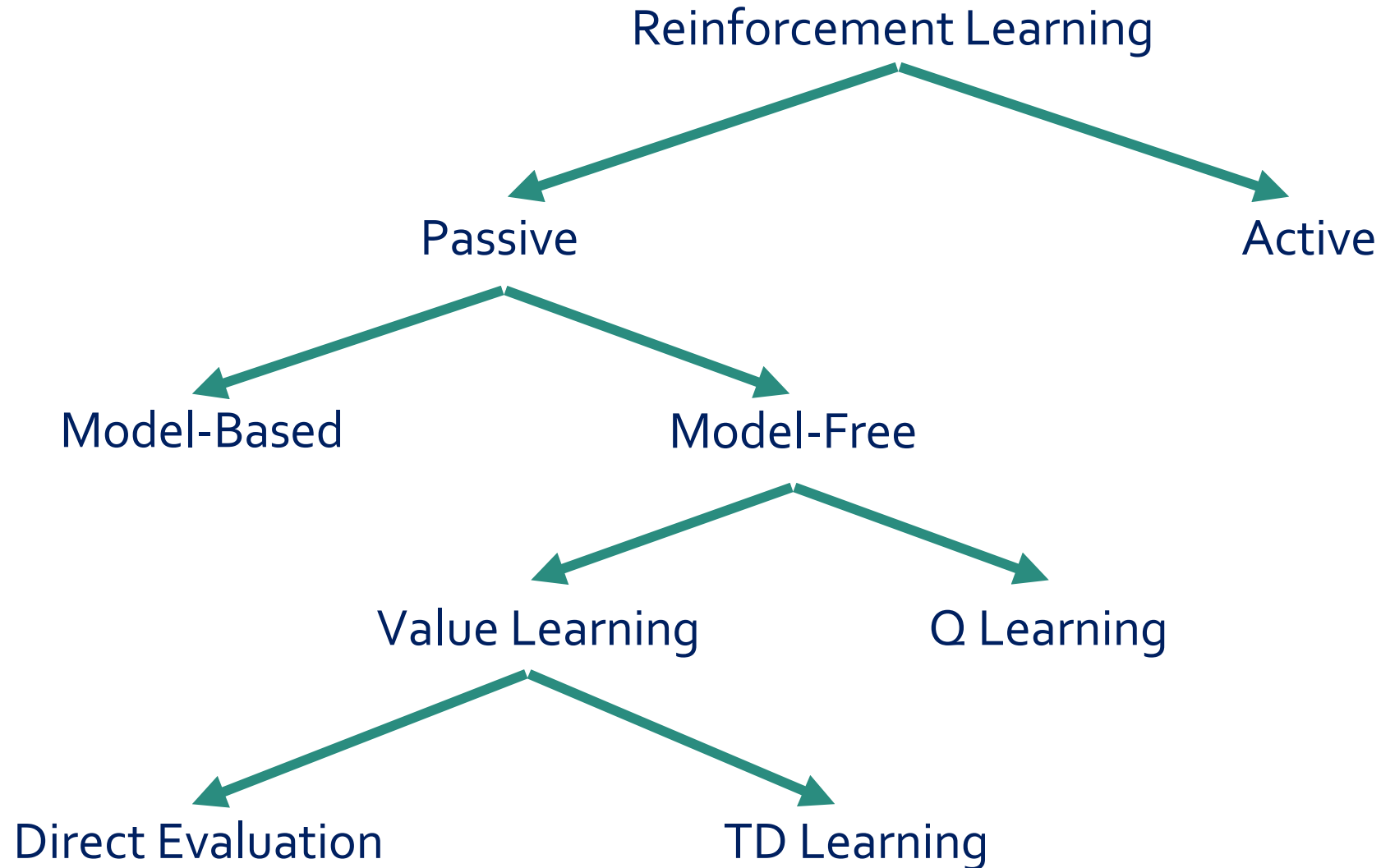
# Reinforcement Learning Taxonomy

# Temporal Difference Value Learning

- Big idea: learn from every experience!
  - Update V(s) each time we experience a transition (s, a, s', r)
  - Likely outcomes s' will contribute updates more often

