

# ECE276B: Planning & Learning in Robotics

## Lecture 1: Introduction

Instructor:

Nikolay Atanasov: [natanasov@ucsd.edu](mailto:natanasov@ucsd.edu)

Teaching Assistant:

Thai Duong: [tduong@eng.ucsd.edu](mailto:tduong@eng.ucsd.edu)

**UC San Diego**

**JACOBS SCHOOL OF ENGINEERING**  
Electrical and Computer Engineering

# What is this class about?

- ▶ **ECE276A**: sensing and estimation in robotics:
  - ▶ how to model robot motion and observations
  - ▶ how to estimate (a distribution of) the robot/environment state  $\mathbf{x}_t$  from the history of observations  $\mathbf{z}_{0:t}$  and control inputs  $\mathbf{u}_{0:t-1}$
- ▶ **ECE276B**: planning and decision making in robotics:
  - ▶ how to select control inputs  $\mathbf{u}_{0:t-1}$  to accomplish a task
- ▶ **References** (not required):
  - ▶ Dynamic Programming and Optimal Control: Bertsekas
  - ▶ Planning Algorithms: LaValle (<http://planning.cs.uiuc.edu>)
  - ▶ Reinforcement Learning: Sutton & Barto (<http://incompleteideas.net/book/the-book.html>)
  - ▶ Calculus of Variations and Optimal Control Theory: Liberzon (<http://liberzon.csl.illinois.edu/teaching/cvoc.pdf>)

# Logistics

- ▶ Course website: <https://natanaso.github.io/ece276b>
- ▶ Includes links to (**sign up!**):
  - ▶ **Canvas**: Zoom meeting schedule and lecture recordings
  - ▶ **Piazza**: discussion – check Piazza regularly because class announcements, updates, etc., will be posted there
  - ▶ **Gradescope**: homework submission and grades
- ▶ Assignments:
  - ▶ 3 theoretical homework sets (16% of grade)
  - ▶ 3 programming assignments in **python** + project report:
    - ▶ Project 1: Dynamic Programming (18% of grade)
    - ▶ Project 2: Motion Planning (18% of grade)
    - ▶ Project 3: Optimal Control (18% of grade)
  - ▶ Final exam (30% of grade)
- ▶ Grades:
  - ▶ assigned based on the class performance, i.e., there will be a curve
  - ▶ **no late policy**: work submitted past the deadline will receive 0 credit

## Prerequisites

- ▶ **Probability theory:** random vectors, probability density functions, expectation, covariance, total probability, conditioning, Bayes rule
- ▶ **Linear algebra/systems:** eigenvalues, positive definiteness, linear systems of ODEs, matrix exponential
- ▶ **Optimization:** gradient descent
- ▶ **Programming:** experience with at least one language (python/C++/Matlab), classes/objects, data structures (e.g., queue, list), data input/output, plotting
- ▶ It is up to you to judge if you are ready for this course!
  - ▶ Consult with your classmates who took ECE276A
  - ▶ Take a look at the material from last year:  
<https://natanaso.github.io/ece276b2020>
  - ▶ If the first assignment seems hard, the rest will be hard as well

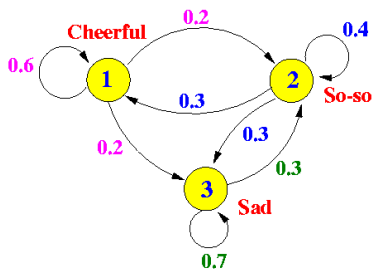
# Syllabus (Tentative)

