



# Reinforcement Learning 101

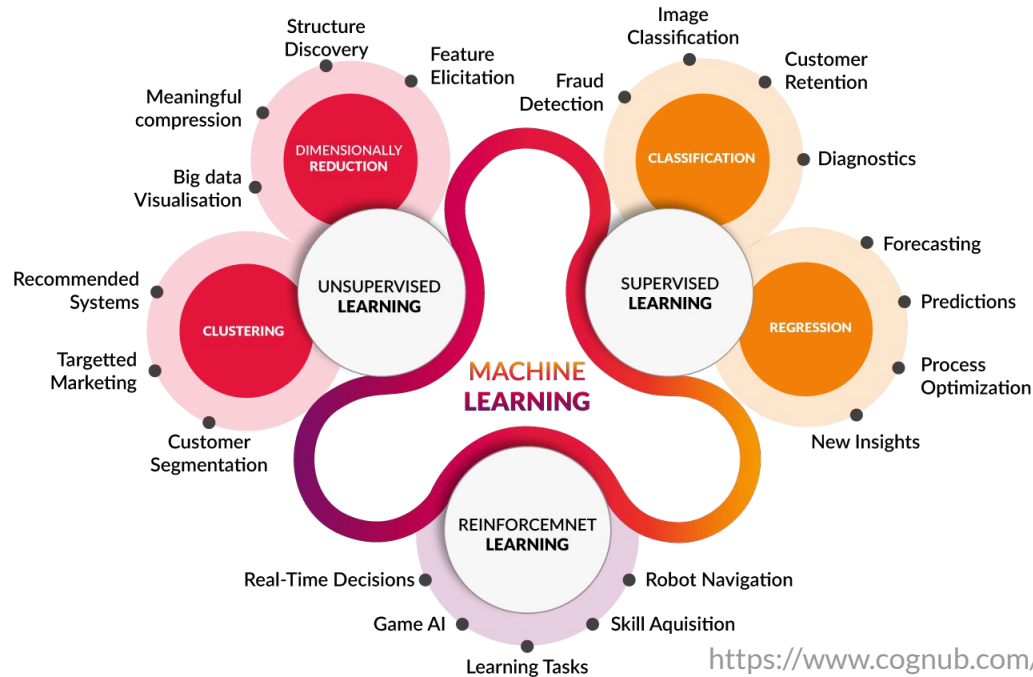
Kim Alyona



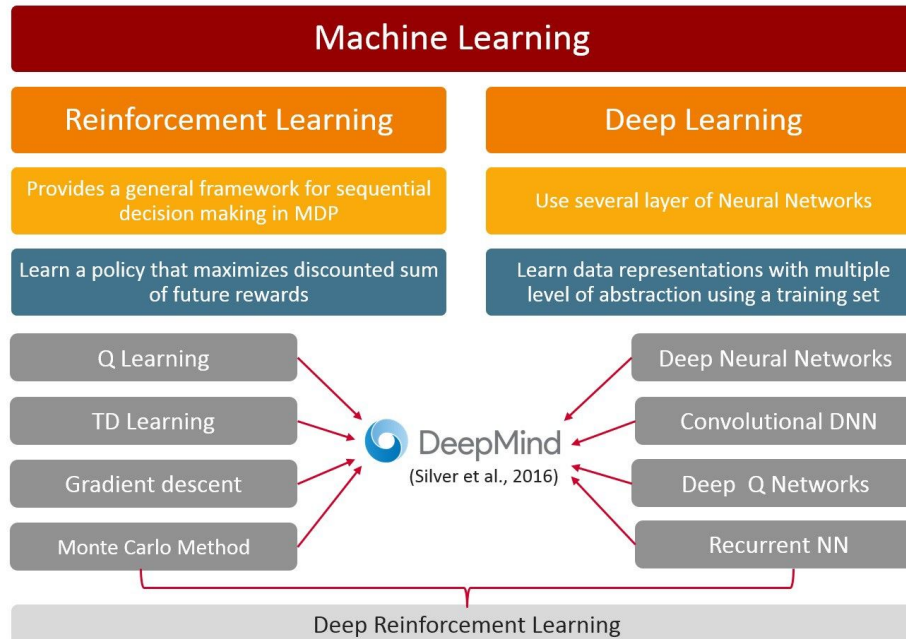
# Summary

1. Introduction
2. Key concepts
3. Approaches
  - a. common
  - b. extra
4. Known problems
5. Applications

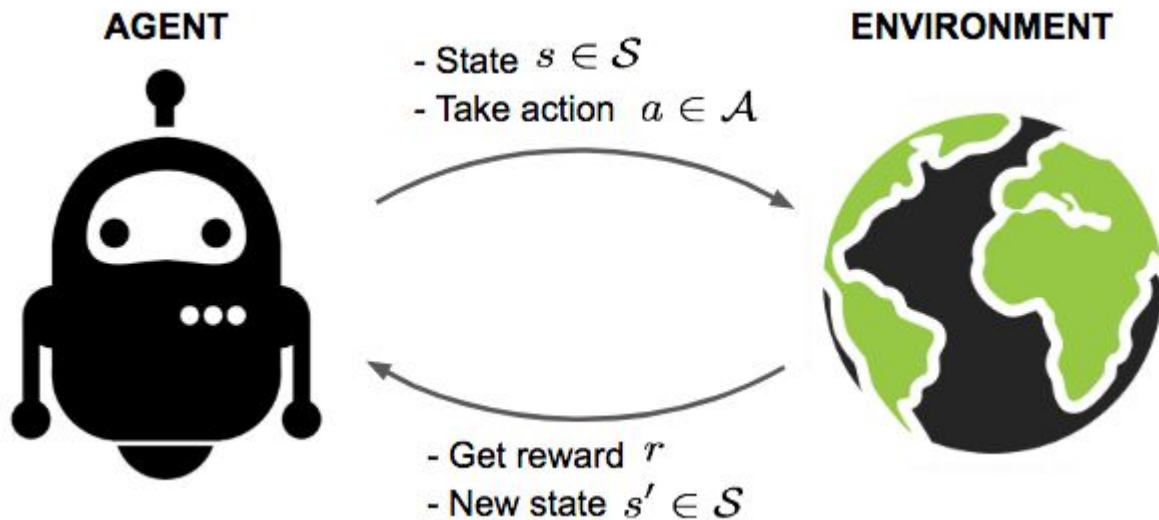
# What is Reinforcement Learning?



# Reinforcement vs Deep



# Key Concepts





## Key Concepts

**Policy**  $\pi(s)$  tries to maximize the sum of rewards. Agent's "brain"

- deterministic  $\pi(s) = a$
- stochastic  $\pi(a|s) = \mathbb{P}_\pi[A = a|S = s]$ 
