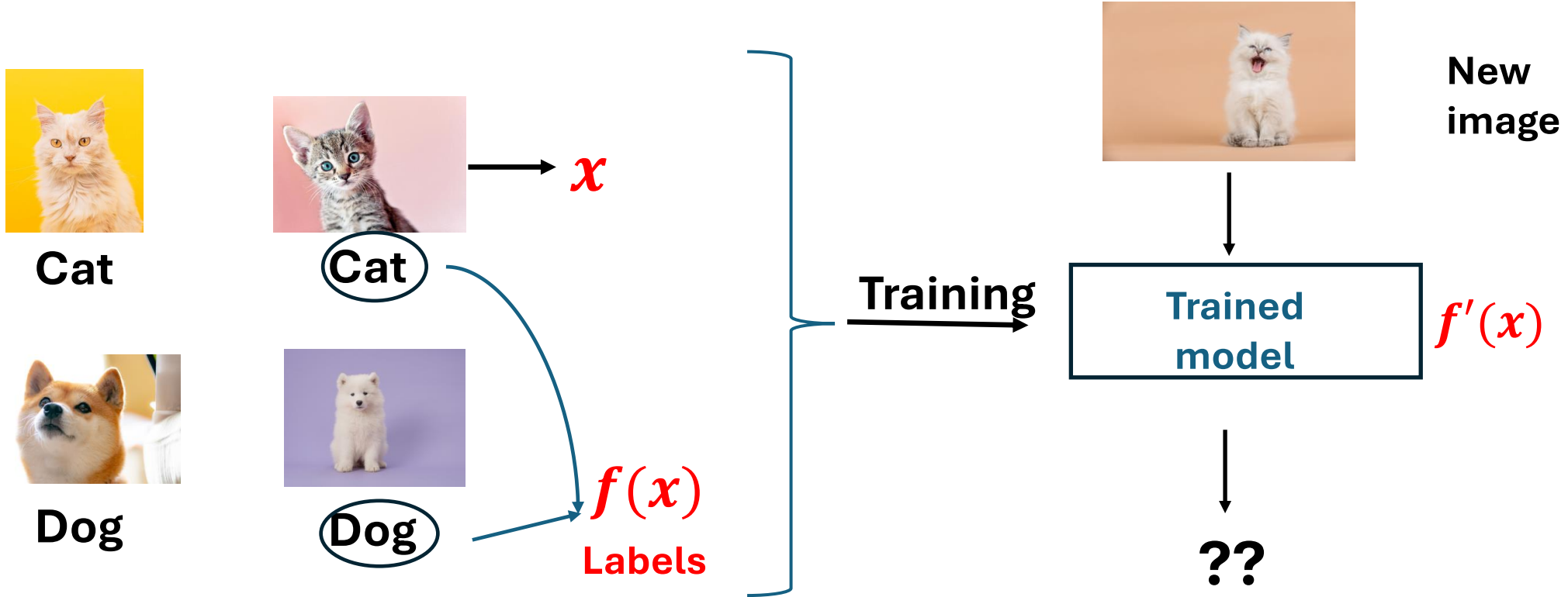# Introduction to Reinforcement Learning

Subrahmanya Swamy Peruru

# Paradigms of Machine Learning

- Supervised Learning

- Unsupervised Learning

- Reinforcement Learning
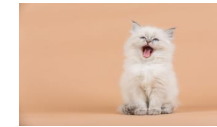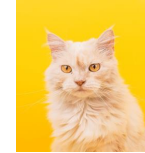
# Supervised Learning

**Labeled Training Data**



Cat

Cat → $x$

Dog

Dog → $f(x)$

Labels

Training

New image

Trained model  $f'(x)$

??

