# Hands-on Reinforcement Learning for RecSys - from Bandits to Offline RL with Ray RLlib

Kourosh Hakhamaneshi - kourosh@anyscale.com
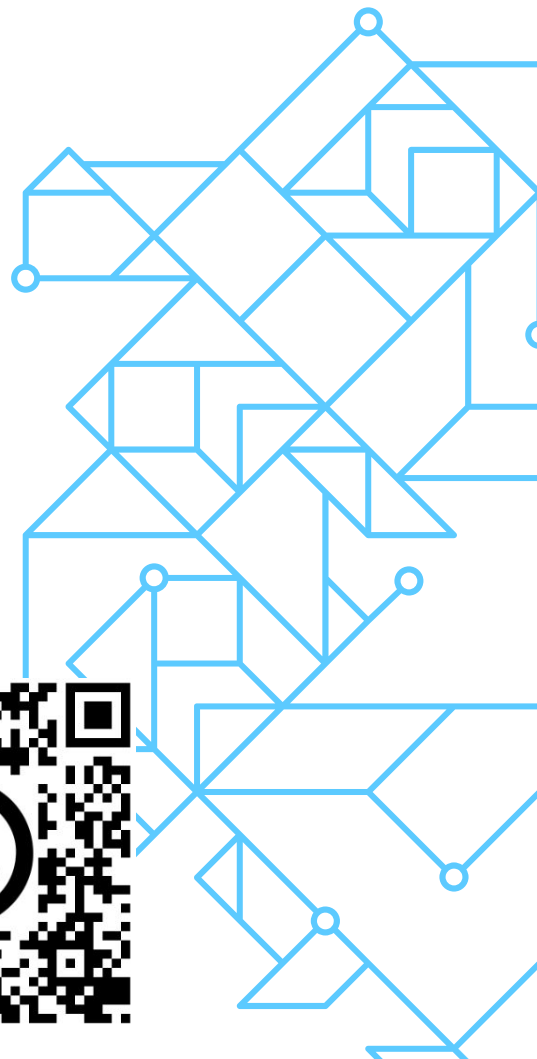Christy Bergman - christy@anyscale.com

# Few Important URLs

Keep these URLs open in your browser tabs

→ **GitHub: https://bit.ly/rllib_recsys_2022_github**

→ Q & A Doc: https://bit.ly/rllib_recsys_2022-qa

→ Logins+passwords: https://bit.ly/rlib_recsys-logins

→ Anyscale: console.anyscale.com
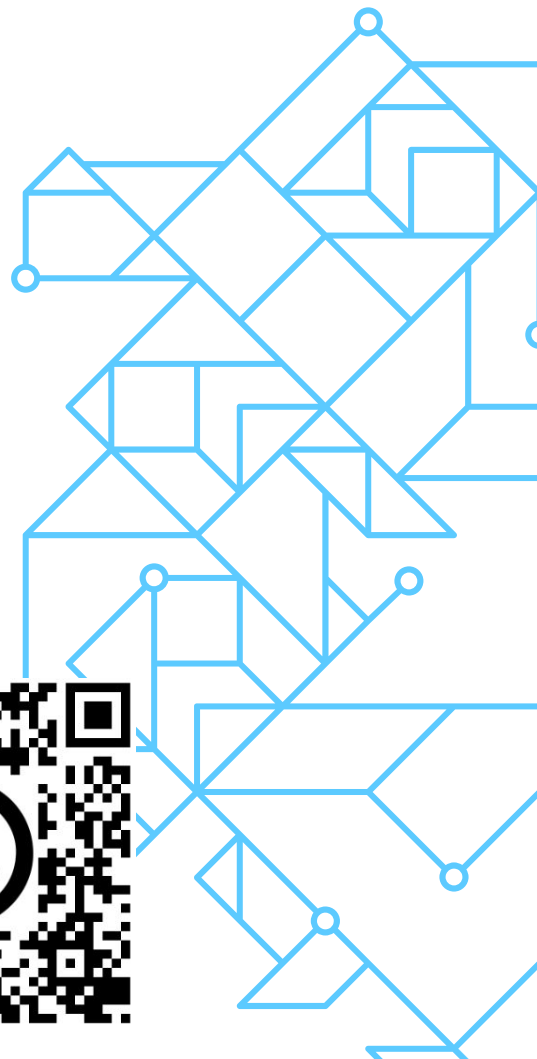
→ Tutorial Survey: https://bit.ly/rlib_recsys_2022

