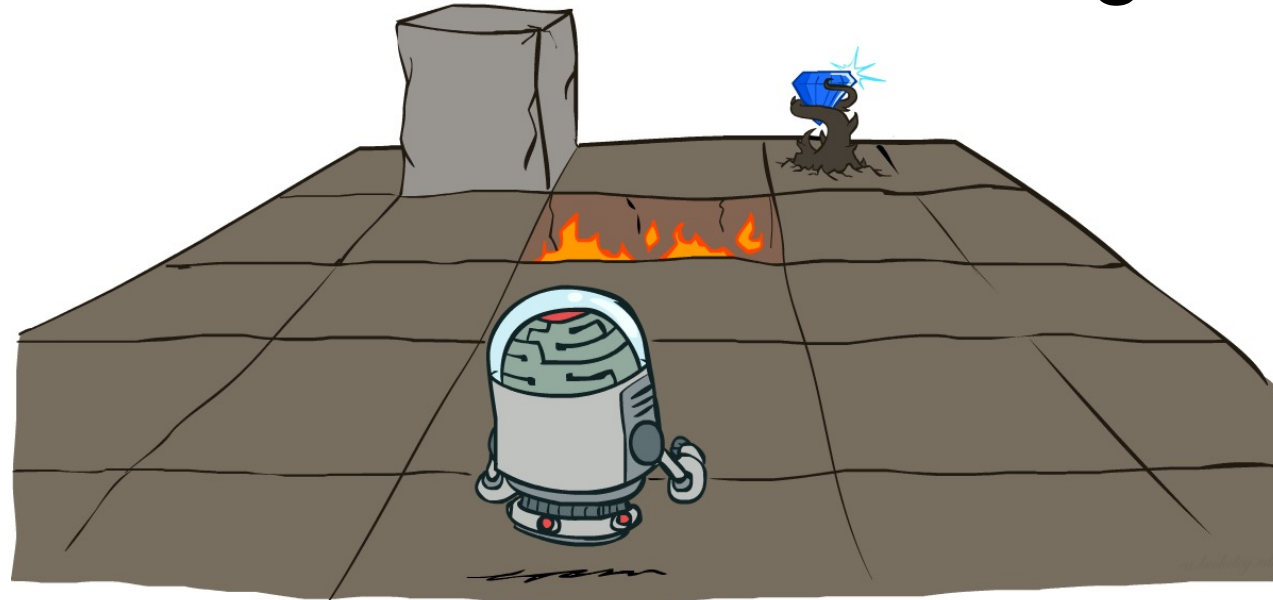# CS 3568: Intelligent Systems

## Reinforcement Learning
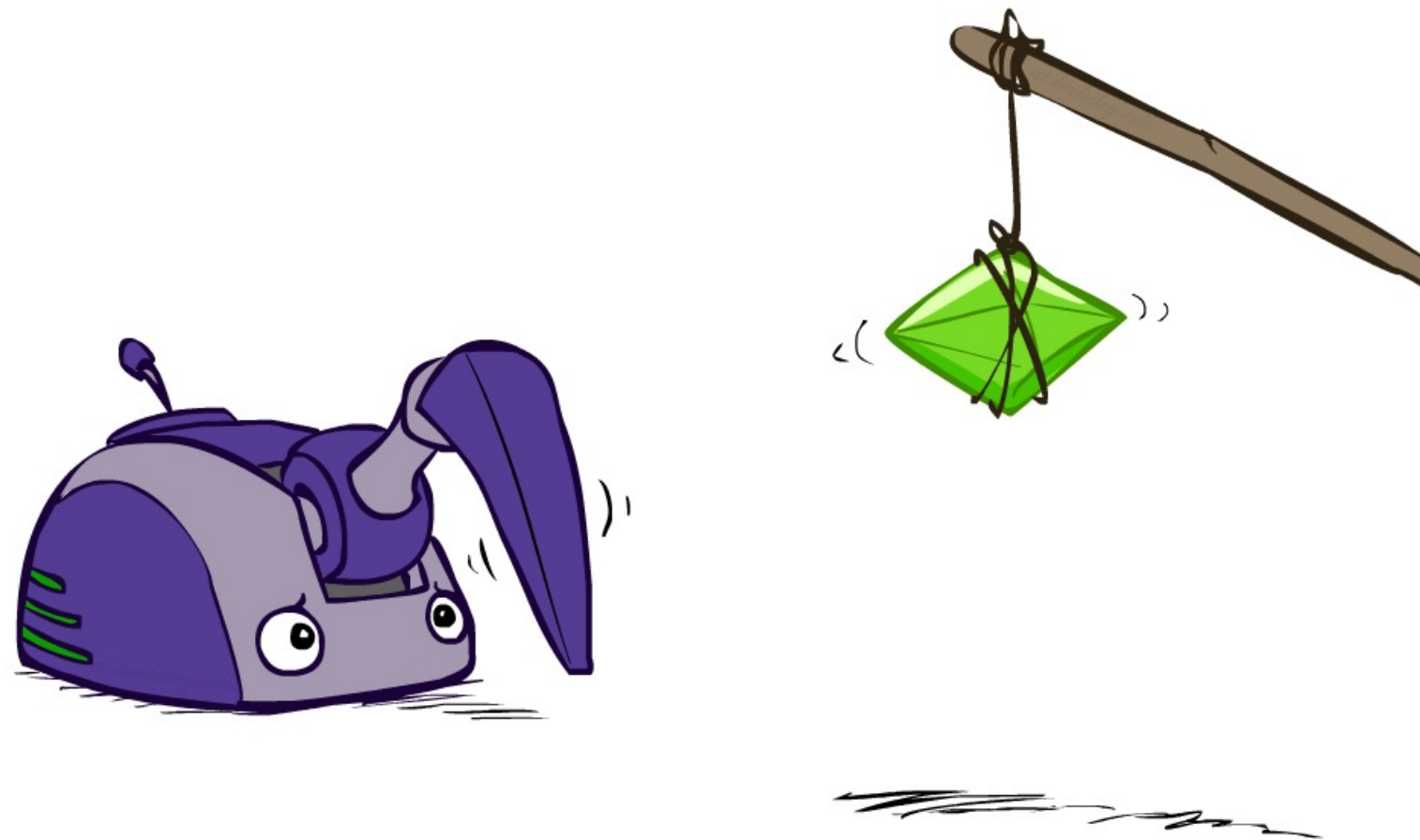
Instructor: Tara Salman

Texas Tech University