Date	Lecture	Materials	Assignments
Mar 30	Introduction, Markov Chains	<a href="#">Grinstead-Snell-Ch11</a>	
Apr 01	Markov Decision Processes	<a href="#">Bertsekas 1.1-1.2</a>	
Apr 06	Dynamic Programming	<a href="#">Bertsekas 1.3-1.4</a>	<a href="#">HW1, PR1</a>
Apr 08	Deterministic Shortest Path	<a href="#">Bertsekas 2.1-2.3</a>	
Apr 13	Configuration Space	<a href="#">LaValle 4.3, 6.2-6.3</a>	
Apr 15	Catch-up		
Apr 20	Search-based Planning	<a href="#">LaValle 2.1-2.3, JPS</a>	
Apr 22	Catch-up		
Apr 27	Anytime Incremental Search	<a href="#">RTAA*, ARA*, AD*, Journal Paper</a>	<a href="#">HW2, PR2</a>
Apr 29	Sampling-based Planning	<a href="#">LaValle 5.5-5.6</a>	
May 04	Stochastic Shortest Path	<a href="#">Bertsekas 7.1-7.3</a>	
May 06	Bellman Equations I	<a href="#">Sutton-Barto 4.1-4.4</a>	
May 11	Bellman Equations II	<a href="#">Sutton-Barto 4.5-4.8</a>	
May 13	Catch-up		
May 18	Model-free Prediction	<a href="#">Sutton-Barto 6.1-6.3</a>	<a href="#">HW3, PR3</a>
May 20	Model-free Control	<a href="#">Sutton-Barto 6.4-6.7</a>	
May 25	Value Function Approximation	<a href="#">Sutton-Barto Ch.9</a>	
May 27	Continuous-time Optimal Control	<a href="#">Bertsekas 3.1-3.2</a>	
Jun 01	Pontryagin's Minimum Principle	<a href="#">Bertsekas 3.3-3.4, Liberzon Ch. 2.4 and Ch. 4</a>	
Jun 03	Linear Quadratic Control	<a href="#">Bertsekas 4.1</a>	
Jun 10	Final Exam		

- Check the course website for updates:  
<https://natanaso.github.io/ece276b>

# Markov Chain

- ▶ A **Markov Chain** (MC) is a probabilistic model used to represent the evolution of a robot system
- ▶ The state  $x_t$  can be discrete or continuous and is fully observed
- ▶ The state transitions are random and uncontrolled, determined by a transition matrix or function
- ▶ A **Markov Decision Process** (MDP) is a Markov chain, whose transitions are controlled

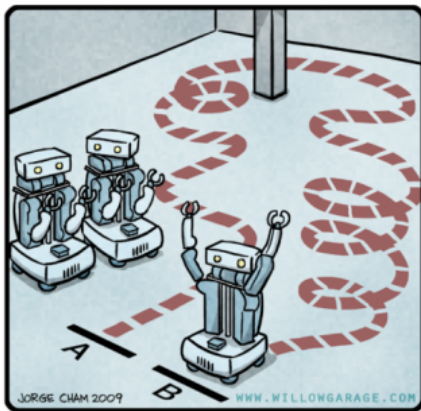


$$P = \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.3 & 0.4 & 0.3 \\ 0.0 & 0.3 & 0.7 \end{bmatrix}$$

$$P_{ij} = \mathbb{P}(x_{t+1} = j \mid x_t = i)$$

# Motion Planning

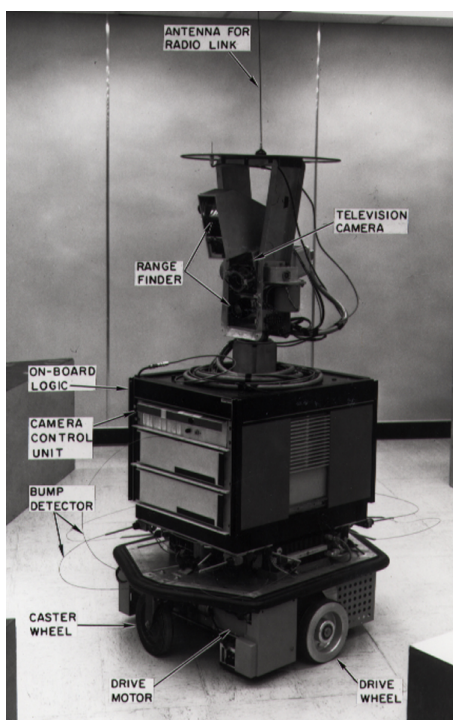
R.O.B.O.T. Comics



"HIS PATH-PLANNING MAY BE  
SUB-OPTIMAL, BUT IT'S GOT FLAIR."