# Few Important URLs

Keep these URLs open in your browser tabs

→  GitHub: https://bit.ly/rllib_recsys_2022_github
→  **Q & A Doc: https://bit.ly/rllib_recsys_2022-qa**
→  Logins+passwords: https://bit.ly/rlib_recsys-logins
→  Anyscale: console.anyscale.com
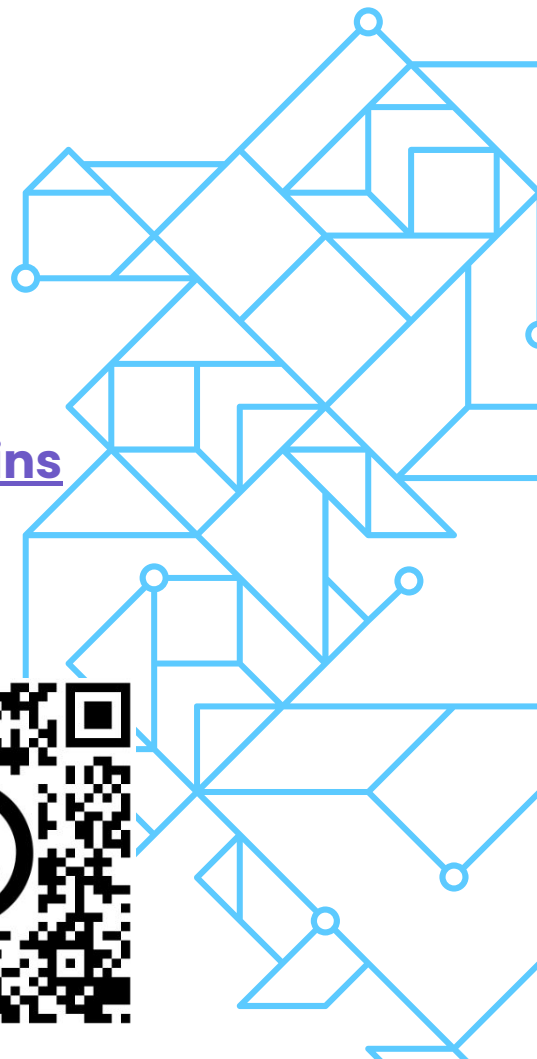→  Tutorial Survey: https://bit.ly/rlib_recsys_2022

# Few Important URLs

Keep these URLs open in your browser tabs

→ GitHub: https://bit.ly/rllib_recsys_2022_github

→ Q & A Doc: https://bit.ly/rllib_recsys_2022-qa

→ **Logins+passwords: https://bit.ly/rlib_recsys-logins**

→ **Anyscale: console.anyscale.com**
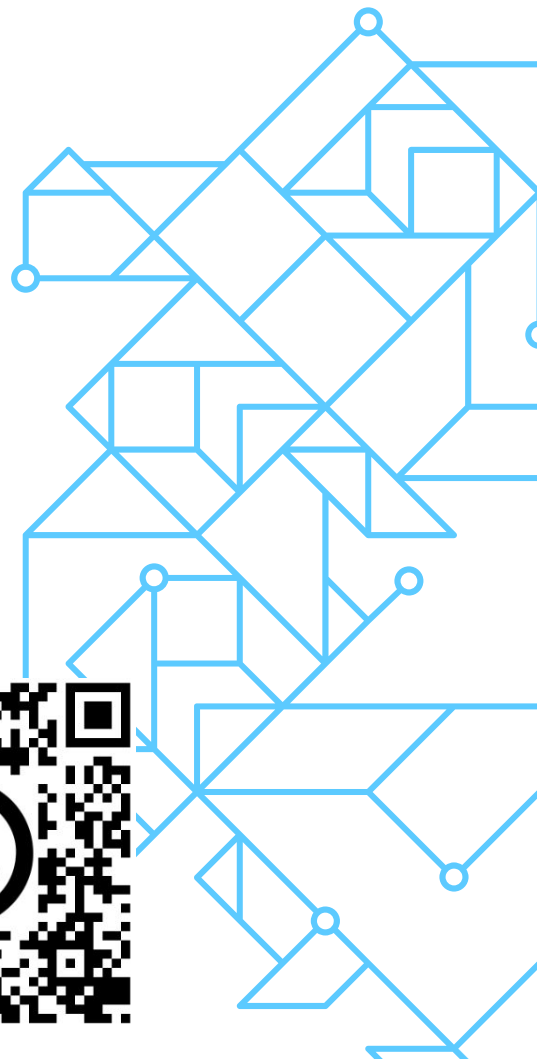
→ Tutorial Survey: https://bit.ly/rlib_recsys_2022