Computer Science Department

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley (ai.berkeley.edu).]
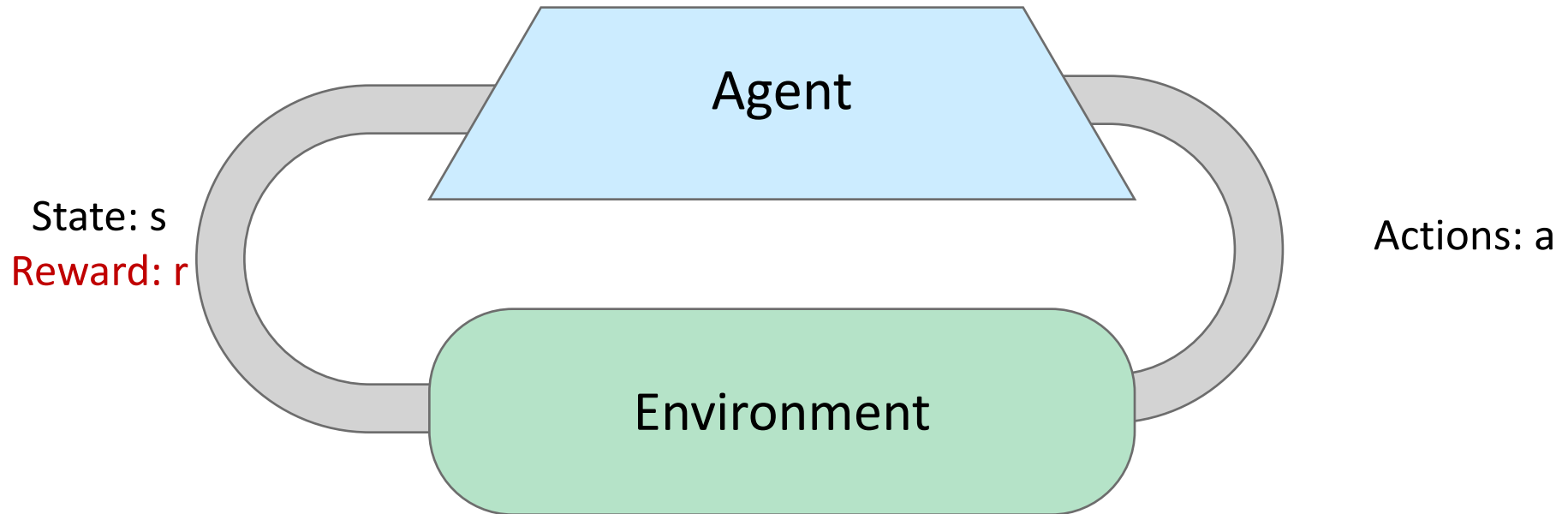
# Reinforcement Learning