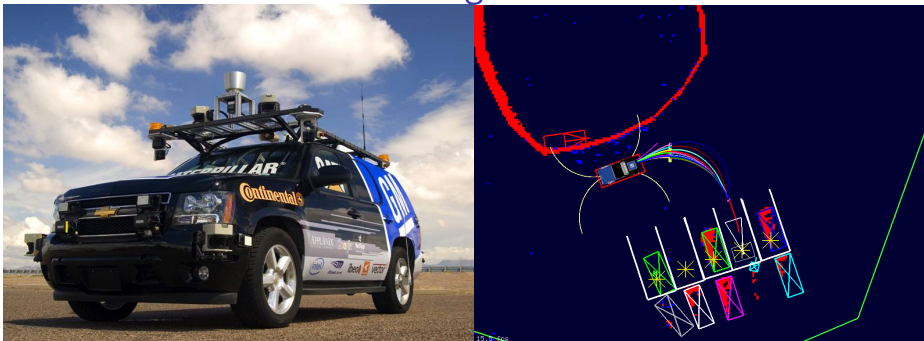
## A\* Search

- ▶ Invented by Hart, Nilsson and Raphael of Stanford Research Institute in 1968 for the Shakey robot
- ▶ Video: <https://youtu.be/qXdn6ynwpiI?t=3m55s>



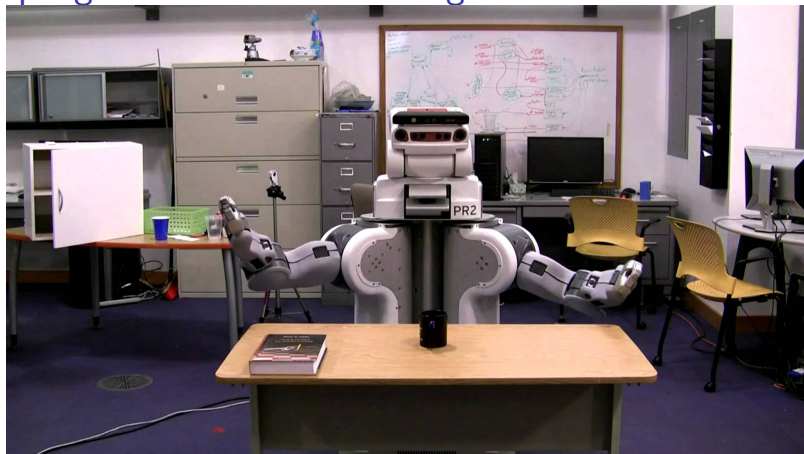


# Search-based Motion Planning



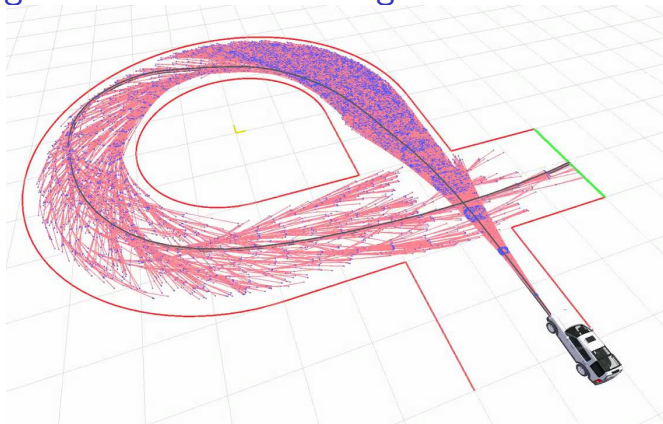
- ▶ CMU's autonomous car used search-based motion planning in the DARPA Urban Challenge in 2007
- ▶ Likhachev and Ferguson, "Planning Long Dynamically Feasible Maneuvers for Autonomous Vehicles," IJRR'09
- ▶ Video: <https://www.youtube.com/watch?v=4hFh100i8KI>
- ▶ Video: <https://www.youtube.com/watch?v=qXZt-B7iUyw>
- ▶ Paper: <http://journals.sagepub.com/doi/pdf/10.1177/0278364909340445>

