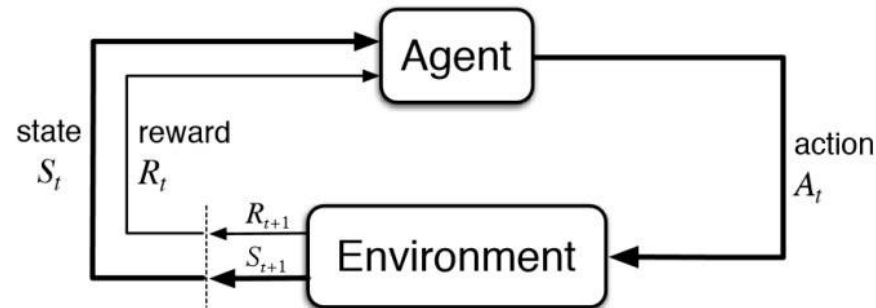# Inverse RL/Optimal Control w/ Energy based model (EBM)

Imola Fodor

- Reinforcement Learning (RL) high-level and challenges
- Inverse Reinforcement Learning (IRL) idea
- IRL: Maximum Entropy → Guided Cost Learning
- Electrolux context

# The Reinforcement Learning system



- Markov Decision Process is a tuple <S, A, p, R, gamma> , policy π
- Goal/return – maximize discounted cumulative reward

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+(k+1)}$$

- Value function – gives long term value of state s

$$V(s) = E[G_t | S_t = s]$$

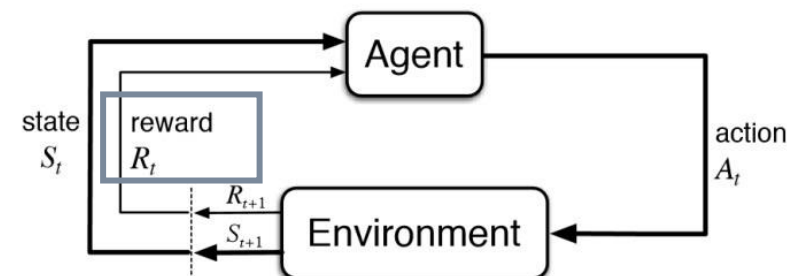- Overall Objective, Expected Reward: $argmax_\pi [E(G_t)]$
- Recursion relationship

$$V_\pi (s) = E[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s]$$

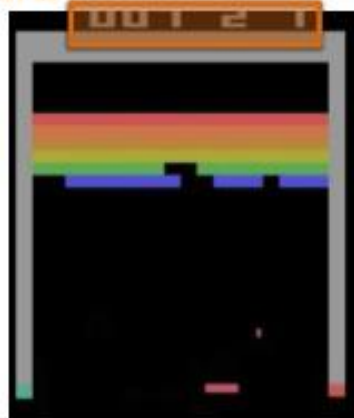- Reward function has a crucial role to quantify how good the agent is.

# Reward function

- Reward function has a crucial role to quantify how good the agent is.
- Can be hand-crafted and learnt



state $S_t$   reward $R_t$   Agent   action $A_t$

$R_{t+1}$

$S_{t+1}$ Environment

**Computer Games**

reward

**Real World Scenarios**

robotics    dialog    autonomous driving

Mnih et al. '15

what is the reward?

# Inverse RL

## Forward RL

Given:

states s from S, actions a from A

(sometimes) transitions $p(s'|s, a)$

reward function $r(s, a)$

**learn** $\pi * (a|s)$

## Inverse RL

Given:

expert state, action pairs

(sometimes) transitions $p(s'|s, a)$
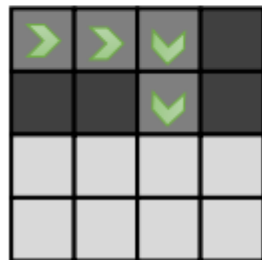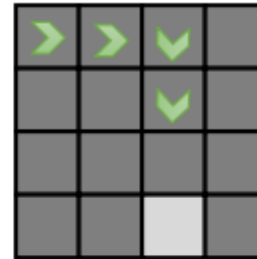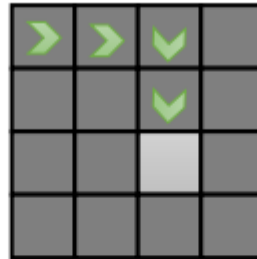
**learn:** $r_\phi(s, a)$