  - different sampling procedures



## Key Concepts

### Return

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- discounted sum
- more uncertainty in the future
- future actions do not result in immediate benefits
- math convenience



## Key Concepts

State-value function  $V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$

Action-value function  $Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$

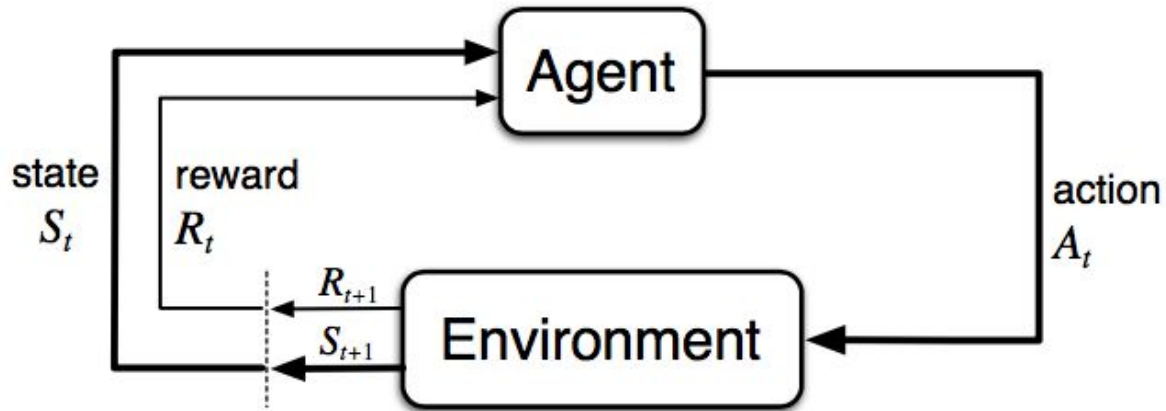
math. property  $V_{\pi}(s) = \sum_{a \in A} Q_{\pi}(s, a) \pi(a|s)$