# Sampling-based Motion Planning



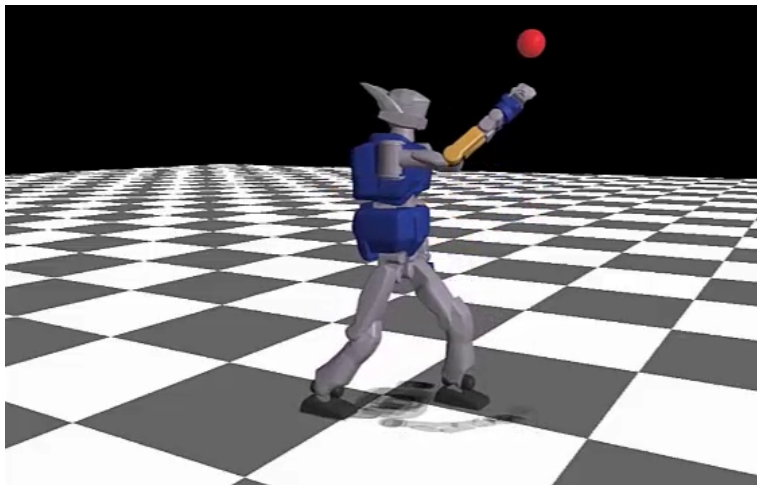
- ▶ RRT algorithm on the PR2 – planning with both arms (12 DOF)
- ▶ Karaman and Frazzoli, “Sampling-based algorithms for optimal motion planning,” IJRR’11
- ▶ Video: <https://www.youtube.com/watch?v=vW74bC-Ygb4>
- ▶ Paper: <http://journals.sagepub.com/doi/pdf/10.1177/0278364911406761>

# Sampling-based Motion Planning



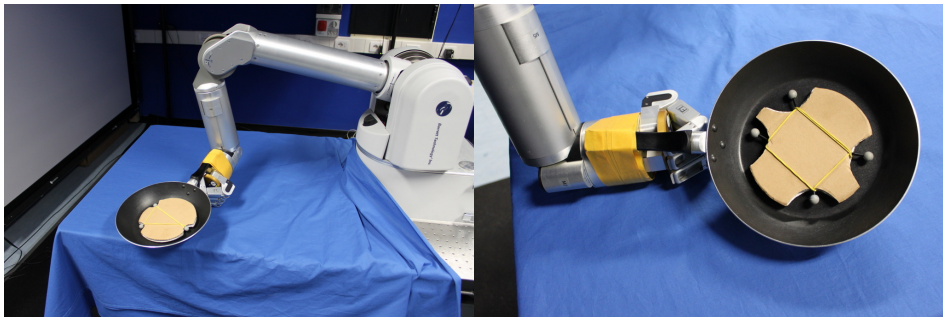
- ▶ RRT\* algorithm on a high-fidelity car model – 270 degree turn
- ▶ Karaman and Frazzoli, “Sampling-based algorithms for optimal motion planning,” IJRR’11
- ▶ Video: <https://www.youtube.com/watch?v=p3nZHn0Whrg>
- ▶ Video: <https://www.youtube.com/watch?v=LKL5qRBiJaM>
- ▶ Paper: <http://journals.sagepub.com/doi/pdf/10.1177/0278364911406761>

