



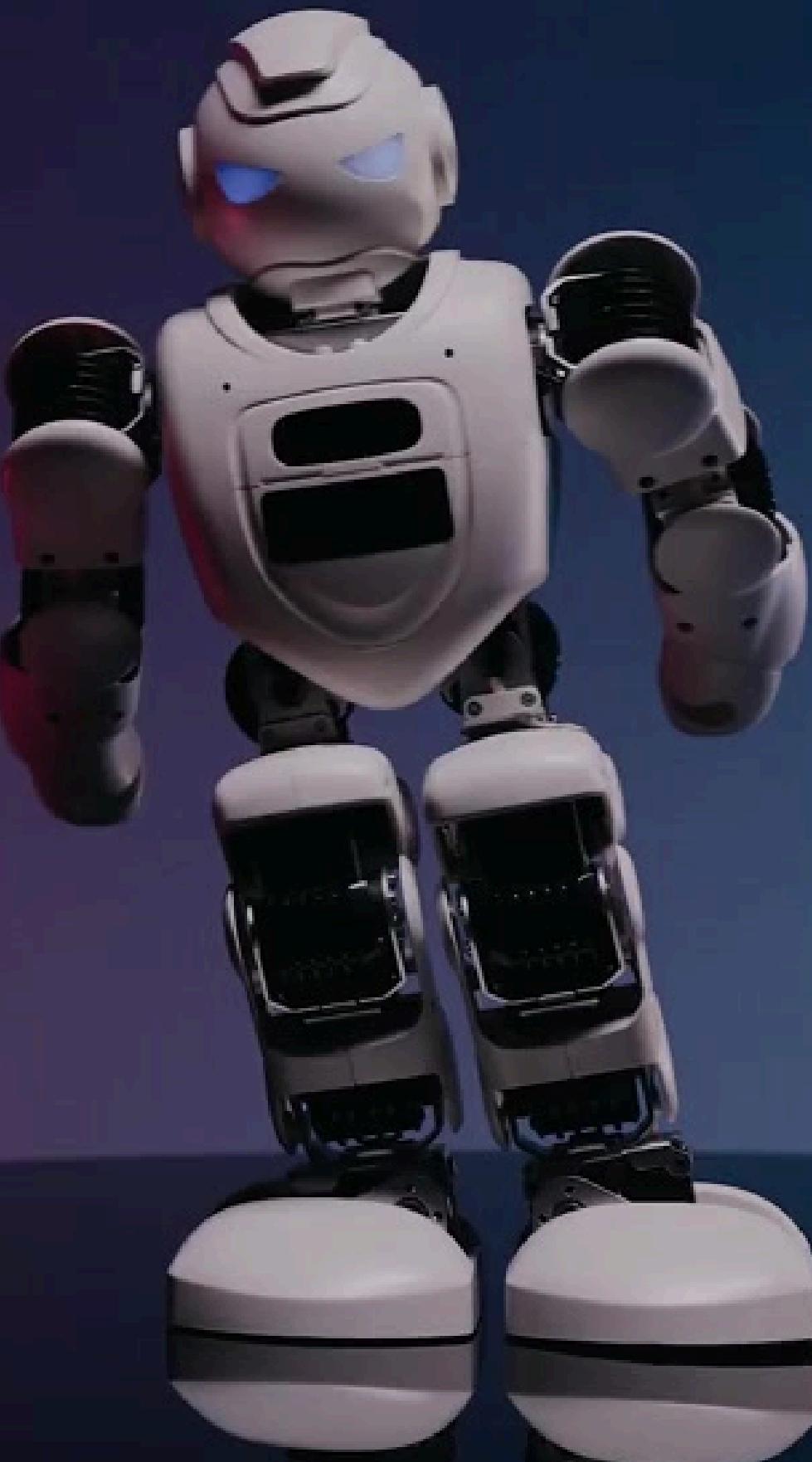
WEEK 11 RL