- Temporal difference learning of values
  - Policy still fixed, still doing evaluation!
  - Move values toward value of whatever successor occurs: running average

$\pi(s)$

s

s, $\pi(s)$

s'

Sample of V(s):     $\text{sample} = R(s, \pi(s), s') + \gamma V^{\pi}(s')$

Update to V(s):     $V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha \cdot \text{sample}$

Same update:        $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha\big(\text{sample} - V^{\pi}(s)\big)$

# Example: TD Value Learning

## States



Assume: $\gamma = 1$, $\alpha = 1/2$

## Observed Transitions

B, east, C, -2

C, east, D, -2



$$V^{\pi}(s) \leftarrow (1 - \alpha)V^{\pi}(s) + \alpha[R(s, \pi(s), s') + \gamma V^{\pi}(s')]$$

# TD Learning in the Brain

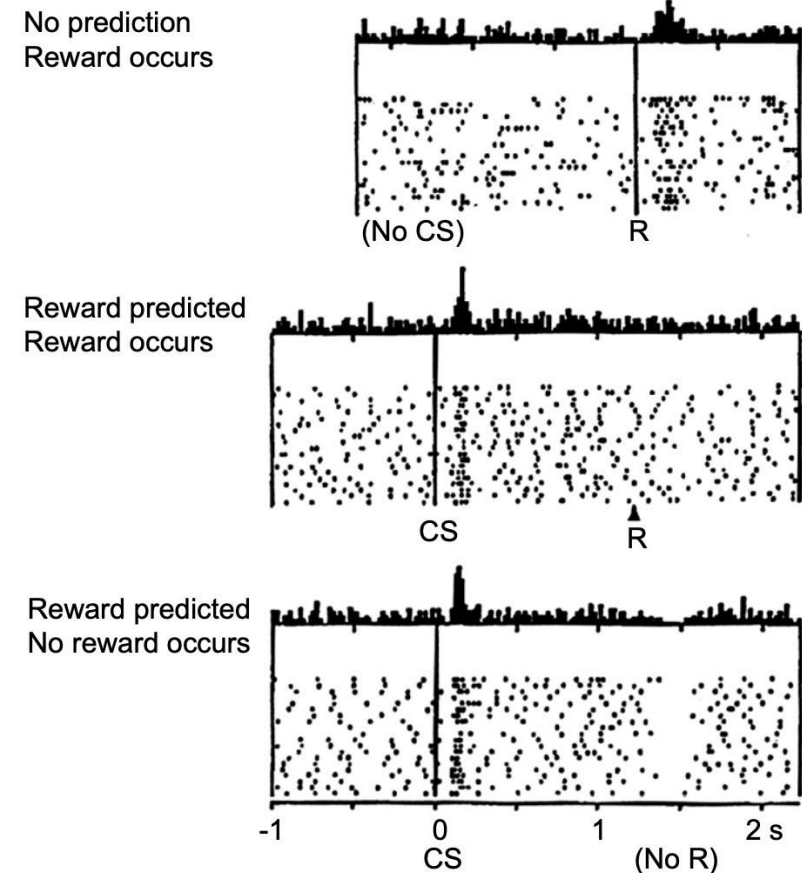- Neurons transmit Dopamine to encode reward or value prediction error
  - $V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha\big(\text{sample} - V^{\pi}(s)\big)$

- Example of Neuroscience & RL informing each other

- For more examples, see
  [AI and Neuroscience: A virtuous circle]
  - https://www.deepmind.com/blog/ai-and-neuroscience-a-virtuous-circle



**Do dopamine neurons report an error in the prediction of reward?**

No prediction
Reward occurs

(No CS)          R

Reward predicted
Reward occurs

CS          R

Reward predicted
No reward occurs

-1          0          1          2 s
            CS                    (No R)