## Dynamic Programming and Optimal Control



- ▶ Tassa, Mansard and Todorov, "Control-limited Differential Dynamic Programming," ICRA'14
- ▶ Video: <https://www.youtube.com/watch?v=tCQSSkBH2NI>
- ▶ Paper: <http://ieeexplore.ieee.org/document/6907001/>

# Model-free Reinforcement Learning



- ▶ A robot learns to flip pancakes
- ▶ Kormushev, Calinon and Caldwell, "Robot Motor Skill Coordination with EM-based Reinforcement Learning," IROS'10
- ▶ Video: [https://www.youtube.com/watch?v=W\\_gxLKSsSIE](https://www.youtube.com/watch?v=W_gxLKSsSIE)
- ▶ Paper: <http://www.dx.doi.org/10.1109/IROS.2010.5649089>

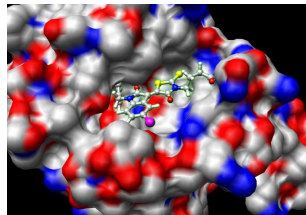
# Applications of Optimal Control & Reinforcement Learning



(a) Autonomous Driving



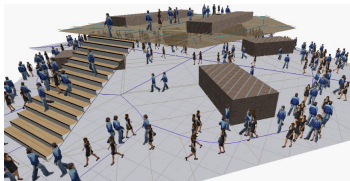
(b) Marketing



(c) Computational Biology



(d) Games



(e) Character Animation



(f) Robotics

# Models

- ▶ **Motion model:** specifies how a dynamical system evolves

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t, \mathbf{w}_t) \sim p_f(\cdot \mid \mathbf{x}_t, \mathbf{u}_t), \quad t = 0, \dots, T - 1$$

- ▶ discrete time  $t \in \{0, \dots, T\}$
  - ▶ state  $\mathbf{x}_t \in \mathcal{X}$
  - ▶ control  $\mathbf{u}_t \in \mathcal{U}(\mathbf{x}_t)$  and  $\mathcal{U} := \bigcup_{\mathbf{x} \in \mathcal{X}} \mathcal{U}(\mathbf{x})$
  - ▶ motion noise  $\mathbf{w}_t$ : random vector with known probability density function (pdf) and assumed conditionally independent of other disturbances  $\mathbf{w}_\tau$  for  $\tau \neq t$  for given  $\mathbf{x}_t$  and  $\mathbf{u}_t$
  - ▶ the motion model is specified by the nonlinear function  $f$  or equivalently by the pdf  $p_f$  of  $\mathbf{x}_{t+1}$  conditioned on  $\mathbf{x}_t$  and  $\mathbf{u}_t$
- ▶ **Observation model:** the state  $\mathbf{x}_t$  might not be observable but perceived through measurements:

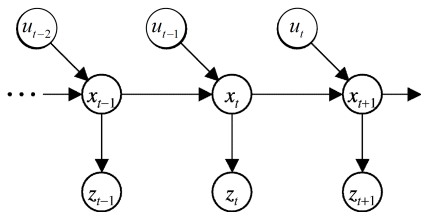
$$\mathbf{z}_t = h(\mathbf{x}_t, \mathbf{v}_t) \sim p_h(\cdot \mid \mathbf{x}_t), \quad t = 0, \dots, T$$

- ▶ measurement noise  $\mathbf{v}_t$ : random vector with known pdf and conditionally independent of other disturbances  $\mathbf{v}_\tau$  for  $\tau \neq t$  and  $\mathbf{w}_t$  for all  $t$  for given  $\mathbf{x}_t$
- ▶ the observation model is specified by the nonlinear function  $h$  or equivalently by the pdf  $p_h$  of  $\mathbf{z}_t$  conditioned on  $\mathbf{x}_t$