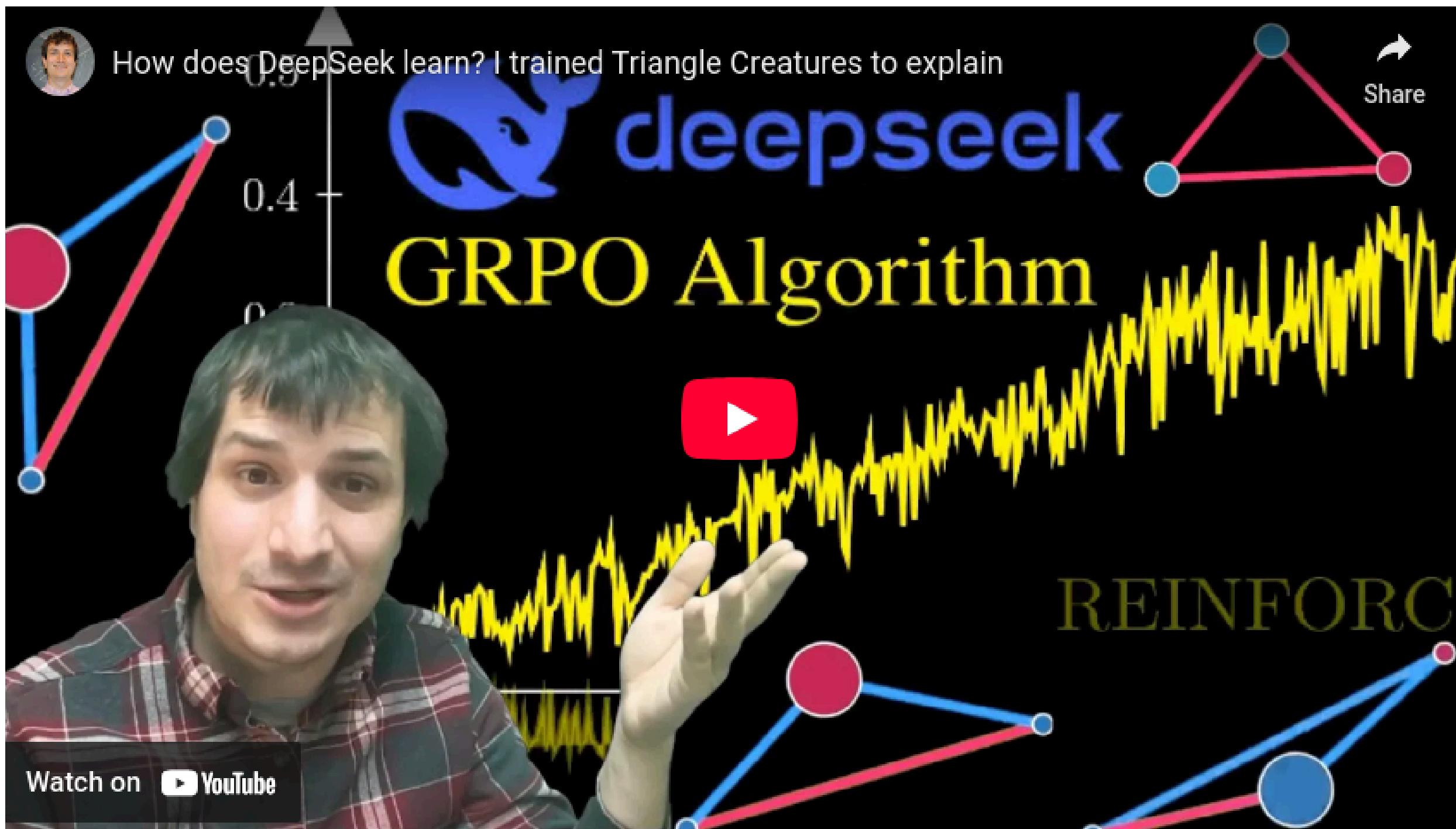
DEEP LEARNING FOR COMPUTER VISION