Photos provided by Pexels

# Unsupervised Learning

**Unlabeled Data**

Identify patterns
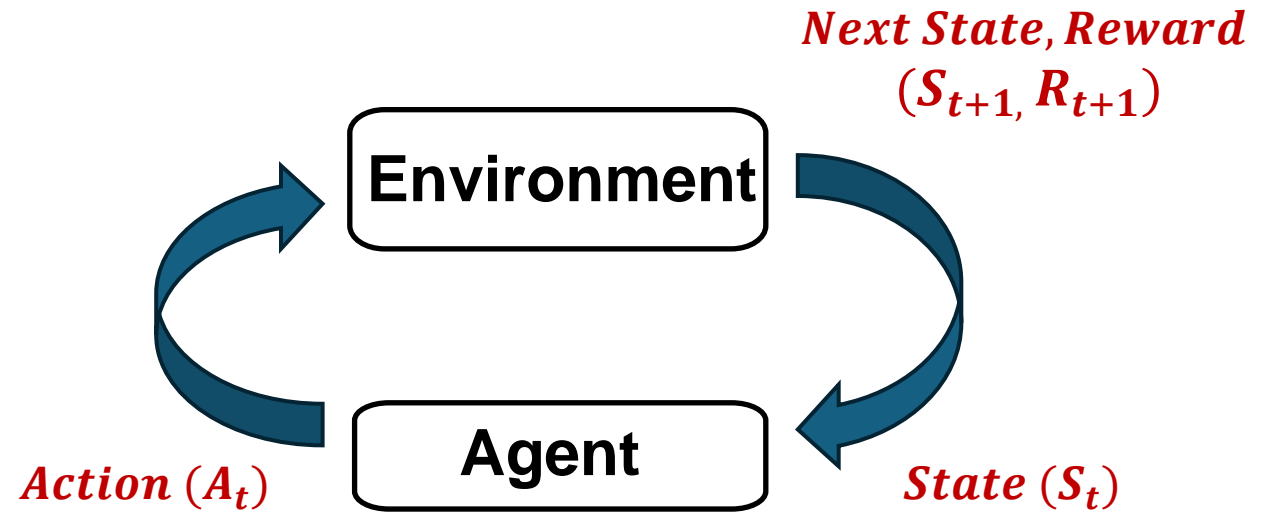
# Reinforcement Learning

**Feedback:**
**Score,**
**new display**

**State:**
**Display**

**Actions:**
**UP / Left /**
**Right / Down**

**Learn by Trial and Error**

*Next State, Reward*
$(S_{t+1,}\ R_{t+1})$

**Environment**

**Agent**

*Action* $(A_t)$

*State* $(S_t)$

1. Agent observes the state and takes action
2. Environment puts the agent in a new state &
3. Gives a reward based on the action taken

**GOAL:** **Learn policy to maximize**
**the cumulative reward**

$$\sum_t R\_t$$

# Paradigms of Machine Learning

- ## Supervised Learning
    - Fitting a function for the given labeled data $(x, y)$
    - $y \approx f(x)$

- ## Unsupervised Learning
    - Identifying patterns in unlabeled data
    - E.g. Clustering
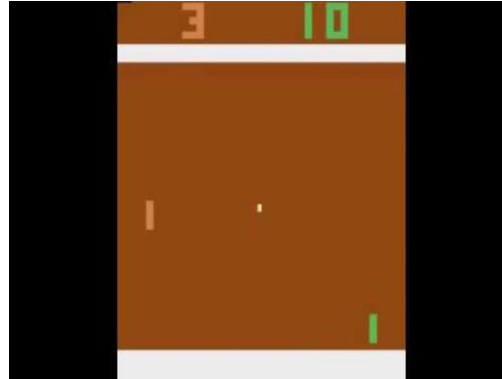
- ## Reinforcement Learning
    - Learning sequential tasks through trial and error
    - Feedback through reward/penalty

# RL Demonstrations

**Autonomous Helicopter**

**Pong game**

**AlphaGo by DeepMind**

# One State RL: Multi-arm Bandits

- Simplified version of RL problem: "Multi-arm Bandit" problem.
  - Only one state
  - Multiple actions (a.k.a. arms)
    - $\mathcal{A}$ – Action set

- A reward distribution corresponding to each arm
  - $\mathcal{R}_a$ – Reward distrution for action a
  - $\mu_a = \mathbb{E}[\mathcal{R}_a]$ – Expected reward for action $a$