# Few Important URLs

Keep these URLs open in your browser tabs

→ GitHub: https://bit.ly/rllib_recsys_2022_github

→ Q & A Doc: https://bit.ly/rllib_recsys_2022-qa

→ Logins+passwords: https://bit.ly/rlib_recsys-logins

→ Anyscale: console.anyscale.com

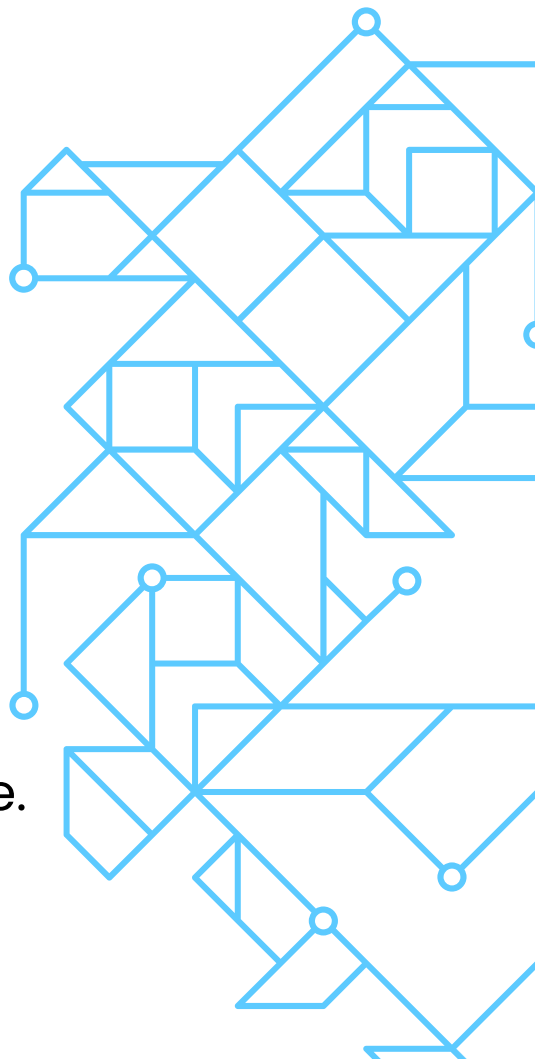→ **Tutorial Survey: https://bit.ly/rlib_recsys_2022**

# $whoami (Christy)

→   AI/ML DevAdvocate @Anyscale.
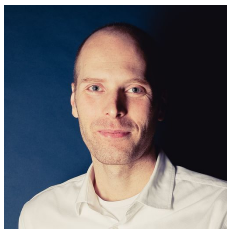→   Previously: AI/ML Solutions Architect at AWS, before that data scientist real-time fraud detection

# $whoami (Kourosh)

→   ML Engineer working on RL and RLlib @Anyscale.
→   Previously: PhD student at UC Berkeley working on RL in Robotics and design optimization

RL Team @ Anyscale

Sven    Jun    Avnish    Artur    Kourosh    Christy (devAdvocate)
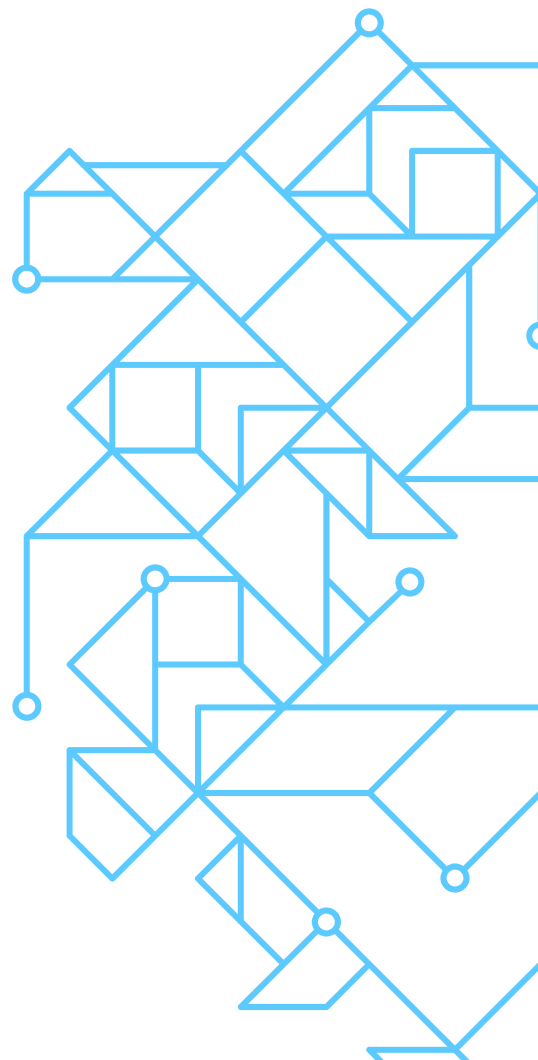
# Anyscale

**Who we are:** Original creators of Ray, a unified framework for scalable, distributed computing. Part of that framework are our libraries for ML and data processing.