Advantage  $A(s, a) = Q_{\pi}(s, a) - V_{\pi}(s)$

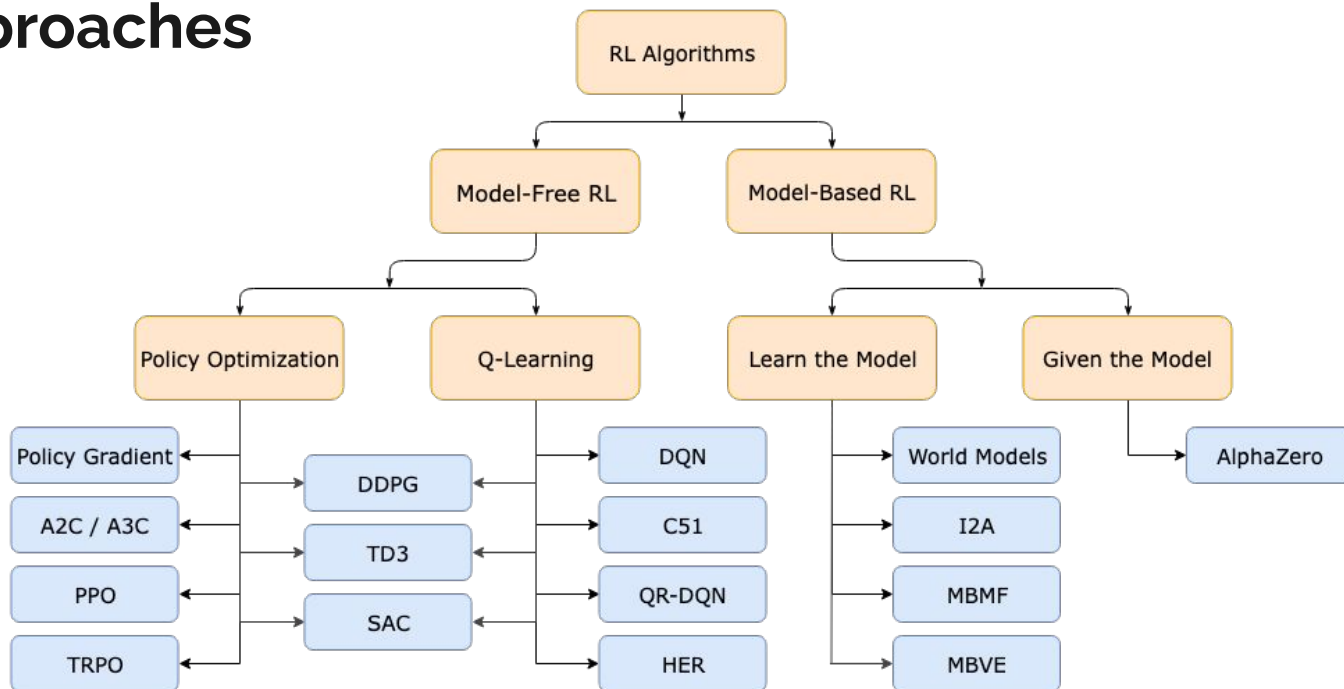


## Markov Decision Process

$$\mathbb{P}[S_{t+1} | S_t] = \mathbb{P}[S_{t+1} | S_1, \dots, S_t]$$



# Approaches





## Approaches

**Model-based:** needs environment. States & Rewards are known or learnt

**Model-free:** does not depend on environment model during training

**On-policy:** use the deterministic **outcomes** or samples from the target policy

**Off-policy:** training on a distribution of **transitions** or episodes produced by a different behavior policy rather than that produced by the target policy.



## Policy Optimization

- model-free & on-policy
- represent a policy explicitly  $\pi_{\theta}(a|s)$
- usually involves learning an approximator for the value-function

Example:

- Asynchronous Advantage Actor-Critic (A3C)
- **Critic:** updates value function  $V(s; w)$  parameters  $w$ .
- **Multiple Actors:** updates policy parameters  $\theta$ , in the direction suggested by the critic,  $\pi(a|s; \theta)$
- $J_V(w) = (G_t - V(s; w))^2$



# Q-learning

- model-free & off-policy
- Bellman equations
- policy:  $a(s) = \arg \max_a Q_\theta(s, a)$
- temporal difference

Example:

- Deep Q-Network (DQN)
  - Experience Replay
  - Periodically Updated Target



## Trade-offs

<b>Policy Optimization</b>	<b>Q-learning</b>
<ul style="list-style-type: none"><li>● you directly optimize for the thing you want (policy)</li><li>● stable and reliable</li><li>● less sample efficient</li></ul>	<ul style="list-style-type: none"><li>● indirectly optimize for agent performance (training to satisfy a self-consistency equation)</li><li>● less stable</li><li>● more sample efficient, can re-use data</li></ul>