..and then learn $\pi *(a|s)$

# Ambiguity

- IRL is essentially an ill-posed (unspecified) problem because multiple reward functions can explain the same expert's behavior
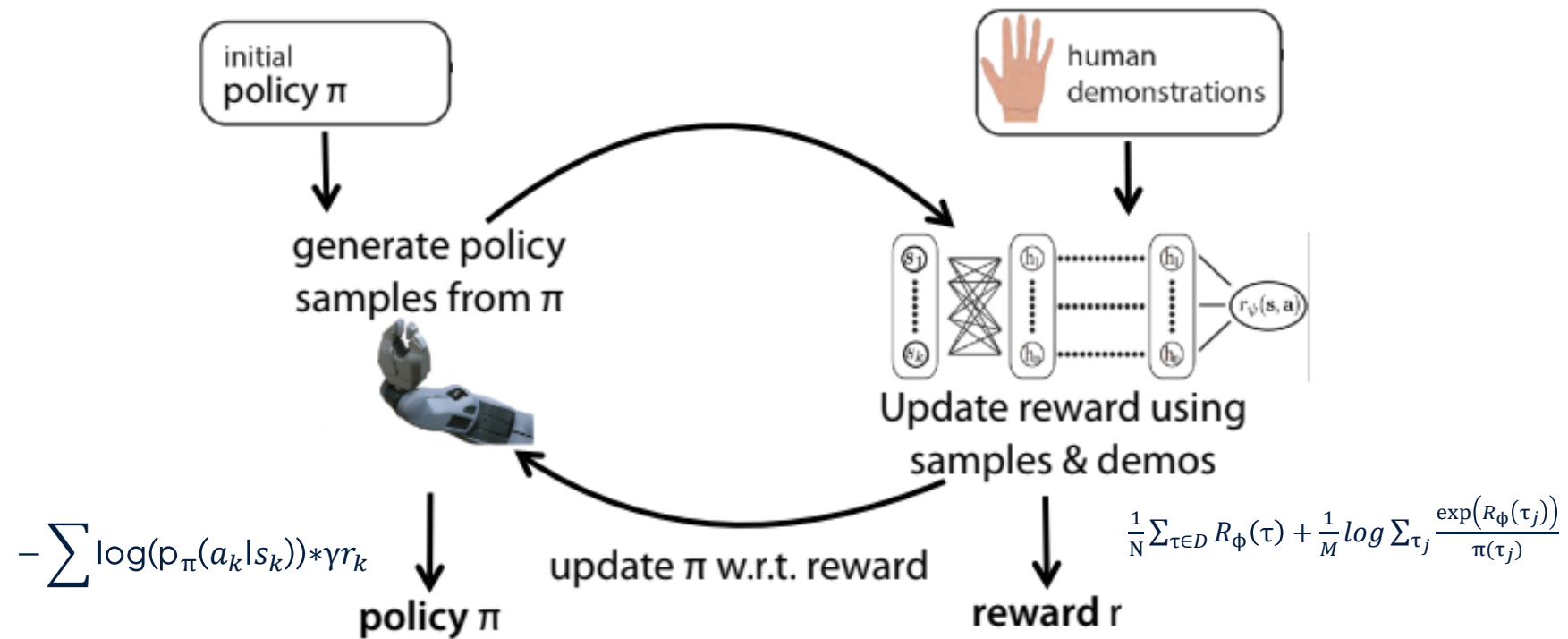
# Energy Based Models

- Originally models the physical behaviour of interacting particles – physics
- Class of generative machine learning models
- Governed by an energy function that describes the probability of a certain state
- Guaranteed to sample a valid probability distribution (learn to sample)

- What probability distribution $p(x)$ can be assigned to a state $x$ with energy $E(x)$?