**What we do:** Scalable compute for AI and Python

**Why we do it:** Scaling is a necessity, scaling is hard; make distributed computing easy and simple for all developers.

# Some of RLlib's Industry Users

# Overview of the tutorial

→ Brief intro RL

→ Brief intro RecSys

  + Traditional Approaches
  + Defining RecSys as an RL problem

→ Online RL vs Offline RL

→ Hands-on coding with python notebooks and scripts

**Goals – Understand:**

➢ What are the advantages of using RL in RecSys?
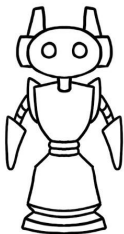➢ What are the pros and cons of offline RL in practical scenarios?

# Overview of the tutorial

→ 4 min: Welcome

→ 5 min: Very brief intro RL

→ 5 min: Very brief intro RecSys

+ Machine learning (ML) approach
+ Challenges with current ML approach
+ Map RecSys problem into MDP for RL

→ 5 min: Intro Online RL vs Offline RL

→ 1 hour: Hands-on with Google Colab

+ 15min: Introduction to the environment
+ 10min: Run baselines, bandit, and RL algorithm
+ 5min: Conclusion so far TODO ADD slide with results
+ 10min: Run offline RL on expert, random, greedy data
+ 5min: Conclusion so far TODO ADD slide with results
+ 5min: Deploy a policy to production using Ray Serve

# Brief intro RL
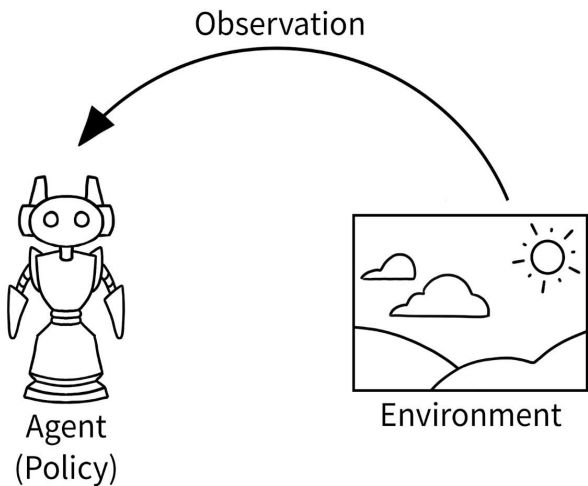
Conversation between an agent and an environment.

Agent
(Policy)

Environment

# Brief intro RL



Conversation between an agent and an environment.

# Brief intro RL

Observation
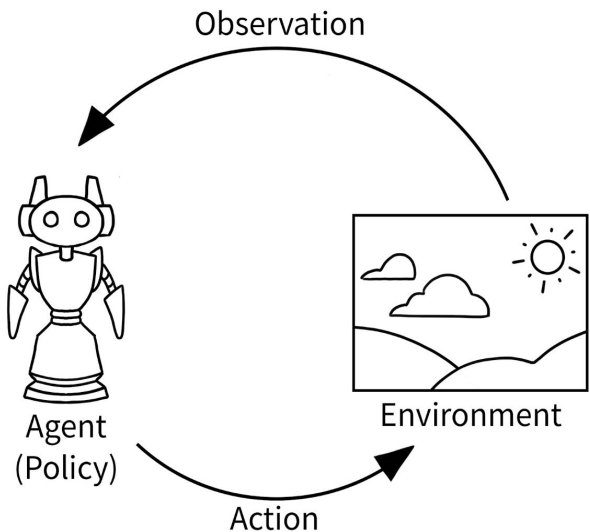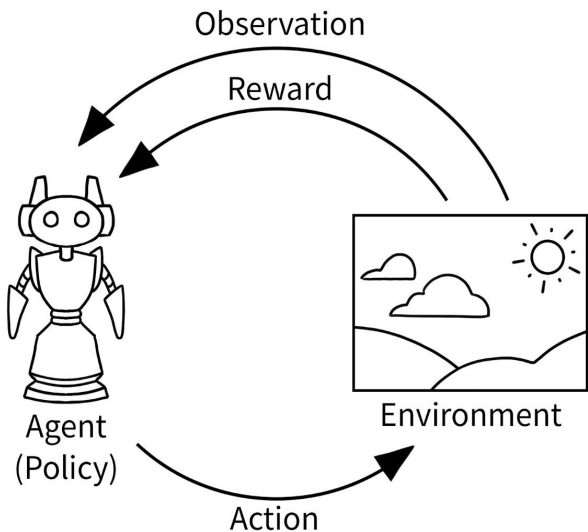
Agent
(Policy)

