

# Path following with reinforcement learning for autonomous cars

- Mozzam Motiwala (IAS)

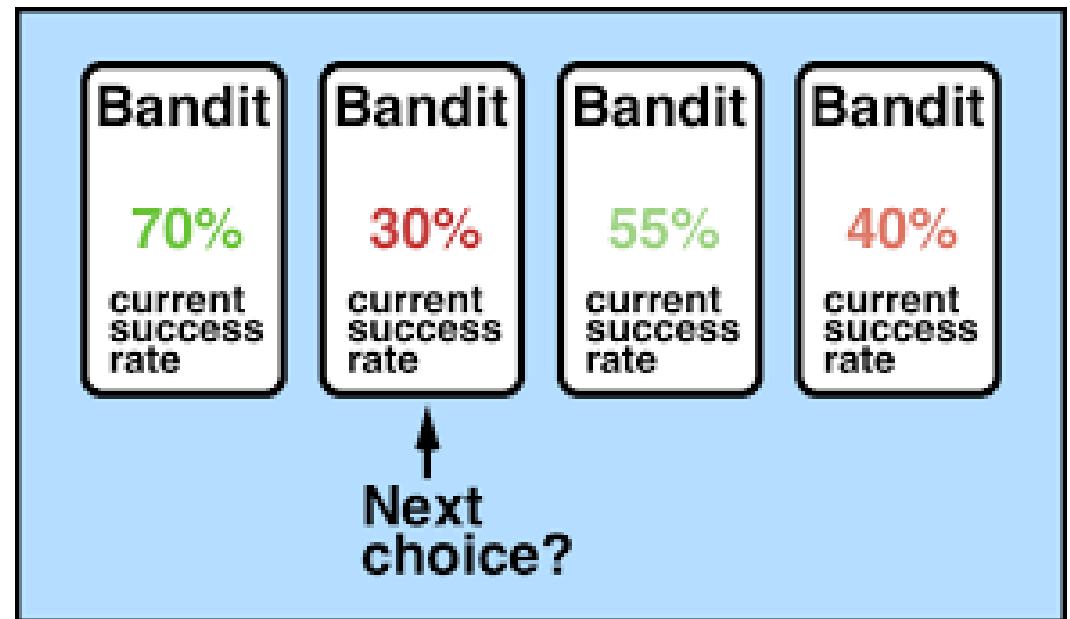


# Index

- **Basics of Reinforcement Learning**
- Model Based vs Model Free Reinforcement Learning
- Autonomous Car collision avoidance

# What is Reinforcement Learning?

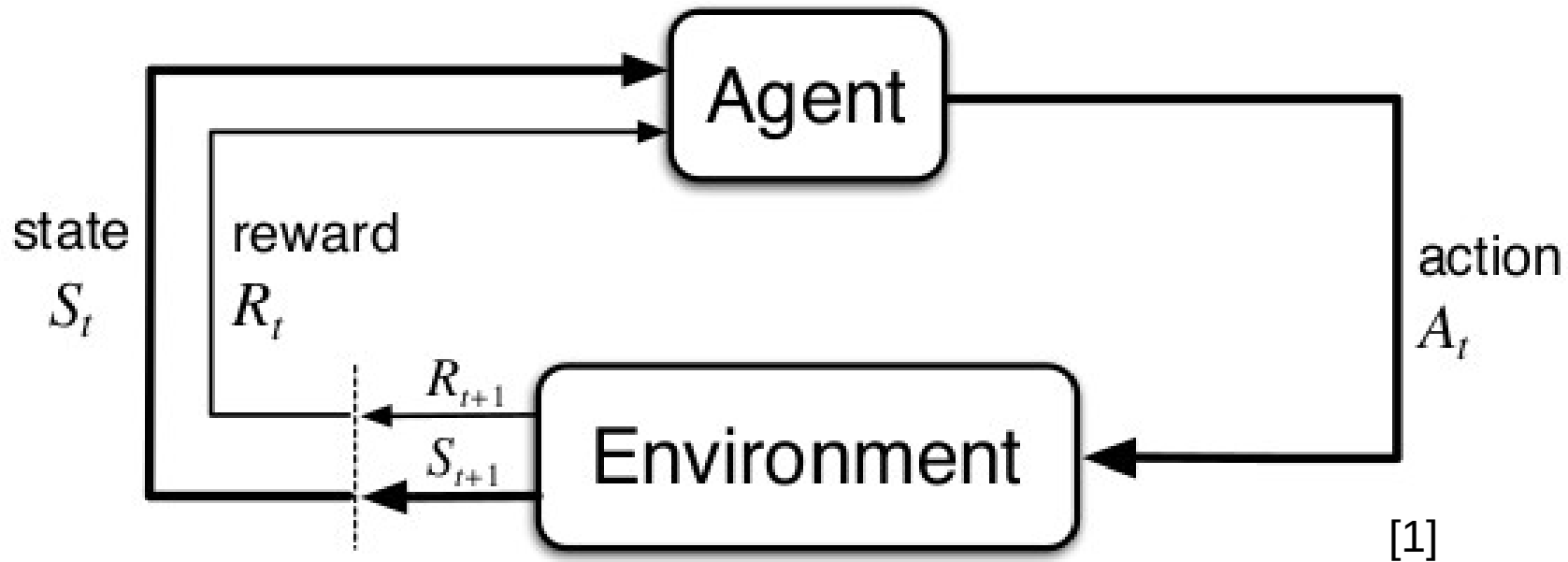
- Learning by trial and error only based on a reward signal[1]