Texas Tech University

Tara Salman

# Reinforcement Learning

Agent

State: s
Reward: r

Actions: a

Environment

❑ Basic idea:
  ➢ Receive feedback in the form of rewards
  ➢ Agent's utility is defined by the reward function
  ➢ Must (learn to) act so as to maximize expected rewards
  ➢ All learning is based on observed samples of outcomes!

Texas Tech University

Tara Salman

# Example: Learning to Walk



Initial

A Learning Trial

After Learning [1K Trials]

[Kohl and Stone, ICRA 2004]

# Example: Learning to Walk



Initial

[Kohl and Stone, ICRA 2004]

Texas Tech University

Tara Salman

# Example: Learning to Walk

Training

6

Texas Tech University

Tara Salman

# Example: Learning to Walk



Finished

[Kohl and Stone, ICRA 2004]

Texas Tech University

Tara Salman

# Example: Toddler Robot

Texas Tech University

Tara Salman

# The Crawler!



[You, in Project 3]

Texas Tech University

Tara Salman

# Video of Demo Crawler Bot

Texas Tech University

Tara Salman

# Reinforcement Learning

❑ Still assume a Markov decision process (MDP):

  ➢ A set of states s ∈ S

  ➢ A set of actions (per state) A
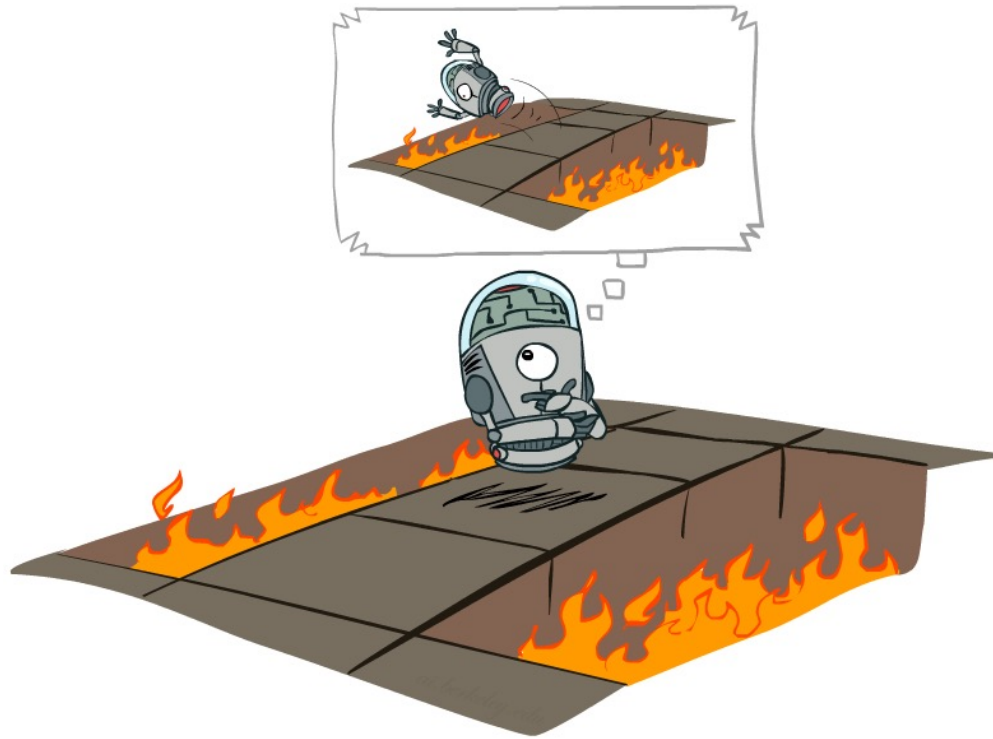
  ➢ A model T(s,a,s')