- Boltzmann distribution
    - The distribution shows that states with lower energy will always have a higher probability of being occupied.
- Through maximum entropy..

# guided cost learning algorithm
## (Finn et al. ICML '16)

initial **policy π**

human demonstrations

generate policy samples from π

Update reward using samples & demos

update π w.r.t. reward

$-\sum \log(p_\pi(a_k|s_k))*\gamma r_k$

**policy π**

**reward r**

$\frac{1}{N}\Sigma_{\tau\in D}\,R_\phi(\tau) + \frac{1}{M}log\,\Sigma_{\tau_j}\frac{\exp\left(R_\phi(\tau_j)\right)}{\pi(\tau_j)}$

*Chelsea Finn*
*Deep RL Bootcamp*

# Maximum entropy

- The principle of maximum entropy states that the probability distribution which best represents the current state of knowledge about a system is the one with largest entropy, where the uncertainty is maximal.

- The formulas show the connection between energy, maximum entropy and Boltzmann distribution

$$\max_{P(x)} \sum_x -P(x) \log P(x)$$

$$\text{s.t.} \sum_x P(x)E(x) = \langle E \rangle,$$

giving [22]

$$P(x) = \frac{1}{Z} \exp\left[-E(x)/T\right],$$

- The Boltzmann distribution establishes a concrete relationship between energy and probability: low-energy states are the most likely to be observed.

- In low energy states are observed, all other in high temperatures are equally likely

# IRL MaxEntropy – for known dynamics

- Boltzmann data distribution (Energy based model):

$$p(\tau) = \frac{1}{Z} \exp(R_\phi(\tau))$$

where $\tau = \{x_1, a_1, \ldots, x_T, a_T\}$ is a trajectory sample, $R(\tau) = \sum_t r_\phi(s_t, a_t)$ is an unknown reward function parameterized by $\phi$, and $x_t$ and $a_t$ are the state and action at time step t. Z partition function

- For inferring the reward function we need:

$$max_\phi \, L(\phi) = \sum_{\tau \in D} log p_{r\psi}(\tau) = N \sum_{\tau \in D} r_\phi(\tau) - M log Z = N \sum_{\tau \in D} r_\phi(\tau) - M log \sum_\tau \exp\left(r_\phi(\tau)\right)$$

- To optimize, taking it's derivative

# Guided Cost Learning GCL – Cost

- Not learning the optimal policy for given reward
- Not calculating partition function Z
- Importance sampling ideally
- Adaptive sampling (from a created policy)

$$\frac{1}{N}\sum_{\tau \in D} R_\phi(\tau) - \frac{1}{M} log \sum_{\tau_j} \frac{\exp\left(R_\phi(\tau_j)\right)}{\pi(\tau_j)}$$

$$\nabla_\phi L \approx \frac{1}{N}\sum_{i=1}^{N} \nabla_\phi r_\phi(\tau_i) - \frac{1}{\sum_j w_j} \sum_{j=1}^{M} w_j \nabla_\phi r_\phi(\tau_j)$$