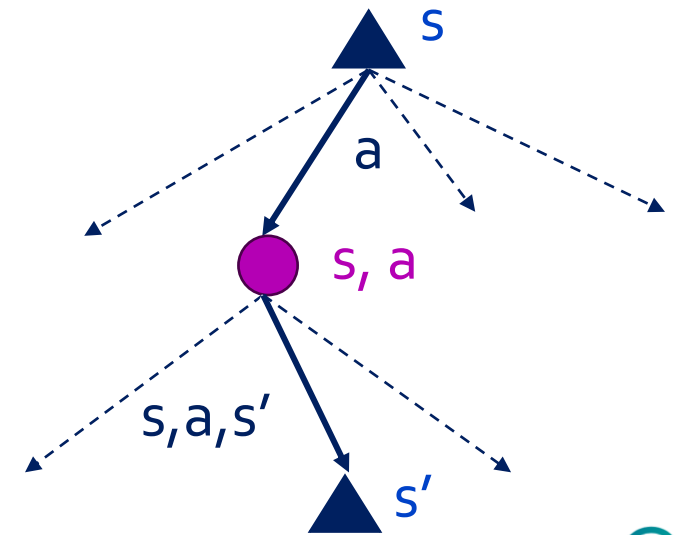[A Neural Substrate of Prediction and Reward. Schultz, Dayan, Montague. 1997]

# Problems with TD Value Learning

- TD value leaning is a model-free way to do policy evaluation, mimicking Bellman updates with running sample averages

- However, if we want to turn values into a (new) policy, we're sunk:
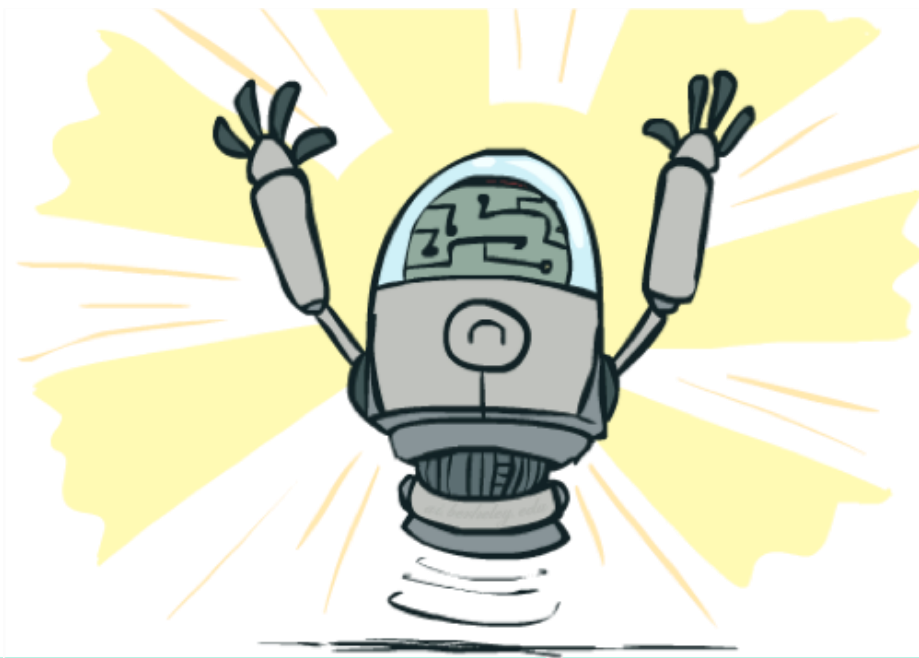
$$\pi(s) = \arg\max_a Q(s, a)$$

$$Q(s, a) = \sum_{s'} T(s, a, s')[R(s, a, s') + \gamma V(s')]$$

- Idea: learn Q-values, not values

- Makes action selection model-free too!

# Advanced Topics in AI

## Next: Q-Learning



Instructor: Prof. Dr. techn. Wolfgang Nejdl

Leibniz University Hannover