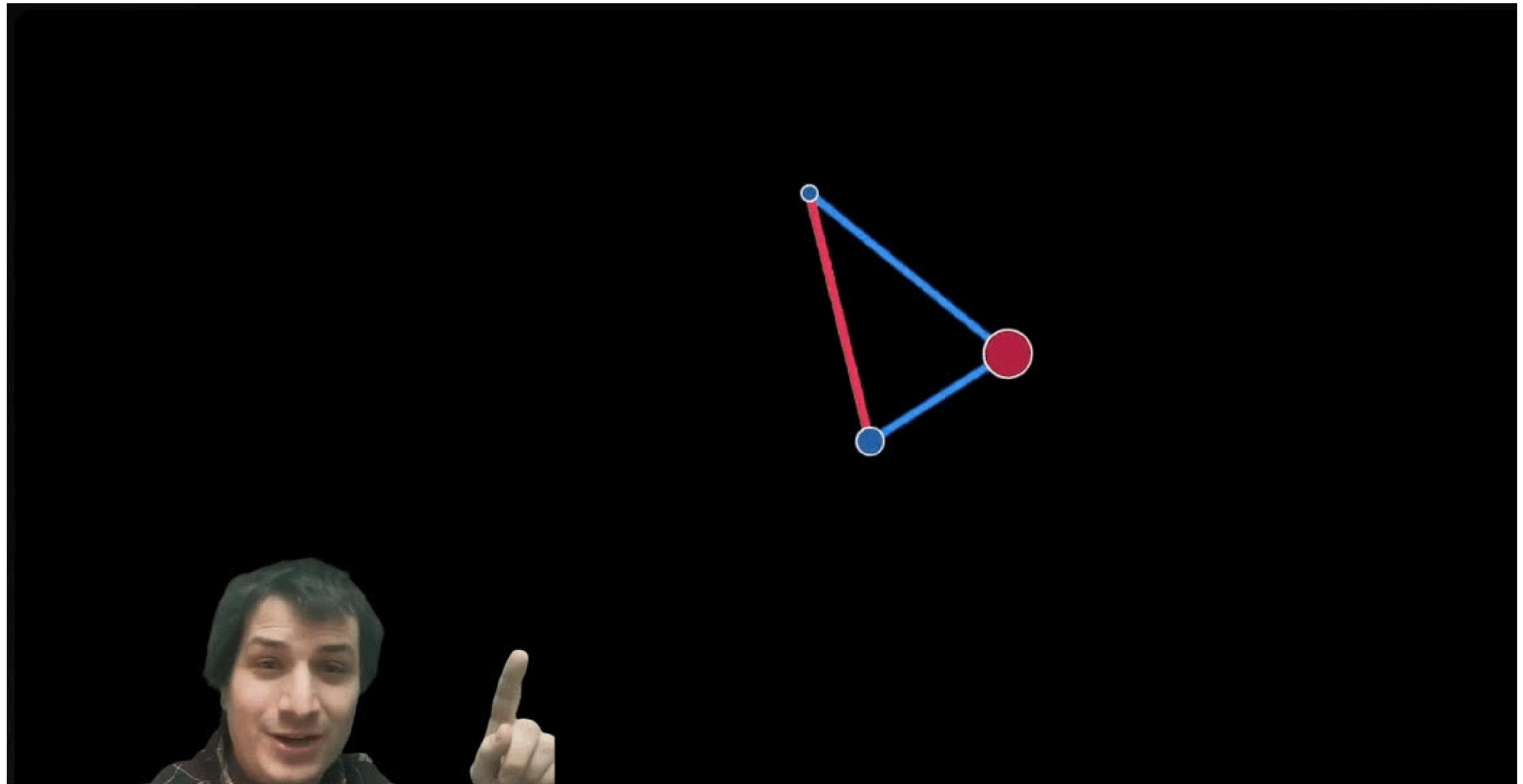


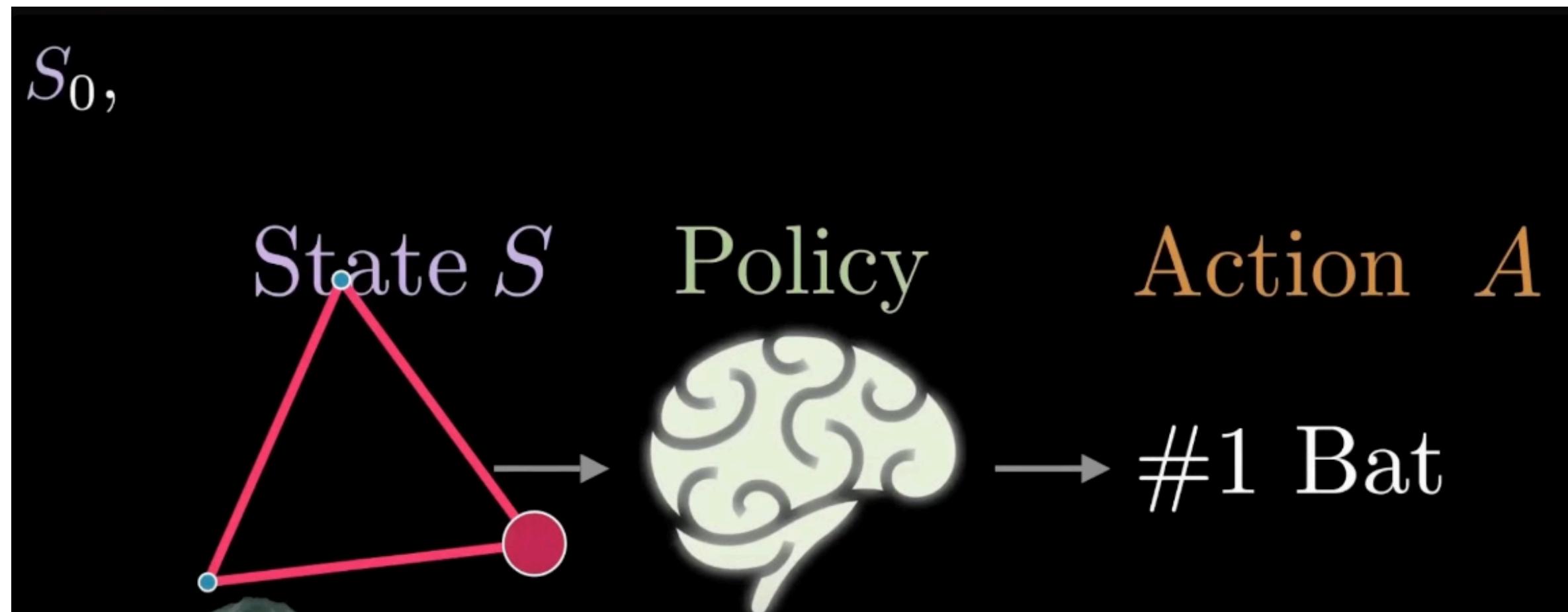