$$w_j = \frac{p(\tau)\exp(r_\phi(\tau_j))}{\pi(\tau_j)} = \frac{\exp(\sum_t r_\phi(s_t, a_t))}{\prod_t \pi(a_t|s_t)}$$
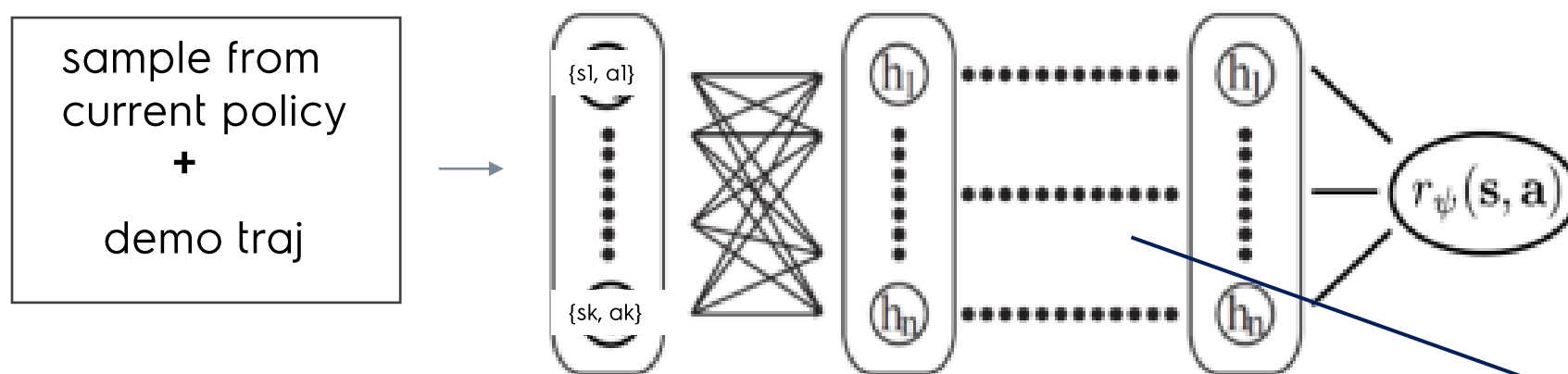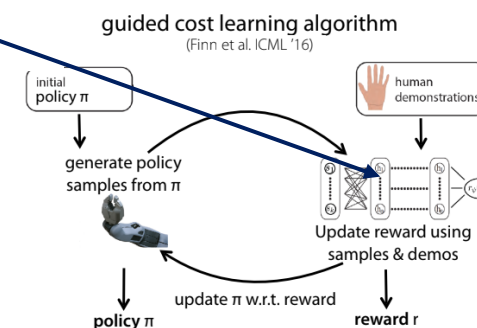
# GCL – Reward

- The neural network for cost function as loss takes:

$$\frac{1}{N}\sum_{\tau \in D} R_\phi(\tau) - \frac{1}{M}\log \sum_{\tau_j} \frac{\exp\left(R_\phi(\tau_j)\right)}{\pi(\tau_j)}$$

and performs the backward pass with gradient descent to improve
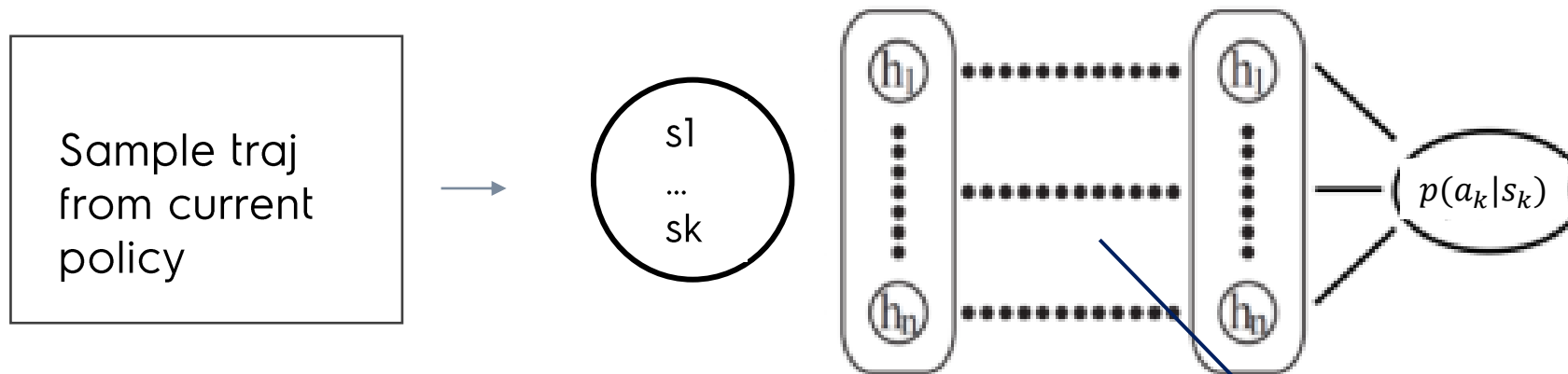