# Problem Structure

## ▶ Markov Assumptions

- ▶ The state  $\mathbf{x}_{t+1}$  only depends on the previous input  $\mathbf{u}_t$  and state  $\mathbf{x}_t$
- ▶ The observation  $\mathbf{z}_t$  only depends on the state  $\mathbf{x}_t$



- ▶ **Problem structure:** due to the Markov assumptions, the joint distribution of the robot states  $\mathbf{x}_{0:T}$ , observations  $\mathbf{z}_{0:T}$ , and controls  $\mathbf{u}_{0:T-1}$  satisfies:

$$p(\mathbf{x}_{0:T}, \mathbf{z}_{0:T}, \mathbf{u}_{0:T-1}) =$$

$$\underbrace{p_0(\mathbf{x}_0)}_{\text{prior}} \prod_{t=0}^T \underbrace{p_h(\mathbf{z}_t | \mathbf{x}_t)}_{\text{observation model}} \prod_{t=1}^T \underbrace{p_f(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{u}_{t-1})}_{\text{motion model}} \prod_{t=0}^{T-1} \underbrace{p(\mathbf{u}_t | \mathbf{x}_t)}_{\text{control policy}}$$



## Problem Statement

- ▶ In general, the states  $\mathbf{x}_t$  are **partially observable** via the noisy observations  $\mathbf{z}_t$
- ▶ A partially observable problem can always be converted to a fully observed one by changing the state from  $\mathbf{x}_t$  to the probability density function  $p_{t|t}(\mathbf{x}_t) := p(\mathbf{x}_t \mid \mathbf{z}_{0:t}, \mathbf{u}_{0:t-1})$
- ▶ Without loss of generality, consider **fully observable** states  $\mathbf{x}_t$
- ▶ Given an initial state  $\mathbf{x}_0$ , determine control inputs  $\mathbf{u}_{0:T-1}$  that minimize (maximize) the expected long-term cost (reward) along the state trajectory  $\mathbf{x}_{1:T}$  determined by the motion model  $p_f$ :

$$V_0^{\mathbf{u}_{0:T-1}}(\mathbf{x}_0) := \mathbb{E}_{\mathbf{x}_{1:T}} \left[ \underbrace{q(\mathbf{x}_T)}_{\text{terminal cost}} + \sum_{t=0}^{T-1} \underbrace{\ell(\mathbf{x}_t, \mathbf{u}_t)}_{\text{stage cost}} \mid \mathbf{x}_0, \mathbf{u}_{0:T-1} \right]$$

## Problem Solution

- ▶ The solution to an optimal control problem is a **policy**  $\pi$
- ▶ Let  $\pi_t$  be a **function** that maps a state  $\mathbf{x}_t \in \mathcal{X}$  to a feasible control input  $\mathbf{u}_t \in \mathcal{U}(\mathbf{x}_t)$
- ▶ An **admissible control policy** is a sequence of functions  $\pi_{0:T-1} := \{\pi_0, \pi_1, \dots, \pi_{T-1}\}$
- ▶ To simplify notation, we informally denote  $\pi_{0:T-1}$  by  $\pi$
- ▶ The expected long-term cost (reward)  $V_t^\pi(\mathbf{x})$  of a policy  $\pi$  starting at time  $t$  at state  $\mathbf{x}$  is called the **value function** of  $\pi$ :

$$V_t^\pi(\mathbf{x}) := \mathbb{E}_{\mathbf{x}_{t+1:T}} \left[ q(\mathbf{x}_T) + \sum_{\tau=t}^{T-1} \ell(\mathbf{x}_\tau, \pi_\tau(\mathbf{x}_\tau)) \mid \mathbf{x}_t = \mathbf{x} \right]$$

- ▶ A policy  $\pi^*$  is **optimal** if  $V_0^{\pi^*}(\mathbf{x}) \leq V_0^\pi(\mathbf{x})$  for all admissible  $\pi$  and all  $\mathbf{x}$  and its value function is denoted  $V_0^*(\mathbf{x}) := V_0^{\pi^*}(\mathbf{x})$