<https://towardsdatascience.com/solving-the-multi-armed-bandit-problem-b72de40db97c>

Exploration vs Exploitation?

# Markov-Decision Process



Reward Function?

Transition Function?

Policy?

Optimal Policy?

# Some terminology

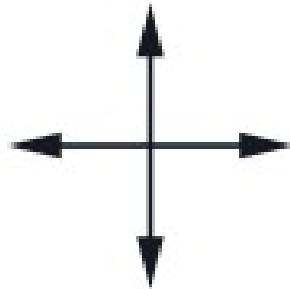
- **Value Function:**  $v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s]$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_{\pi}(s') \right], \quad \text{for all } s \in \mathcal{S},$$

- **Action Value Function:**  $q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a]$

Why Discounting Factor?

# Gridworld

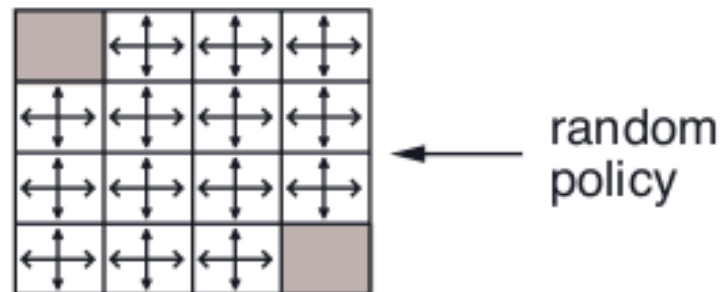


actions

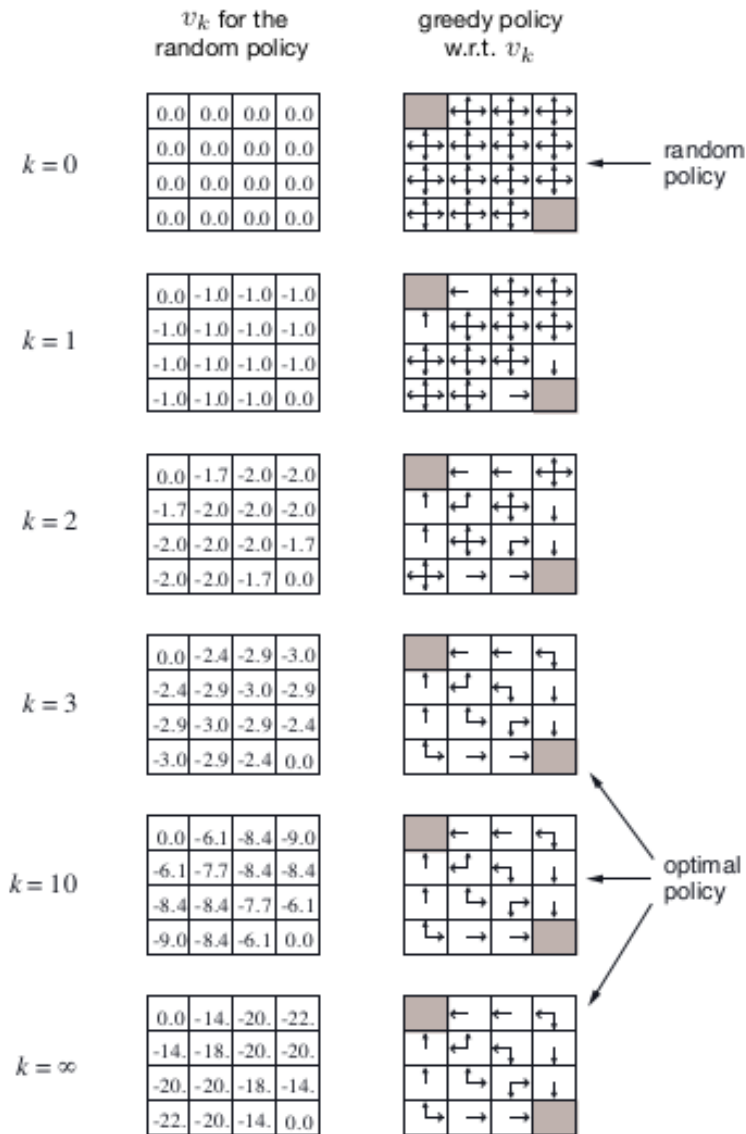
	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$   
on all transitions