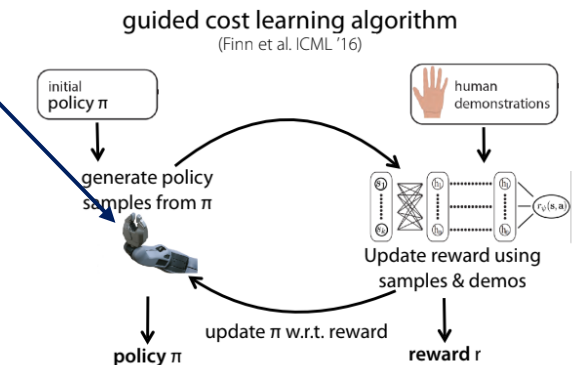The policy gets updated w.r.t to reward

# GCL – Policy

- At each improvement of the cost function, the policy performs a single Gradient Policy step (eg. REINFORCE algorithm)



Sample traj from current policy

$s1$
...
$sk$

$p(a_k|s_k)$

REINFORCE Algorithm:
1. Run forward pass of policy network
2. Get cumulative reward $r$ for each state, action based on current cost function
3. Calculate the expected reward $G = \log(p_\pi(a|s)) * \gamma r$
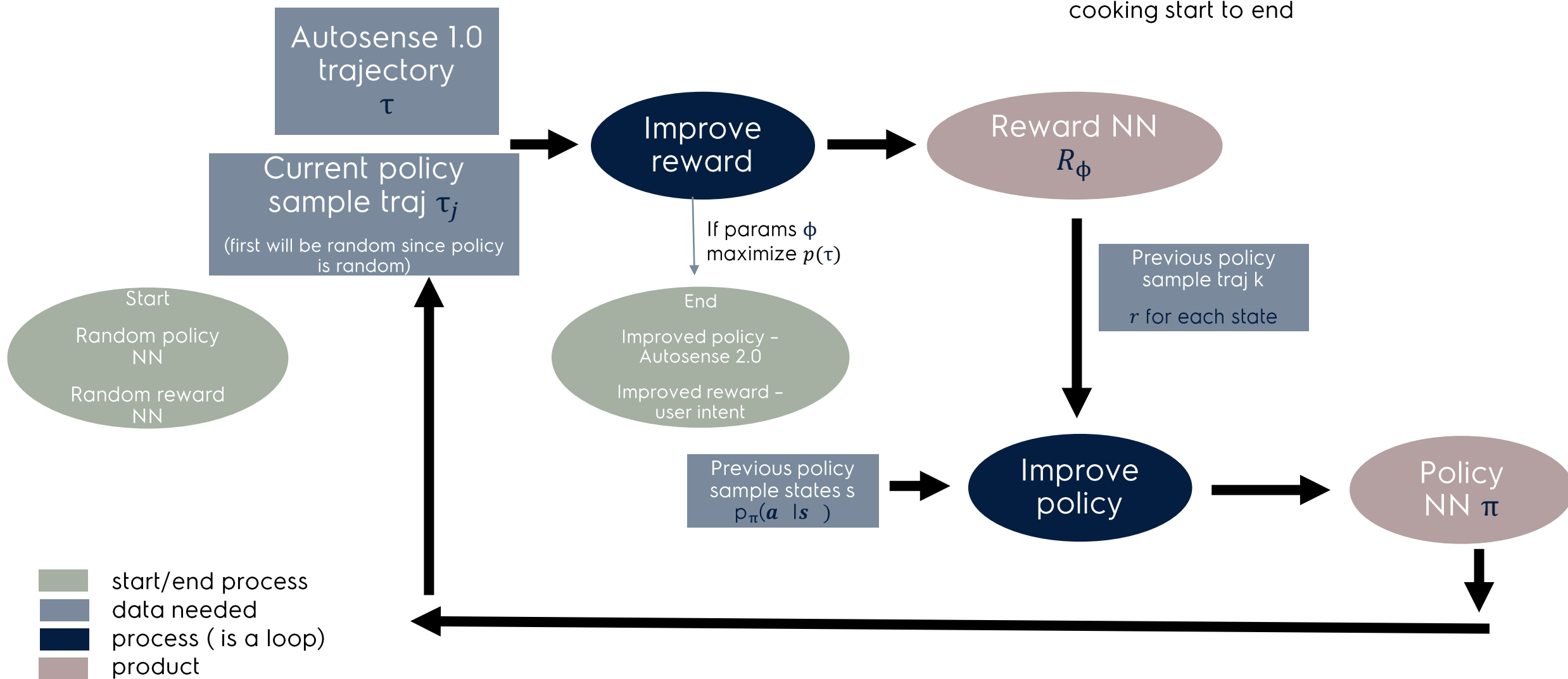4. Adjust weights by back-propagating the error $(-\sum G)$