# Naming Conventions

- ▶ The problem of acting optimally is called:
  - ▶ **Optimal Control** (OC): when the models  $p_f$ ,  $p_h$  and cost functions  $\ell$ ,  $q$  are known
  - ▶ **Reinforcement Learning** (RL): when the models  $p_f$ ,  $p_h$  and cost functions  $\ell$ ,  $q$  are unknown but samples  $\mathbf{x}_t$ ,  $\mathbf{z}_t$ ,  $\ell(\mathbf{x}_t, \mathbf{u}_t)$ ,  $q(\mathbf{x}_t)$  can be obtained from them
- ▶ Conventions differ in optimal control and reinforcement learning:
  - ▶ **OC**: minimization, cost, state  $\mathbf{x}$ , control  $\mathbf{u}$ , policy  $\mu$
  - ▶ **RL**: maximization, reward, state  $\mathbf{s}$ , action  $\mathbf{a}$ , policy  $\pi$
  - ▶ **ECE276B**: minimization, cost, state  $\mathbf{x}$ , control  $\mathbf{u}$ , policy  $\pi$

## Policy Types

- ▶ Controls may have long-term consequences, e.g., delayed cost/reward
- ▶ It may be better to sacrifice immediate rewards to gain long-term rewards:
  - ▶ A financial investment may take months to mature
  - ▶ Re-fueling a helicopter now might prevent a crash in several hours
  - ▶ Blocking an opponent move now might help winning chances many moves from now
- ▶ A policy fully defines, at any given point in time  $t$  and any given state  $\mathbf{x}_t$ , which control  $\mathbf{u}_t$  to apply.
- ▶ A policy can be:
  - ▶ **stationary** ( $\pi \equiv \pi_0 \equiv \pi_1 \equiv \dots$ )  $\subset$  **non-stationary** (time-dependent)
  - ▶ **deterministic** ( $\mathbf{u}_t = \pi_t(\mathbf{x}_t)$ )  $\subset$  **stochastic** ( $\mathbf{u}_t \sim \pi_t(\cdot | \mathbf{x}_t)$ )
  - ▶ **open-loop** (a sequence  $\mathbf{u}_{0:T-1}$  regardless of  $\mathbf{x}_t$ )  $\subset$  **closed-loop** ( $\pi_t$  depends on  $\mathbf{x}_t$ )

# Problem Variations

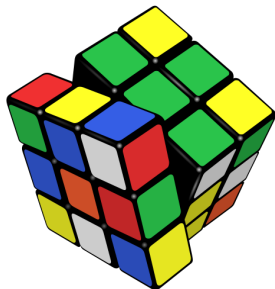
- ▶ **deterministic** (no noise) vs **stochastic**
- ▶ **fully observable** ( $\mathbf{z}_t = \mathbf{x}_t$ ) vs **partially observable** ( $\mathbf{z}_t \sim p_h(\cdot|\mathbf{x}_t)$ )
  - ▶ Markov Decision Process (MDP) vs Partially Observable Markov Decision Process (POMDP)
- ▶ **stationary** vs **nonstationary** (time-dependent  $p_{f,t}, p_{h,t}, \ell_t$ )
- ▶ **discrete** vs **continuous** state space  $\mathcal{X}$ 
  - ▶ tabular approach vs function approximation (linear, SVM, neural nets,...)
- ▶ **discrete** vs **continuous** control space  $\mathcal{U}$ :
  - ▶ tabular approach vs optimization problem to select next-best control
- ▶ **discrete** vs **continuous** time:
  - ▶ finite-horizon discrete time: dynamic programming
  - ▶ infinite-horizon discrete time: Bellman equation (first-exit vs discounted vs average-reward formulation)
  - ▶ continuous time: Hamilton-Jacobi-Bellman (HJB) Partial Differential Equation (PDE)
- ▶ reinforcement learning ( $p_f, p_h, \ell, q$  are unknown):
  - ▶ **Model-based RL**: explicitly approximate the models  $\hat{p}_f, \hat{p}_h, \hat{\ell}, \hat{q}$  from data and apply optimal control algorithms
  - ▶ **Model-free RL**: directly approximate  $V_t^*$  and  $\pi_t^*$  without approximating the motion, observation, or cost models