Action

Environment

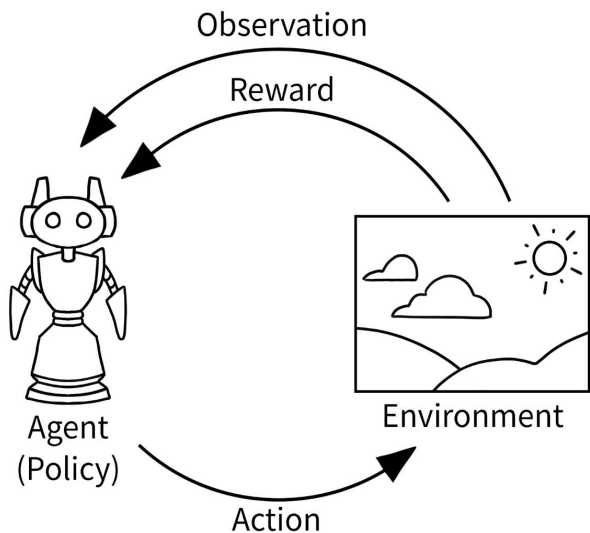Conversation between an agent and an environment.

# Brief intro RL



Conversation between an agent and an environment.

# Brief intro RL



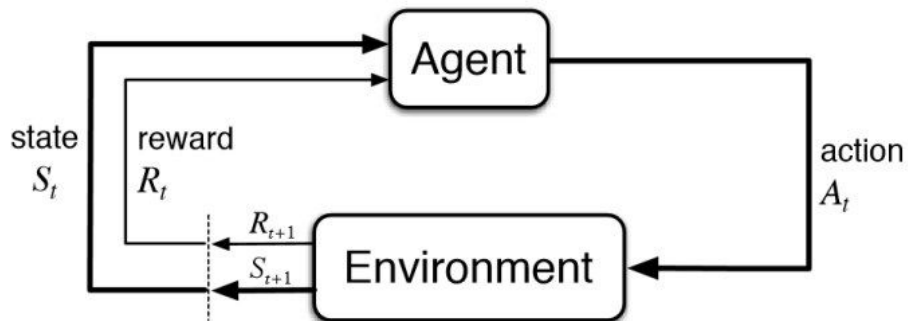Conversation between an agent and an environment.

Learning objectives:
- Maximize sum of rewards.
- Learn from delayed reward.
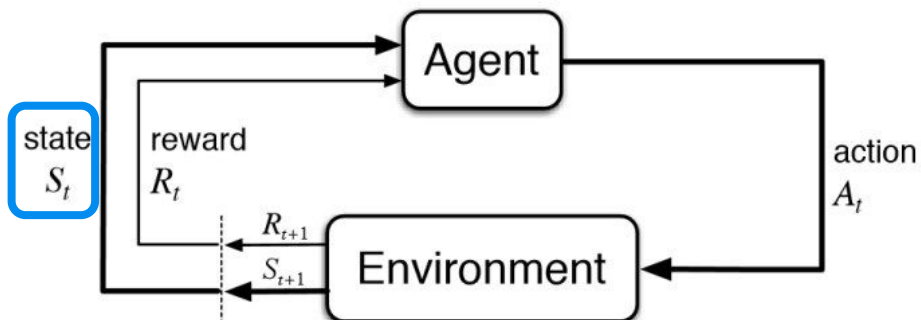- Proper exploration to maximally learn