guided cost learning algorithm
(Finn et al. ICML '16)

initial policy π

human demonstrations

generate policy samples from π

Update reward using samples & demos

update π w.r.t. reward

policy π

reward r

# Autosense context

**Policy** – Control (state->action)
**Reward** – Cleaner air/power consumption
**Trajectory** – tvoc, speed pair from cooking start to end

Autosense 1.0 trajectory
$\tau$

Current policy sample traj $\tau_j$

(first will be random since policy is random)

Start
Random policy NN
Random reward NN

Improve reward

If params $\phi$ maximize $p(\tau)$

End
Improved policy – Autosense 2.0
Improved reward – user intent

Reward NN
$R_\phi$

Previous policy sample traj k
$r$ for each state

Previous policy sample states s
$p_\pi(\boldsymbol{a}\ |\boldsymbol{s}\ )$

Improve policy

Policy NN $\pi$

start/end process
data needed
process ( is a loop)
product

# GCL – Algorithm

0. Input: demo trajectories, arbitrary (or from demos) policy q NN, arbitrary cost NN

1.  Generate sample trajectory for current q

2.  Forward pass of cost function with demo and q trajectory

3.  Calculate c loss $\frac{1}{N}\sum_{\tau \in D} R_\phi(\tau) - \frac{1}{M} log \sum_{\tau_j} \frac{\exp\left(R_\phi(\tau_j)\right)}{q(\tau_j)}$ and backpropagate error

4.  For q trajectory run current q, calculate cost f

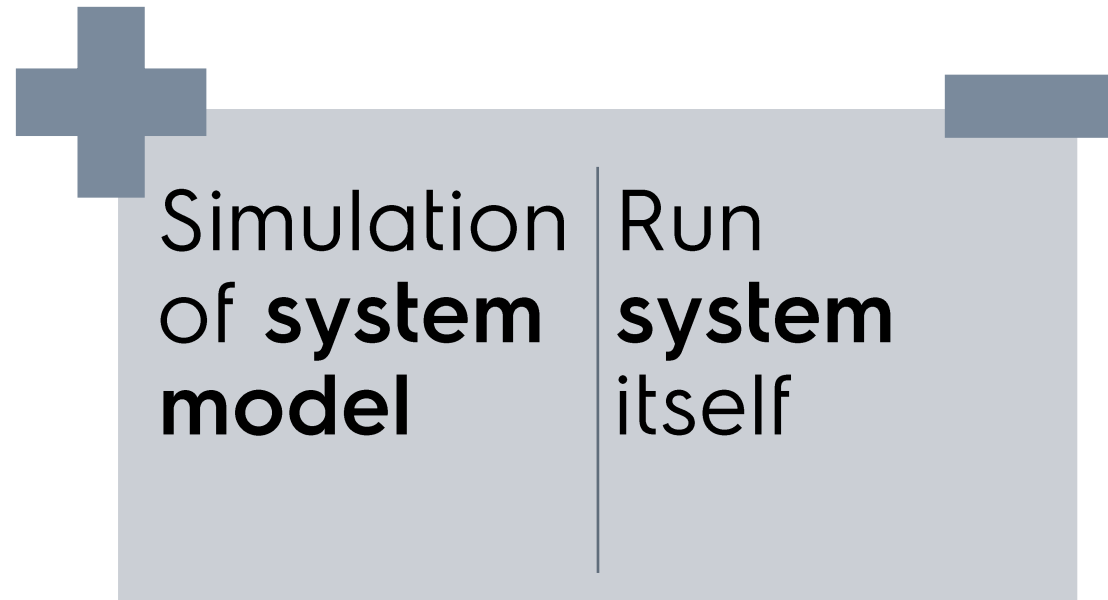5.  Calculate q loss $log(p_\pi(a|s))*\gamma r$ and backpropagate error