  ➢ A reward function R(s,a,s')

❑ Still looking for a policy π(s)

❑ New twist: don't know T or R

  ➢ I.e. we don't know which states are good or what the actions do

  ➢ Must actually try out actions and states to learn

Warm

Cool

Overheated

Texas Tech University

Tara Salman

# Offline (MDPs) vs. Online (RL)



Offline Solution

Online Learning

Texas Tech University

Tara Salman

# Model-Based Learning

Texas Tech University

Tara Salman

# Model-Based Learning

❑ Model-Based Idea:
  ➢ Learn an approximate model based on experiences
  ➢ Solve for values as if the learned model were correct

❑ Step 1: Learn empirical MDP model
  ➢ Count outcomes s' for each s, a
  ➢ Normalize to give an estimate of $\widehat{T}(s, a, s')$
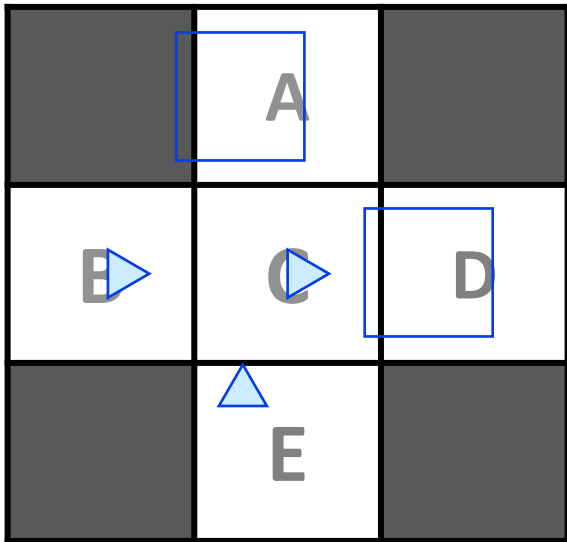  ➢ Discover each $\widehat{R}(s, a, s')$ when we experience (s, a, s')

❑ Step 2: Solve the learned MDP
  ➢ For example, use value iteration, as before

Texas Tech University

Tara Salman

# Example: Model-Based Learning

## Input Policy π



*Assume: γ = 1*

## Observed Episodes (Training)

### Episode 1

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 2

B, east, C, -1
C, east, D, -1
D, exit,  x, +10

### Episode 3

E, north, C, -1
C, east,   D, -1
D, exit,    x, +10

### Episode 4

E, north, C, -1
C, east,   A, -1
A, exit,    x, -10

## Learned Model

$\hat{T}(s, a, s')$

T(B, east, C) = 1.00
T(C, east, D) = 0.75
T(C, east, A) = 0.25
...

$\hat{R}(s, a, s')$

R(B, east, C) = -1
R(C, east, D) = -1
R(D, exit, x) = +10
...

Texas Tech University

Tara Salman

# Example: Expected Age

Goal: Compute expected age of cs188 students

**Known P(A)**

$$E[A] = \sum_a P(a) \cdot a \qquad = 0.35 \times 20 + \dots$$

Without P(A), instead collect samples $[a_1, a_2, \dots a_N]$

**Unknown P(A): "Model Based"**

$$\hat{P}(a) = \frac{\text{num}(a)}{N}$$

$$E[A] \approx \sum_a \hat{P}(a) \cdot a$$

Why does this work?  Because eventually you learn the right model.

**Unknown P(A): "Model Free"**

$$E[A] \approx \frac{1}{N} \sum_i a_i$$

Why does this work?  Because samples appear with the right frequencies.

Texas Tech University

Tara Salman