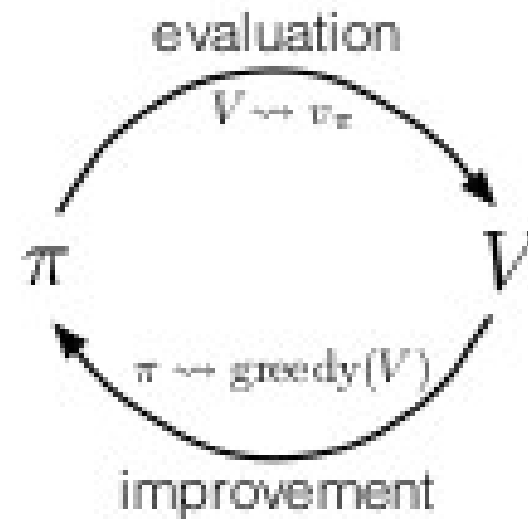
[1]



# Finding Optimal Policy



[1]



# Cart Pole Balancing Problem

1 Observation:

2 Type: Box(4)

3	Num	Observation	Min	Max
4	0	Cart Position	-4.8	4.8
5	1	Cart Velocity	-Inf	Inf
6	2	Pole Angle	-24°	24°
7	3	Pole Velocity At Tip		

8

9 Action:

10 Type: Discrete(2)

11	Num	Action
12	0	Push cart to the left
13	1	Push cart to the right



<https://towardsdatascience.com/cartpole-introduction-to-reinforcement-learning-ed0eb5b58288>