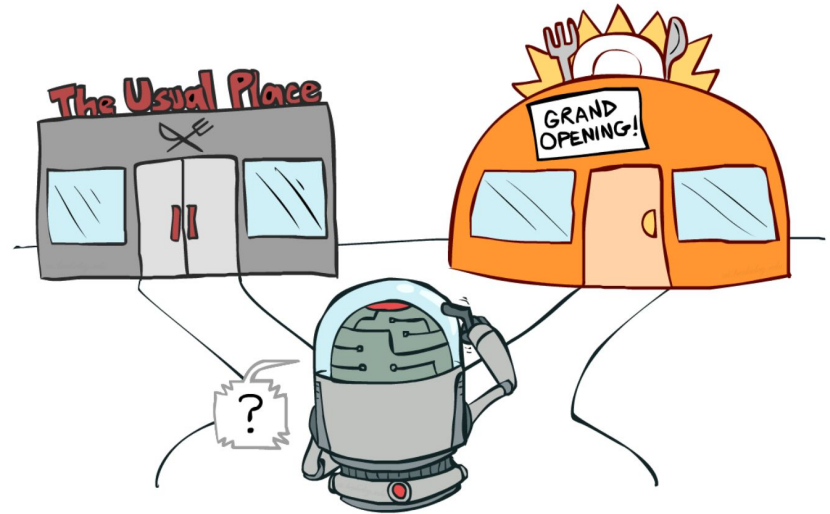


## Model-based RL

- Background: Pure Planning (MBMF)
  - model-predictive control
  - never explicitly represents policy (an optimal plan over fixed time window)
- Expert Iteration (ExIt, AlphaZero)
  - planning algorithm like Monte-Carlo Tree Search
  - sampling from the current policy & evaluate samples with planning algorithm
- Data Augmentation for Model-Free Methods (MBVE)
  - augment real data with synthetic

# Known Problems

- Exploration-Exploitation Dilemma
- Deadly Triad Issue
  - off-policy
  - nonlinear function approximation
  - bootstrapping
  - **unstable learning, does not converge**



[http://ai.berkeley.edu/lecture\\_slides.html](http://ai.berkeley.edu/lecture_slides.html), lecture 11





# Applications

- modeling and explaining neural activity
- linking phasic **dopamine release** with temporal-difference **reward-prediction error** [Niv 19]
- learning complex **robotic skills** from raw **sensory input**
- modeling goal-oriented search **policy** with **top-down factors** as **states**

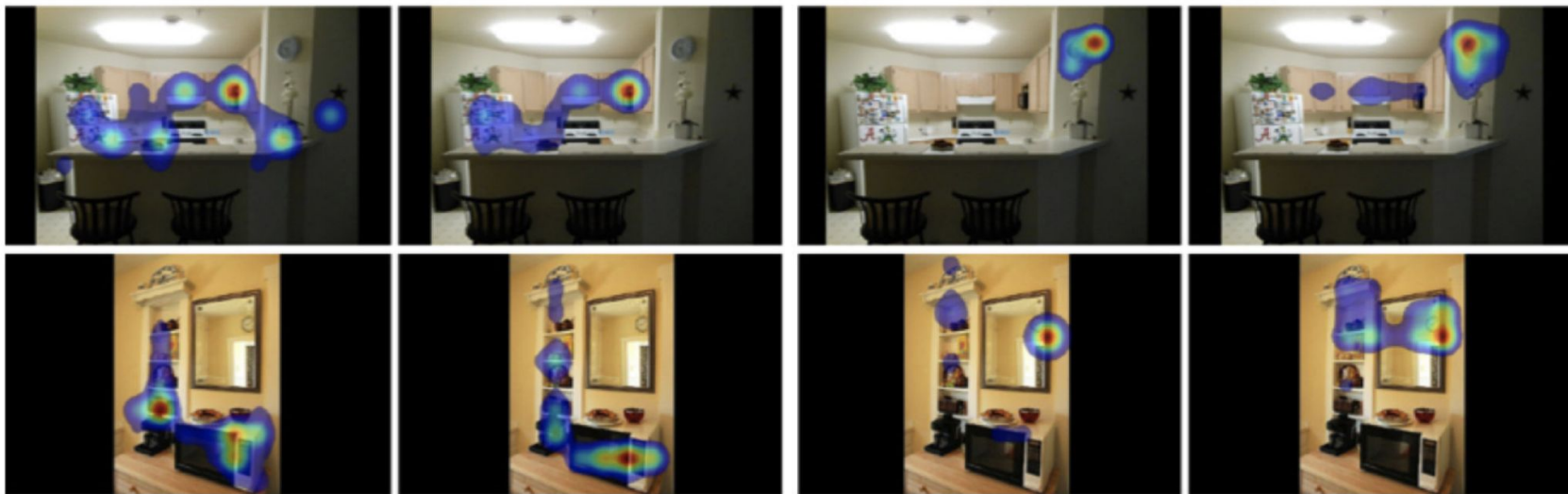
# Spoiler for my next talk

model-micro

human-micro

model-clock

human-clock



## Extra slide

- Curriculum RL
- Meta-Reinforce Learning
- Inverse RL