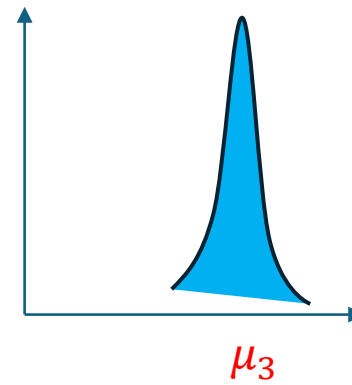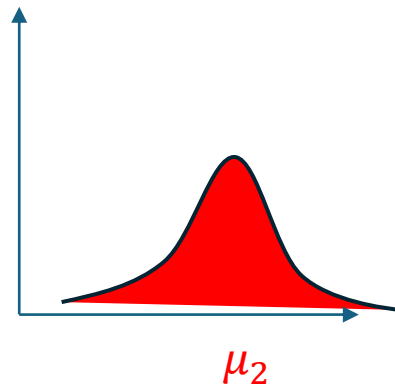- Applications: Recommendation systems, Ad placement, …

# Multi-arm Bandits



Arm 1

Arm 2

Arm 3

Reward Distributions
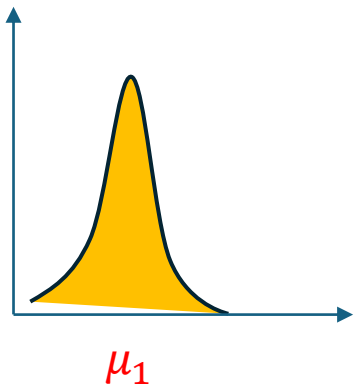
$\mu_1$

$\mu_2$

$\mu_3$

**Problem:**
- Reward distributions are **unknown**

- Given **T chances** to pull the arms

- Which arms should be pulled to **maximize the total reward** in those T rounds

Exploration
Vs
Exploitation dilemma

# ETC (Explore-Then-Commit)

1. **Explore:** Play each arm $N$ times

2. Compute the sample average rewards $\bar{\mu}(a) = \frac{1}{N}\sum_{t=1}^{KN} R_t\, 1\{a_t = a\}$ for each arm $a \in \mathcal{A}$

3. **Commit:** Play the arm with the highest sample average for the remaining $T - KN$ rounds

$\mu^*$ - Optimal arm's expected reward    $R_t$ - Sample reward obtained in round $t$

$a_t$ - Arm played in round $t$    $T$   - Total number of rounds

$K$   - Number of arms

**Performance (ETC Vs Best possible reward) :** $T\mu^* - \sum_{t=1}^{T} \mathbb{E}[R_t]$

**How much to Explore?** $N \approx (\frac{T}{K})^{\frac{2}{3}}$

# $\epsilon$-Greedy (Explore uniformly)

1. Play each arm once

2. In each round $t$:
   - Toss a coin with bias $\epsilon$.

   - If it lands in head: Explore - Play any arm randomly

   - Else: Exploit - Play the arm with the highest sample average so far

**What $\epsilon$ to choose?** $\epsilon \approx \left(\dfrac{K}{T}\right)^{\frac{1}{3}}$

# UCB (Upper Confidence Bound)

## Optimism under UnCertainty

1. Play each arm once in the first $K$ rounds

2. For $t > K$:

    - Play the arm with the highest $UCB_t(a) = \overline{\mu_{t-1}}(a) + \sqrt{\dfrac{2 \log T}{n_{t-1}(a)}}$

      **Exploit** ↑ (points to $\overline{\mu_{t-1}}(a)$)

      **Explore** ↑ (points to $\sqrt{\dfrac{2 \log T}{n_{t-1}(a)}}$)

    - Based on the observed sample reward $R_t$, update $n_t(a_t)$ and $\overline{\mu}_t(a_t)$

        - $n_t(a_t) = n_{t-1}(a_t) + 1$
        - $\overline{\mu}_t(a_t) = \dfrac{1}{n_t(a_t)}\left[(n_t(a_t) - 1)\,\overline{\mu_{t-1}}(a_t) + R_t\right]$

**Exploit:** High sample reward arms are favoured
**Explore:** Least played arms are favoured

# Contextual Bandits – Multiple states

- ## News article Recommendation systems



- Articles – arms

- Like / Dislike – Reward

- User – State

Different users have different preferences to articles