# Brief intro RL

$$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$$

# Brief intro RL

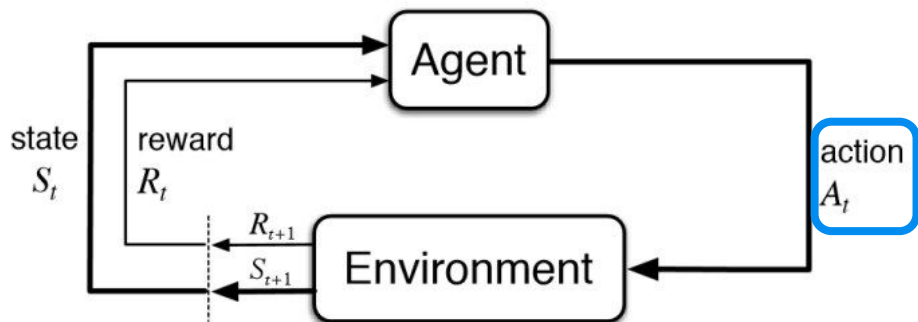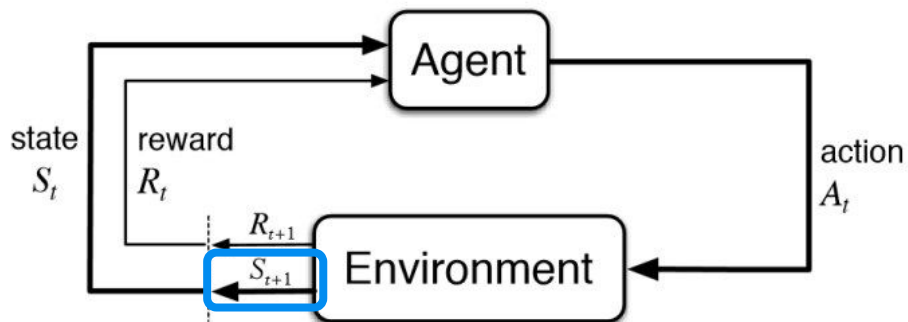$$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$$

# Brief intro RL

$$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$$

# Brief intro RL



$$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$$
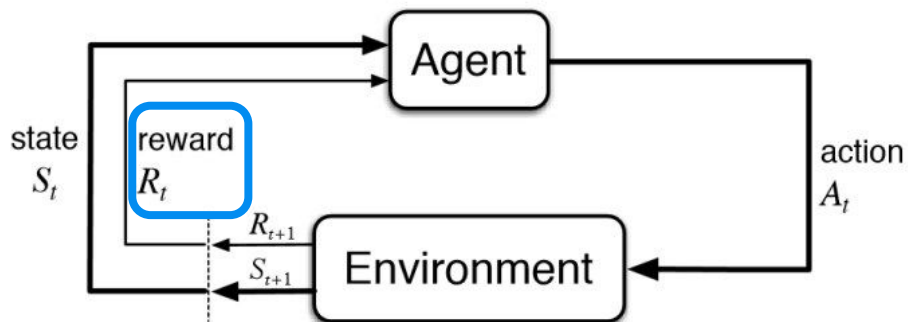
$$S_0 \sim \mathcal{P}(.)$$

$$S_{t+1} \sim \mathcal{P}(.|S_t, A_t)$$

# Brief intro RL



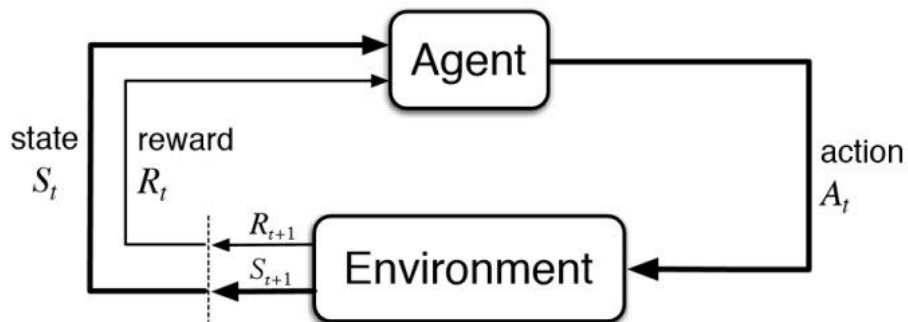$$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \boxed{\mathcal{R},} \gamma)$$

$$R_t = \mathcal{R}(S_t, A_t)$$

# Brief intro RL



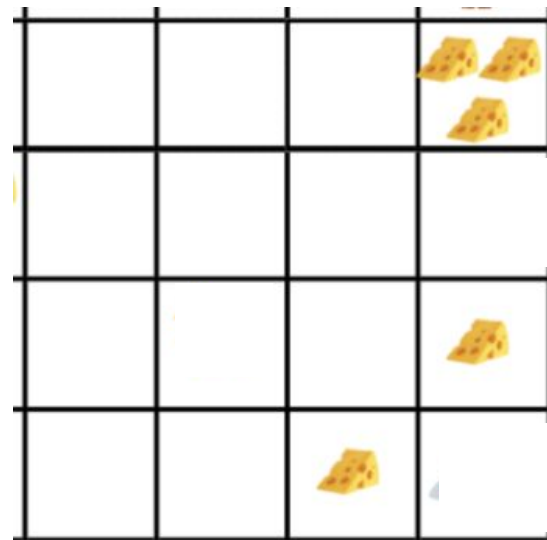$$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \boxed{\gamma})$$

$$R(\tau) = \sum_t \gamma^t R_t$$

# Discount factor $\gamma$ in RL

- If $\gamma = 0$, the algorithm considers **1-step rewards only**.