<https://www.youtube.com/watch?v=Lt-KLtkDIh8>

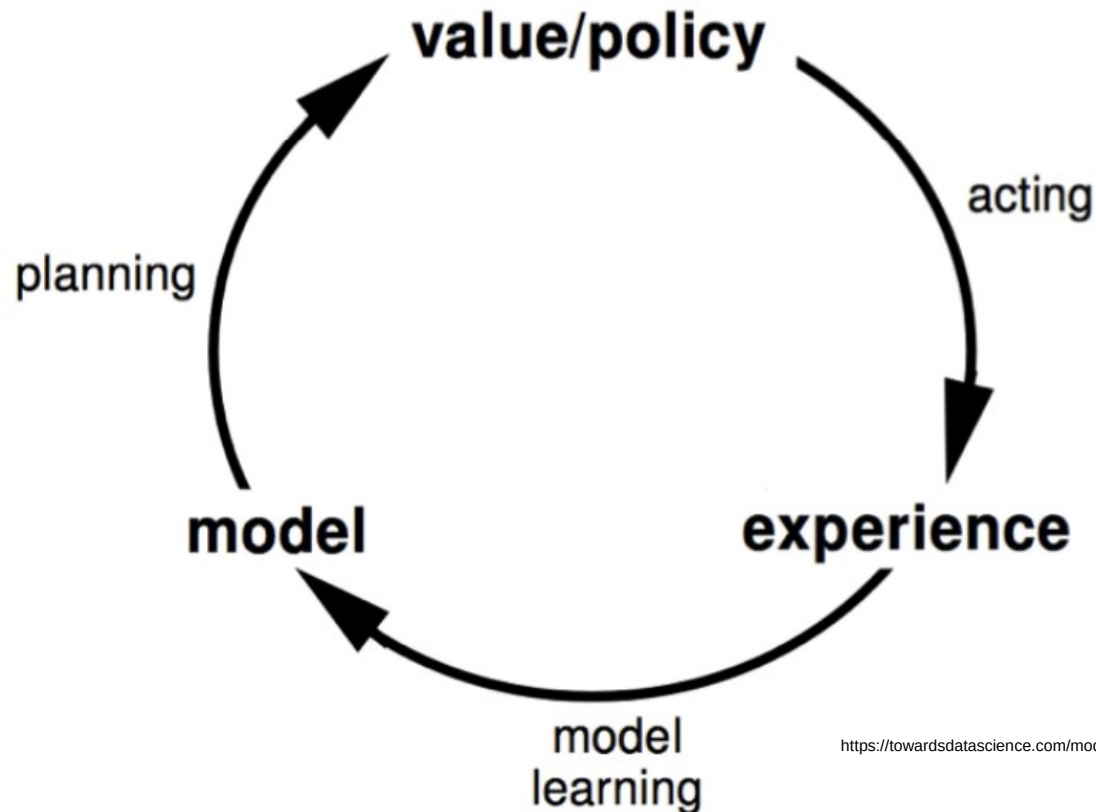


# Index

- Basics of Reinforcement learning
- Model Based vs Model Free Reinforcement Learning
- Autonomous Car collision avoidance

# Model-based

By a model of the environment we mean anything that an agent can use to predict how the environment will respond to its actions[2].



# Example

<https://towardsdatascience.com/model-based-reinforcement-learning-cb9e41ff1f0d>

Two states  $A, B$ ; no discounting; 8 episodes of experience

$A, 0, B, 0$

$B, 1$

$B, 1$

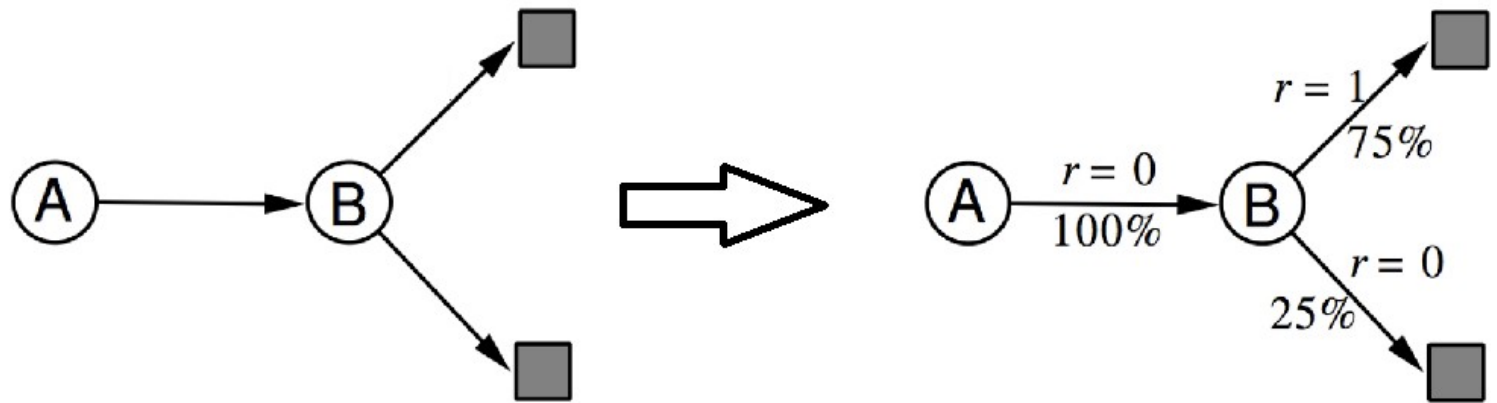
$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 0$



Whats Next? :: Now lets sample from it to adjust the policy..

# Why model-based RL?

Reduced number of interaction with the real environment while learning.