# GCL – (sometimes) environment

In the examples an environment is used, eg. CartPole/LunarLander–v2 with implemented physics to facilitate the simulation of real behavior of the system.

In our internal cases there would be a **third** model that is able to map state + action to new_state.

Simulation of **system model** | Run **system** itself

# gym/gymnasium : Lunar Lander v2

- v2 takes into consideration also the fuel consumption

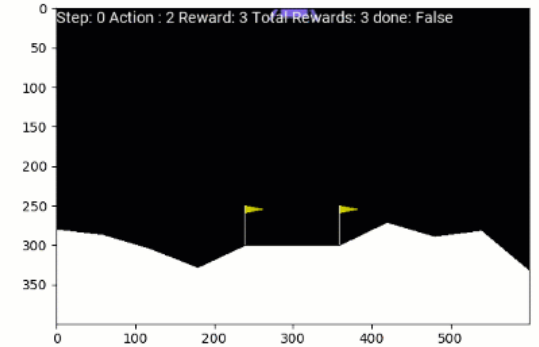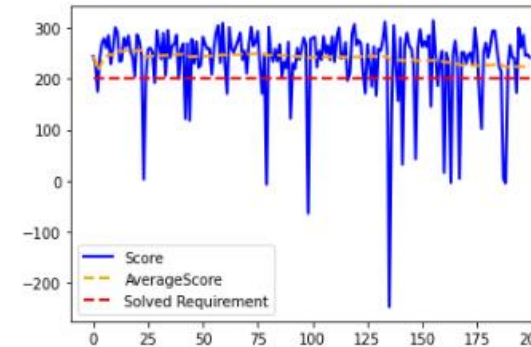| Actions | States | Rewards Rule | Episode termination |
|---------|--------|--------------|---------------------|
| 0 – No action<br>1 – Fire left engine<br>2 – Fire main engine<br>3 – Fire right engine | 0 – Lander horizontal coordinate<br>1 – Lander vertical coordinate<br>2 – Lander horizontal speed<br>3 – Lander vertical speed<br>4 – Lander angle<br>5 – Lander angular speed<br>6 – Bool: 1 if first leg has contact, else 0<br>7 – Bool: 1 if second leg has contact, else 0 | Moving from the top of the screen to the landing pad gives a scalar reward between (100-140)<br>Negative reward if the lander moves away from the landing pad<br>If the lander crashes, a scalar reward of (-100) is given<br>If the lander comes to rest, a scalar reward of (100) is given<br>Each leg with ground contact corresponds to a scalar reward of (10)<br>Firing the main engine corresponds to a scalar reward of (-0.3) per frame<br>Firing the side engines corresponds to a scalar reward of (-0.3) per frame | Lander crashes<br>Lander comes to rest<br>Episode length > 400 |

# GCL Implementations

- Implementation (**I1**) found online for *CartPole-v0* environment
  - Bugs fixed
- I1 adapted to *LunarLander-v2* environment, which has multiple actions on output, ie. four
  - collected expert trajectories from (yet another) implementation (**I2**) of LunarLander-v2, which trained successfully an agent for the problem