- If $\gamma = 1$, the algorithm considers all future rewards equally.

# Brief intro RecSys

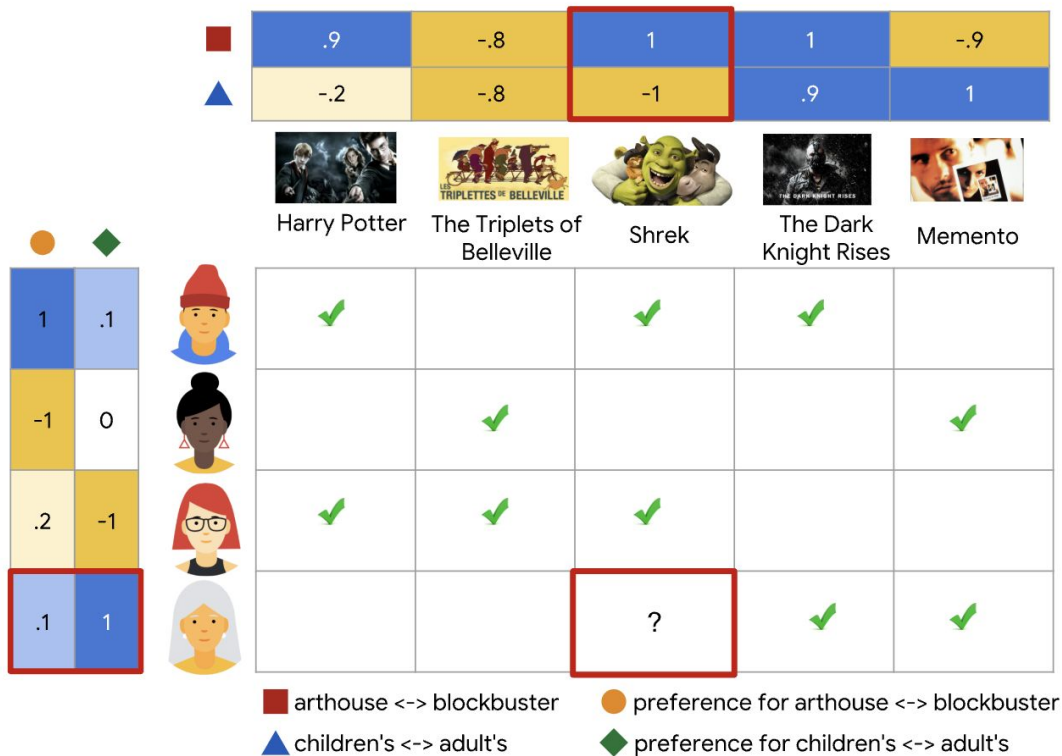Companies want to recommend content.

ML: Pointwise recommendations.

RL: Combine pointwise recommendations with session based data.

# Brief intro RecSys ML

| ■ | .9 | -.8 | 1 | 1 | -.9 |
|---|-----|-----|---|---|-----|
| ▲ | -.2 | -.8 | -1 | .9 | 1 |

| | Harry Potter | The Triplets of Belleville | Shrek | The Dark Knight Rises | Memento |
|---|---|---|---|---|---|

| ● | ◆ | | Harry Potter | The Triplets of Belleville | Shrek | The Dark Knight Rises | Memento |
|---|---|---|---|---|---|---|---|
| 1 | .1 | 🧑 | ✔ | | ✔ | ✔ | |
| -1 | 0 | 👩 | | ✔ | | | ✔ |
| .2 | -1 | 👩 | ✔ | ✔ | ✔ | | |
| .1 | 1 | 👵 | | | ? | ✔ | ✔ |

■ arthouse <-> blockbuster    ● preference for arthouse <-> blockbuster

▲ children's <-> adult's    ◆ preference for children's <-> adult's

Credit:  https://developers.google.com/machine-learning/

# Challenges with traditional ML in RecSys

- Traditional ML (collaborative filtering) models are **static with respect to time**.

- This type of model **ignores time order** in which users did actions.

# Challenges with traditional ML in RecSys

- Traditional ML (collaborative filtering) models are **static with respect to time**.

  - Ignores the **sequence of interactions** with a given user.

# Challenges with traditional ML in RecSys

- Traditional ML (collaborative filtering) models are **static with respect to time**.

  - Ignores the **sequence of interactions** with a given user.

- Static models can be:

  - Too short-sighted and **miss out on Long-term, delayed rewards**

# Challenges with traditional ML in RecSys

- Traditional ML (collaborative filtering) models are **static with respect to time**.

  - Ignores the **sequence of interactions** with a given user.