**Types:** Neural Network Model, Gaussian Process Model.. etc

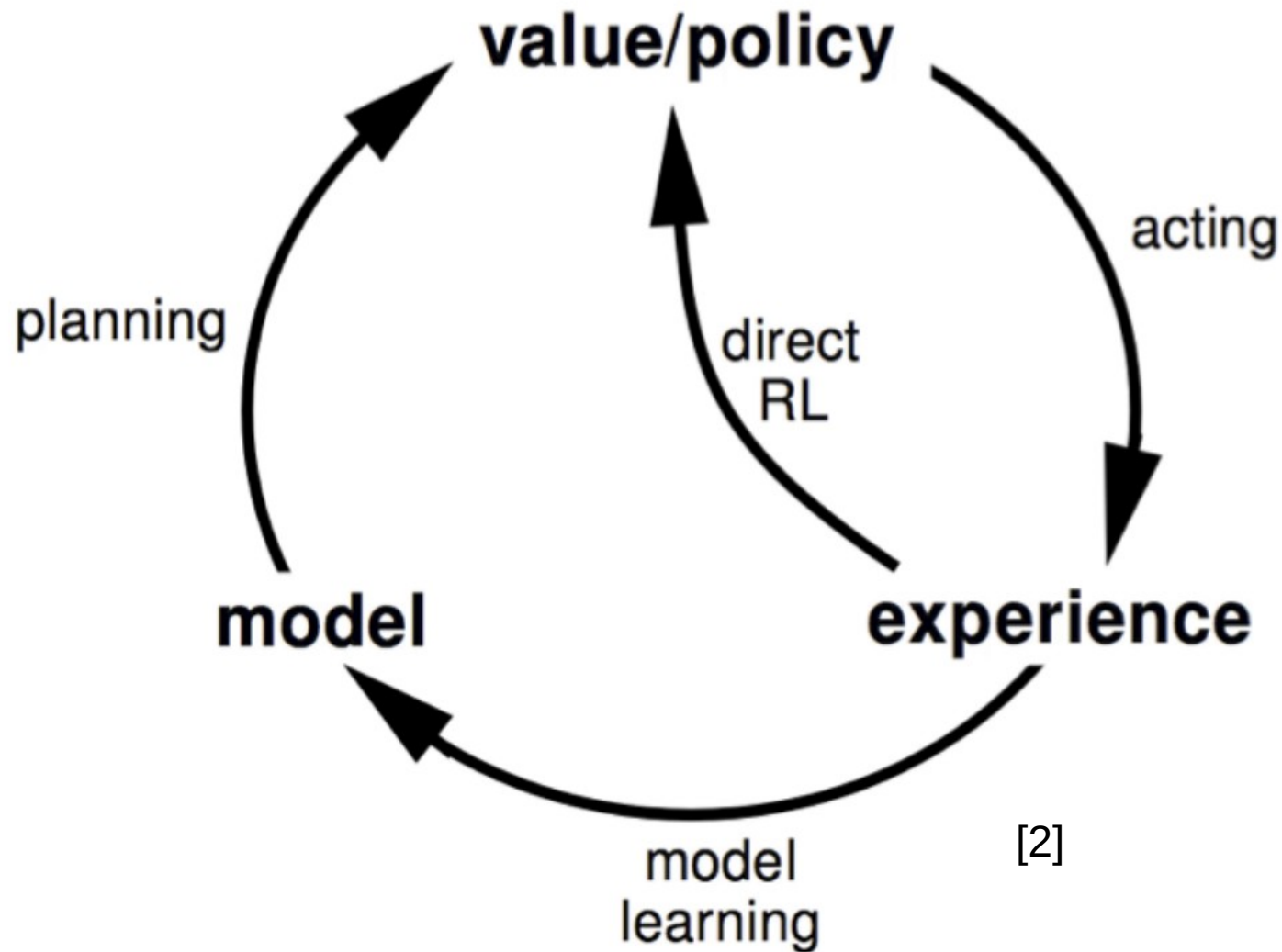
## Advantages?