## Example: Inventory Control

- ▶ Consider keeping an item stocked in a warehouse:
  - ▶ If there is too little, we may run out (not preferred).
  - ▶ If there is too much, the storage cost will be high (not preferred).
- ▶ This scenario can be modeled as a discrete-time system:
  - ▶  $x_t \in \mathbb{R}$ : stock available in the warehouse at the beginning of the  $t$ -th time period
  - ▶  $u_t \in \mathbb{R}_{\geq 0}$ : stock ordered and immediately delivered at the beginning of the  $t$ -th time period (supply)
  - ▶  $w_t$ : random demand during the  $t$ -th time period with known pdf. Note that excess demand is back-logged, i.e., corresponds to negative stock  $x_t$
  - ▶ **Motion model:**  $x_{t+1} = x_t + u_t - w_t$
  - ▶ **Cost function:**  $\mathbb{E} \left[ R(x_T) + \sum_{t=0}^{T-1} (r(x_t) + cu_t - pw_t) \right]$  where
    - ▶  $pw_t$ : revenue
    - ▶  $cu_t$ : cost of items
    - ▶  $r(x_t)$ : penalizes too much stock or negative stock
    - ▶  $R(x_T)$ : remaining items we cannot sell or demand that we cannot meet

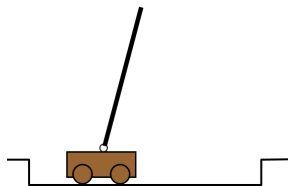
## Example: Rubik's Cube

- ▶ Invented in 1974 by Ernő Rubik
- ▶ Formalization:
  - ▶ State space:  $\sim 4.33 \times 10^{19}$
  - ▶ Actions: 12
  - ▶ Reward:  $-1$  for each time step
  - ▶ Deterministic, Fully Observable
- ▶ The cube can be solved in 20 or fewer moves



## Example: Cart-Pole Problem

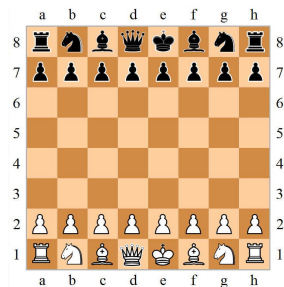
- ▶ Move a cart left and right in order to keep a pole balanced
- ▶ Formalization:
  - ▶ State space: 4-D continuous  $(x, \dot{x}, \theta, \dot{\theta})$
  - ▶ Actions:  $\{-N, N\}$
  - ▶ Reward:
    - ▶ 0 when in the goal region
    - ▶  $-1$  when outside the goal region
    - ▶  $-100$  when outside the feasible region
  - ▶ Deterministic, Fully Observable





# Example: Chess

- ▶ Formalization:
  - ▶ State space:  $\sim 10^{47}$
  - ▶ Actions: from 0 to 218
  - ▶ Reward: 0 each step,  $\{-1, 0, 1\}$  at the end of the game
  - ▶ Deterministic, opponent-dependent state transitions (can be modeled as a game)
- ▶ The size of the game tree is  $10^{123}$



## Example: Grid World Navigation

- ▶ Navigate to a goal without crashing into obstacles
- ▶ Formalization:
  - ▶ State space: robot pose, e.g., 2-D position
  - ▶ Actions: allowable robot movement, e.g.,  $\{left, right, up, down\}$
  - ▶ Reward:  $-1$  until the goal is reached;  $-\infty$  if an obstacle is hit
  - ▶ Can be deterministic or stochastic; fully or partially observable

