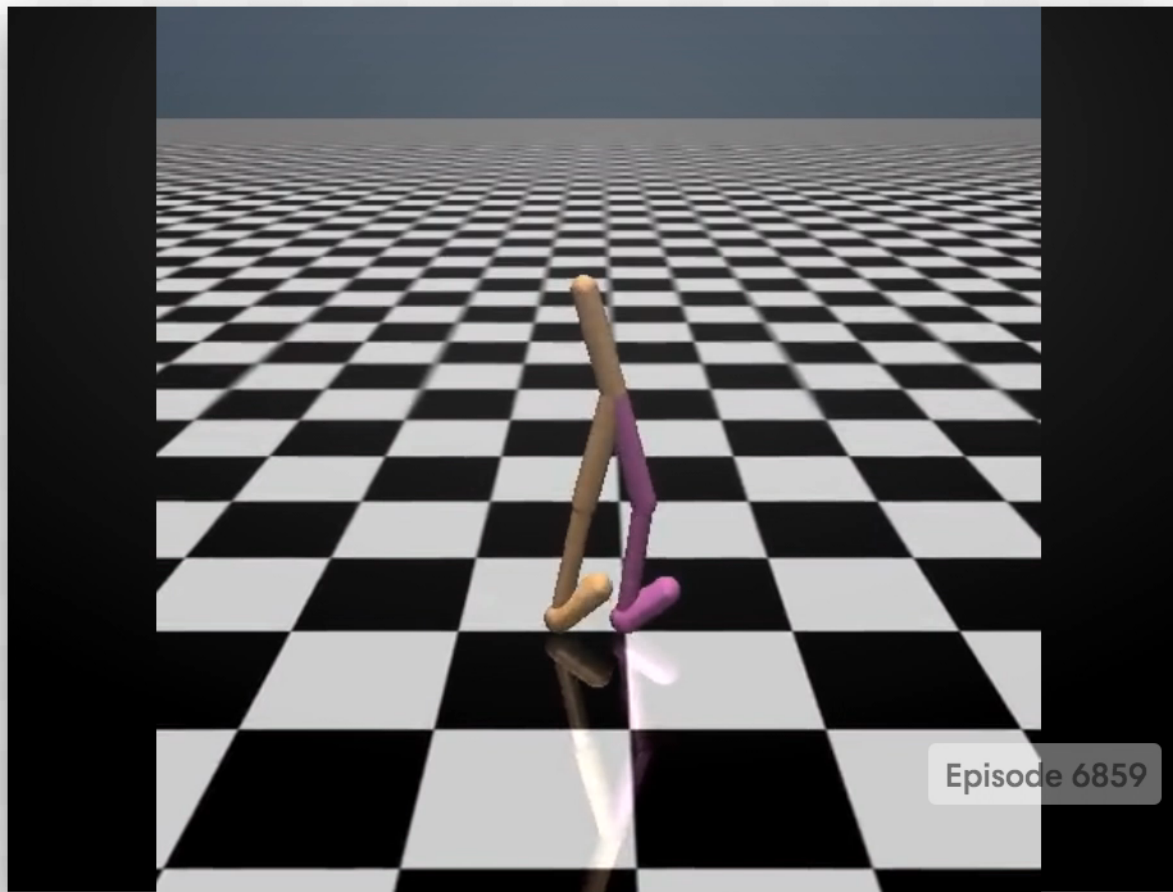# Reinforcement Learning

Min Sun

VSLab

# Vision and Control

- Learning to play game with weak supervision:

  Reinforcement Learning (RL)



joschu's trpo-gae-v0 on Walker2d-v0

Episode 6859

https://gym.openai.com/

# Where It All Begins ...

# Playing Atari with
# Deep Reinforcement Learning

*by DeepMind in NIPS 2013 Deep Learning Wrokshop*



slides by
Yen-Chen Lin

# Control: Learning to Act

Play **Breakout** equals to

- **Input:** screen images
- **Output:** actions
  (**do nothing** | **left** | **right**)

## Supervised Classification

# Supervised Solution

- **Training data:** Record experts game sessions
- **Target label:** Action experts take at every step

⚠️ Problems:

- What if there's no expert?

- This is not how human learns

# How Human Learns

- Don't need somebody to tell us a million times which move to choose at each screen

- Just need **occasional feedback** that we did the right thing



slides by
Yen-Chen Lin

# Reinforcement Learning

- Somewhere between **supervised** and **unsupervised** learning
- **Sparse** and **time-delayed** labels

Based only on those rewards, the agent has to learn to behave in the environment.
A rational agent should optimize total reward.

slides by
Yen-Chen Lin

# AlphaGo



# Mastering the game of Go with deep neural networks and tree search

David Silver[1]*, Aja Huang[1]*, Chris J. Maddison[1], Arthur Guez[1], Laurent Sifre[1], George van den Driessche[1], Julian Schrittwieser[1], Ioannis Antonoglou[1], Veda Panneershelvam[1], Marc Lanctot[1], Sander Dieleman[1], Dominik Grewe[1], John Nham[2], Nal Kalchbrenner[1], Ilya Sutskever[2], Timothy Lillicrap[1], Madeleine Leach[1], Koray Kavukcuoglu[1], Thore Graepel[1] & Demis Hassabis[1]

# Lecture 1: Introduction to Reinforcement Learning

David Silver

# Sequential Decision Making

- Goal: *select actions to maximise total future reward*
- Actions may have long term consequences
- Reward may be delayed
- It may be better to sacrifice immediate reward to gain more long-term reward
- Examples:
    - A financial investment (may take months to mature)
    - Refuelling a helicopter (might prevent a crash in several hours)
    - Blocking opponent moves (might help winning chances many moves from now)