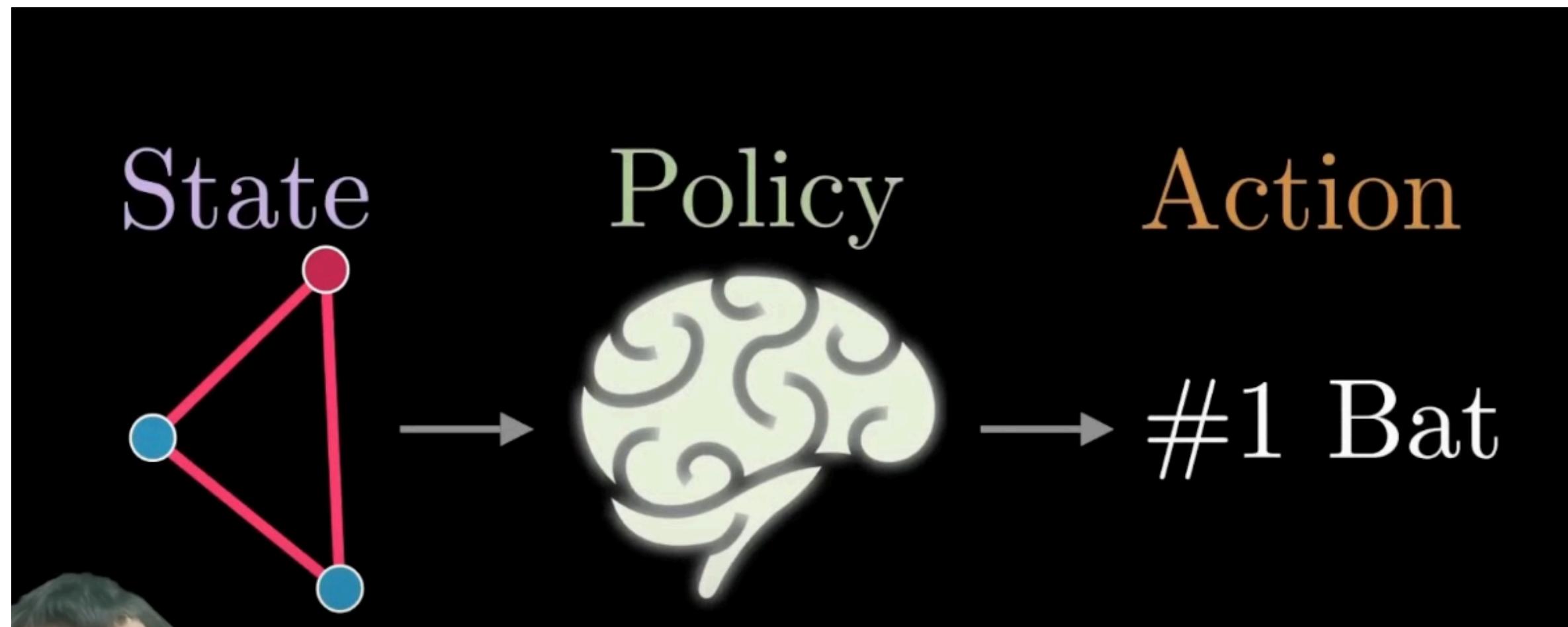
UNLIMITED