- Fast
- Need less data

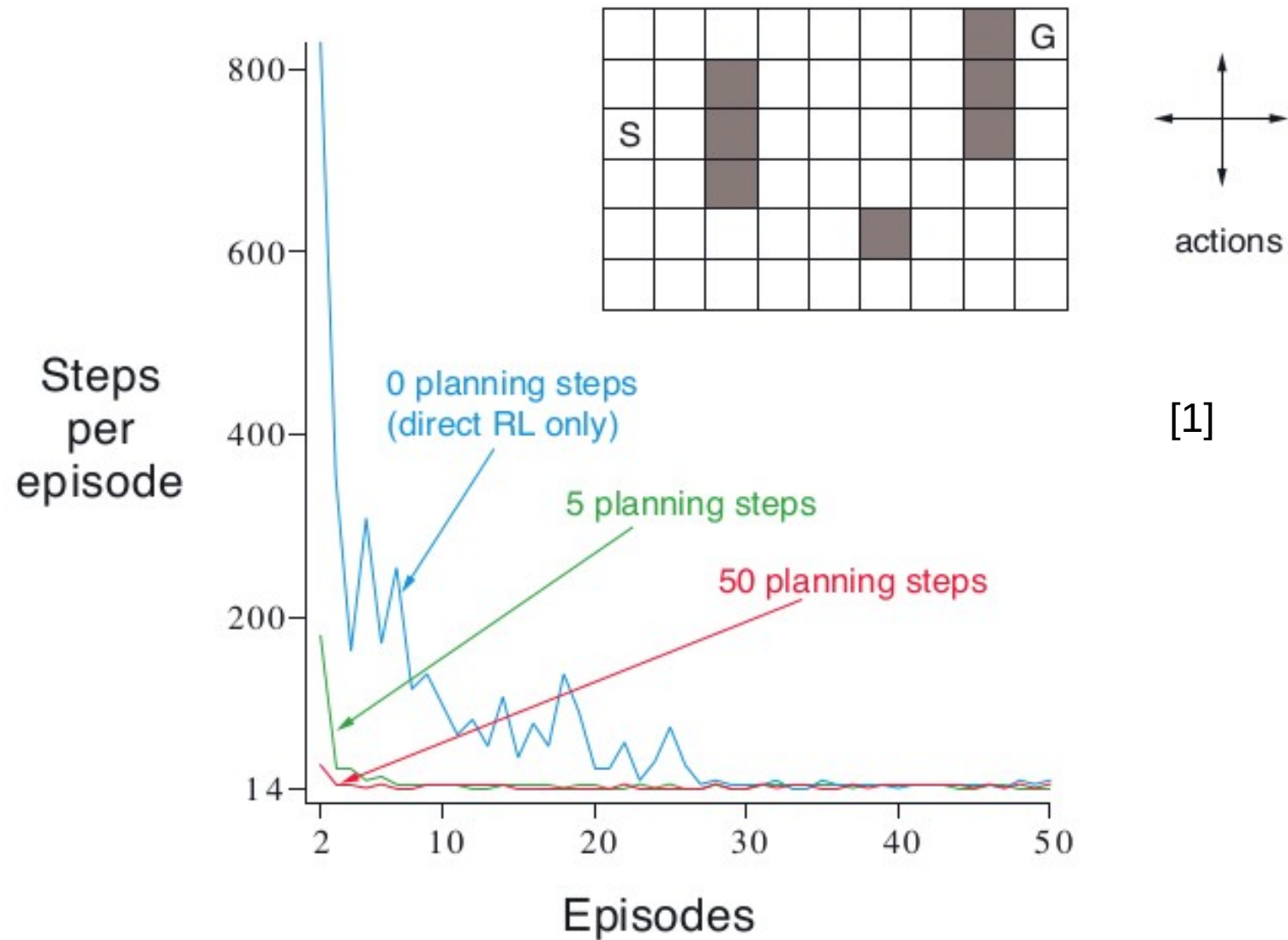
## Problems?

- What if the model is wrong?

# Model Based+ Model Free



# Results



# Why better result?

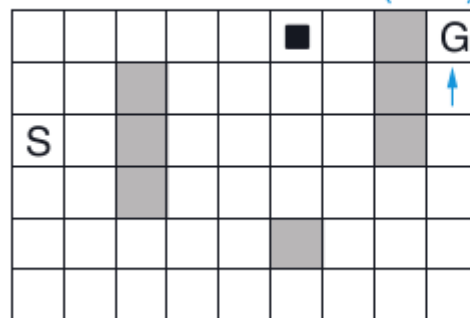
## Tabular Dyna-Q

Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$

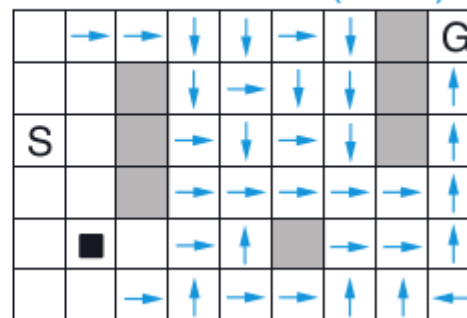
Loop forever:

- (a)  $S \leftarrow$  current (nonterminal) state
- (b)  $A \leftarrow \epsilon$ -greedy( $S, Q$ )
- (c) Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$
- (d)  $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e)  $Model(S, A) \leftarrow R, S'$  (assuming deterministic environment)
- (f) Loop repeat  $n$  times:
  - $S \leftarrow$  random previously observed state
  - $A \leftarrow$  random action previously taken in  $S$
  - $R, S' \leftarrow Model(S, A)$
  - $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

WITHOUT PLANNING ( $n=0$ )



WITH PLANNING ( $n=50$ )



[1]

# Index

- Basics of Reinforcement learning
- Model Based vs Model Free Reinforcement Learning
- **Autonomous Car Collision Avoidance**



# Application: Autonomous Car



## Why Reinforcement Learning?

Problem with traditional methods

- Slow
- Assumptions

Learning in RL

- Adapting to environment
- Learning from mistakes





# Model Details

- Deep RNN as Model
- Model output 1= Current Reward  $\hat{y}$ : Robots speed
- Model output 2= Future Value to go(value of the state)  
 $\hat{b}$ : Distance travelled before collision

- Policy Evaluation Function :

$$J(\mathbf{s}_t, \mathbf{A}_t^H) = \sum_{h=0}^{H-1} -\hat{y}_{t+h} - \hat{b}_{t+H}.$$

- Policy Evaluation by sampling k random action sequence and selecting the one with max reward.

# GCG : Algorithm

---

**Algorithm 1** Reinforcement learning with generalized computation graphs

---

- 1: **input:** computation graph  $G_\theta(\mathbf{s}_t, \mathbf{A}_t^H)$ , error function  $\mathcal{E}_t(\theta)$ , and policy evaluation function  $J(\mathbf{s}_t, \mathbf{A}_t^H)$
  - 2: initialize dataset  $\mathcal{D} \leftarrow \emptyset$
  - 3: **for**  $t = 1$  to  $T$  **do**
  - 4:   get current state  $\mathbf{s}_t$
  - 5:    $\mathbf{A}_t^H \leftarrow \arg \max_{\mathbf{A}} J(\mathbf{s}_t, \mathbf{A})$
  - 6:   execute first action  $\mathbf{a}_t$
  - 7:   receive labels  $y_t$  and  $b_t$
  - 8:   add  $(\mathbf{s}_t, \mathbf{a}_t, y_t, b_t)$  to dataset  $\mathcal{D}$
  - 9:   update  $G_\theta$  by  $\theta \leftarrow \arg \min_{\theta} \mathcal{E}_{t'}(\theta)$  using  $\mathcal{D}$
  - 10: **end for**
-

# Evaluation and Results

<https://www.youtube.com/watch?v=NIFbLVG6LpA>

Distance until crash (m)	Random policy	Double Q-learning with off-policy data	Our approach
Mean	3.4	7.2	52.4
Median	2.8	6.1	29.3
Max	8.0	21.5	197.0

TABLE I: Evaluation of our learned policy navigating at 1.2m/s using only monocular images in a real-world indoor environment after 4 hours of self-supervised training, compared to a random policy and double Q-learning trained with the same data gathered by our approach.

[3]

# Summary

- Benefits of Reinforcement Learning
- Model-Free vs Model-Based
- Combined approach that subsumes Model-free and Model-based

# References

1. R. Sutton and A. Barto, Reinforcement Learning: An Introduction
2. R. Sutton, “Dyna, an Integrated Architecture for Learning, Planning, and Reacting,” in AAI, 1991.
3. G. Kahn, A. Villaflor, B. Ding, P. Abbeel, and S. Levine. Self-Supervised Deep Reinforcement Learning with Generalized Computation Graphs for Robot Navigation. In IEEE International Conference on Robotics and Automation, 2018.



# Question?

**Doubt is not a pleasant  
condition, but certainty is  
absurd.**

Voltaire