# GCL Results for LunarLander-v2

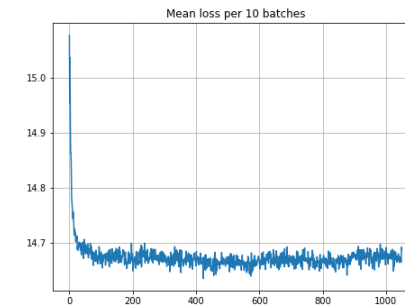- I2: Perform usual RL, train a policy (DQN eg.), and test to collect expert trajectories (demos): deleted bad episodes ~30/200
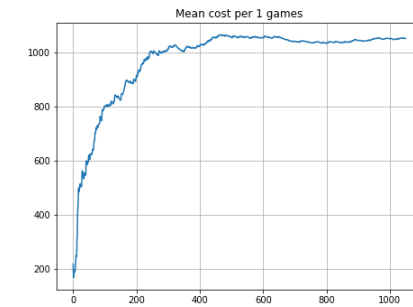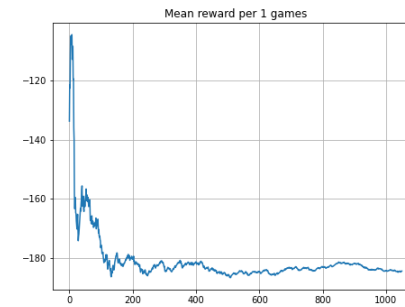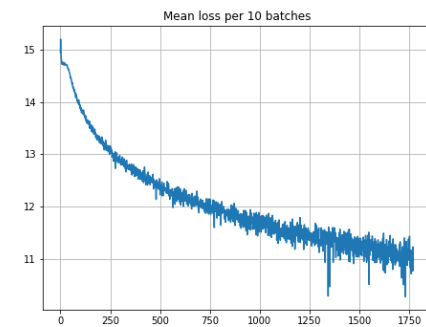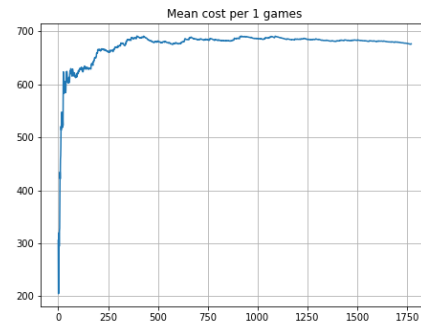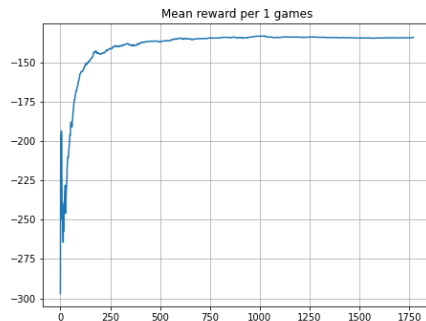




*NOTE: this particular policy trained with reward rules & REINFORCE algo (not I2)*

- I1: **Learn reward** and policy: converge to poor policy
- Next steps:
  Better baseline
  Actor-Critic
  Use the policy from I2 as starting policy to improve in GCL
  Not use env from gym, but generate steps (for trajectories) from learnt transition probabilities

# GCL Results for LunarLander-v2

- I1: Learn reward and policy
- Without weight decay it was diverging
- A. without baseline added, B. with whitening

# GCL Success Stories



Guided Cost Learning:
Deep Inverse Optimal Control via Policy Optimization

Chelsea Finn, Sergey Levine, Pieter Abbeel
UC Berkeley

Dish placement and pouring tasks. The robot learned
to place the plate gently into the correct slot, and to pour almonds, localizing the
target cup using unsupervised visual features.

# Literature

- A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models, C. Finn, P. Christiano, P. Abbeel, S. Levine

- Inverse RL, Andrew Ng et al, 2000

- MaxEntropy Ziebart et al., 2008

- Guided Cost Learning, Finn et al, 2016

- Generative Adversarial Imitation Learning, Jonathan Ho, Stefano Ermon, 2016

- A Survey of Inverse Reinforcement Learning: Challenges, Methods and Progress, Saurabh Arora et al, 2018