- Static models can be:

  - Too short-sighted and **miss out on Long-term, delayed rewards**

  - **Overlook important and changing user intents** or business conditions such as seasonality or promotional campaigns

# New way: RL in RecSys

- By taking each user's session history as a sequence of decisions, **the RecSys problem can be converted into a sequential decision-making problem**.

  - $\Pr[R(t+1)=r_t \mid A(t)=a_t, S(t)=s_t, A(t-1)=a_{t-1}, \ldots S(0)=s_0]$

- A stochastic process is a **Markov Decision Process (MDP)** if the values at time t depend only on the values at time t-1.

  - $$Q_\pi(s, a) = E_\pi \left[ \sum_{j=0}^{T} \gamma^j r_{t+j+1} \mid S_t = s, A_t = a \right]$$

- **RL has become the de-facto ML approach for solving MDPs.**

# New way: RL in RecSys

- By taking each user's session history as a sequence of decisions, **the RecSys problem can be converted into a sequential decision-making problem**.

# New way: RL in RecSys

- By taking each user's session history as a sequence of decisions, **the RecSys problem can be converted into a sequential decision-making problem**.

  States: item features, user feature, history of interactions

# New way: RL in RecSys

- By taking each user's session history as a sequence of decisions, **the RecSys problem can be converted into a sequential decision-making problem**.

  States: item features, user feature, history of interactions

  Actions: the items to recommend

# New way: RL in RecSys

- By taking each user's session history as a sequence of decisions, **the RecSys problem can be converted into a sequential decision-making problem**.

  States: item features, user feature, history of interactions

  Actions: the items to recommend

  Reward: long term satisfaction (explicit or implicit)

# New way: RL in RecSys

- By taking each user's session history as a sequence of decisions, **the RecSys problem can be converted into a sequential decision-making problem**.

  States: item features, user feature, history of interactions

  Actions: the items to recommend

  Reward: long term satisfaction (explicit or implicit)

  Gamma: 0 (bandits) or 1 (RL)

# RL Environment:
# Delayed Rewards &
# Long Term Satisfaction (LTS)

top 20 candidates

1

# RL Environment:
# Delayed Rewards &
# Long Term Satisfaction (LTS)

Recommendation
1 (action)

top 20 candidates

# RL Environment:
# Delayed Rewards &
# Long Term Satisfaction (LTS)



top 20 candidates

Recommendation
1 (action)

User decides to click or not

States, actions, rewards updated.
Policy model is retrained.
Best next action is sampled

Repeat until user quits
or end of time limit

$max_\pi \left( \sum_{j=0}^{T} \gamma^j \hat{R}_j \right)$

# RL Environment: Delayed Rewards & Long Term Satisfaction (LTS)



top 20 candidates

Recommendation 1 (action)

User decides to click or not

3

States, actions, rewards updated.
Policy model is retrained.
Best next action is sampled

Repeat until user quits or end of time limit

4

$$max_\pi \left( \sum_{j=0}^{T} \gamma^j \hat{R}_j \right)$$

5

# RL Environment:
# Delayed Rewards &
# Long Term Satisfaction (LTS)



top 20 candidates

Recommendation
1 (action)

User decides to click or not

Repeat until user quits
or end of time limit

States, actions, rewards updated.
Policy model is retrained.
Best next action is sampled

$$max_\pi \left( \sum_{j=0}^{T} \gamma^j \hat{R}_j \right)$$

# Online RL vs Offline RL

| Online RL | Batch / Offline RL |
|---|---|
|  | |

Online RL diagram: Agent — From State s, Take action a → Environment. Get back next reward r' and next state s'.

# Online RL vs Offline RL

| Online RL | Batch / Offline RL |
|---|---|
|  |   $[s, a, s', r', \dots]$ |

# Your Anyscale Cluster

→   Claim username/password at https://bit.ly/rlib_recsys-logins
    +   Update the "Status" column to "Available" or "Claimed"
→   Go to Console: http://console.anyscale.com/
→   Enter username (for the email) and password

# Your Anyscale Cluster

# Your Anyscale Cluster

# Your Anyscale Cluster

# Your Anyscale Cluster

# Your Anyscale Cluster

**9** → **Navigate to "acm_recsys_tutorial_2022"**



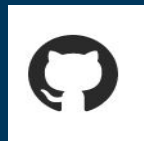**10** → **Open notebook "01_..."**

# Thank you.

We would love to connect with you!

**Twitter - @anyscalecompute | @raydistributed**

**GitHub - https://github.com/ray-project/ray**

**Slack - https://www.ray.io/community**

**Discuss - https://discuss.ray.io/**