Presented by **Asst. Prof. Dr. Tuchsanai Ploysuwan**

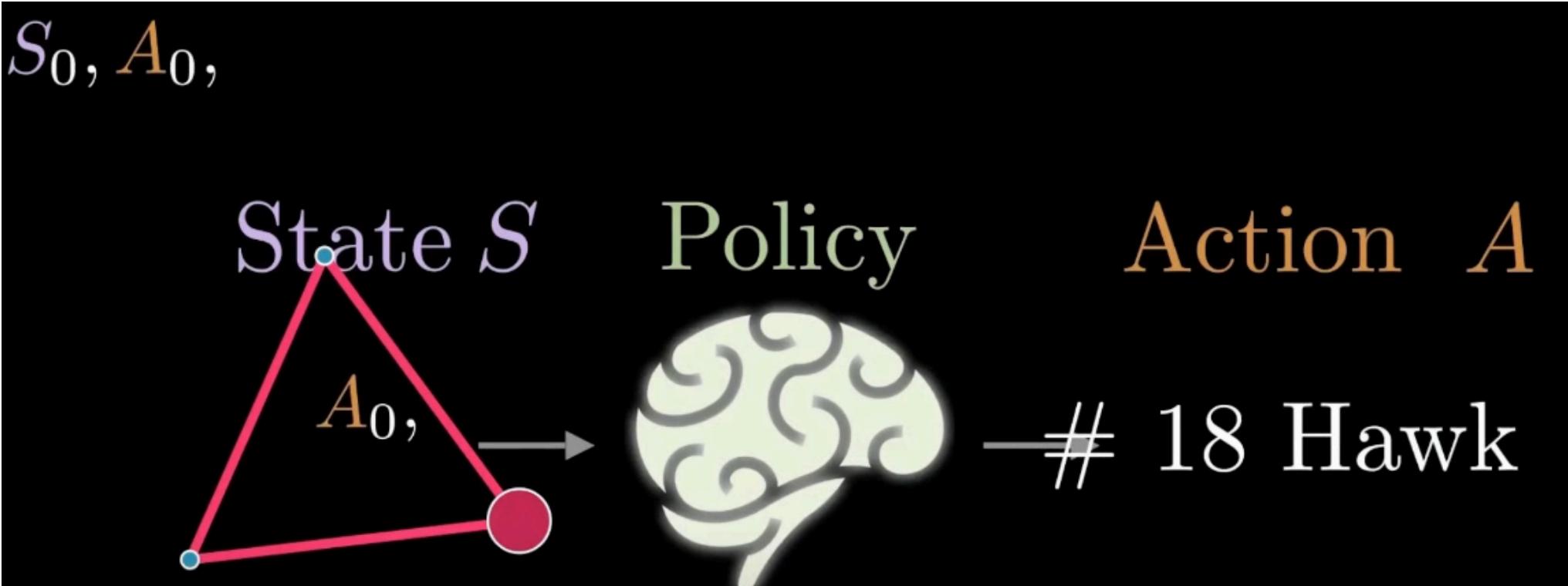




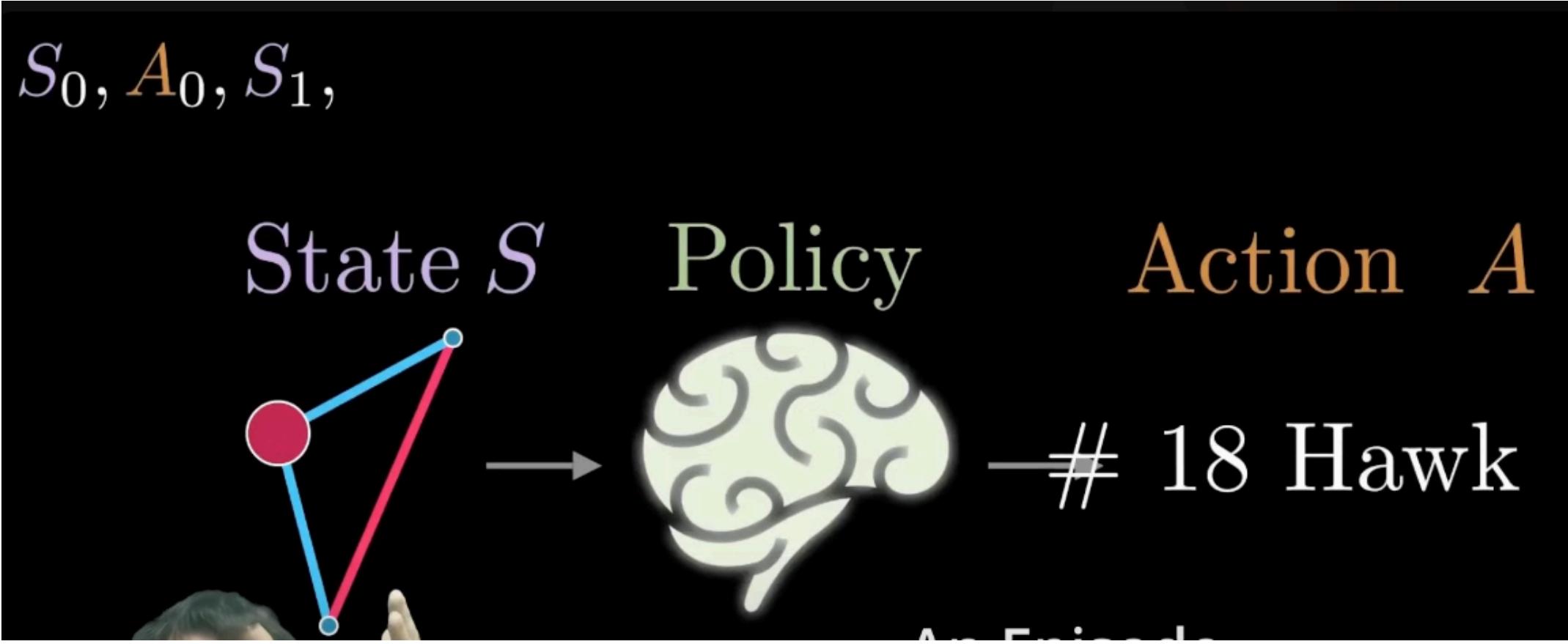


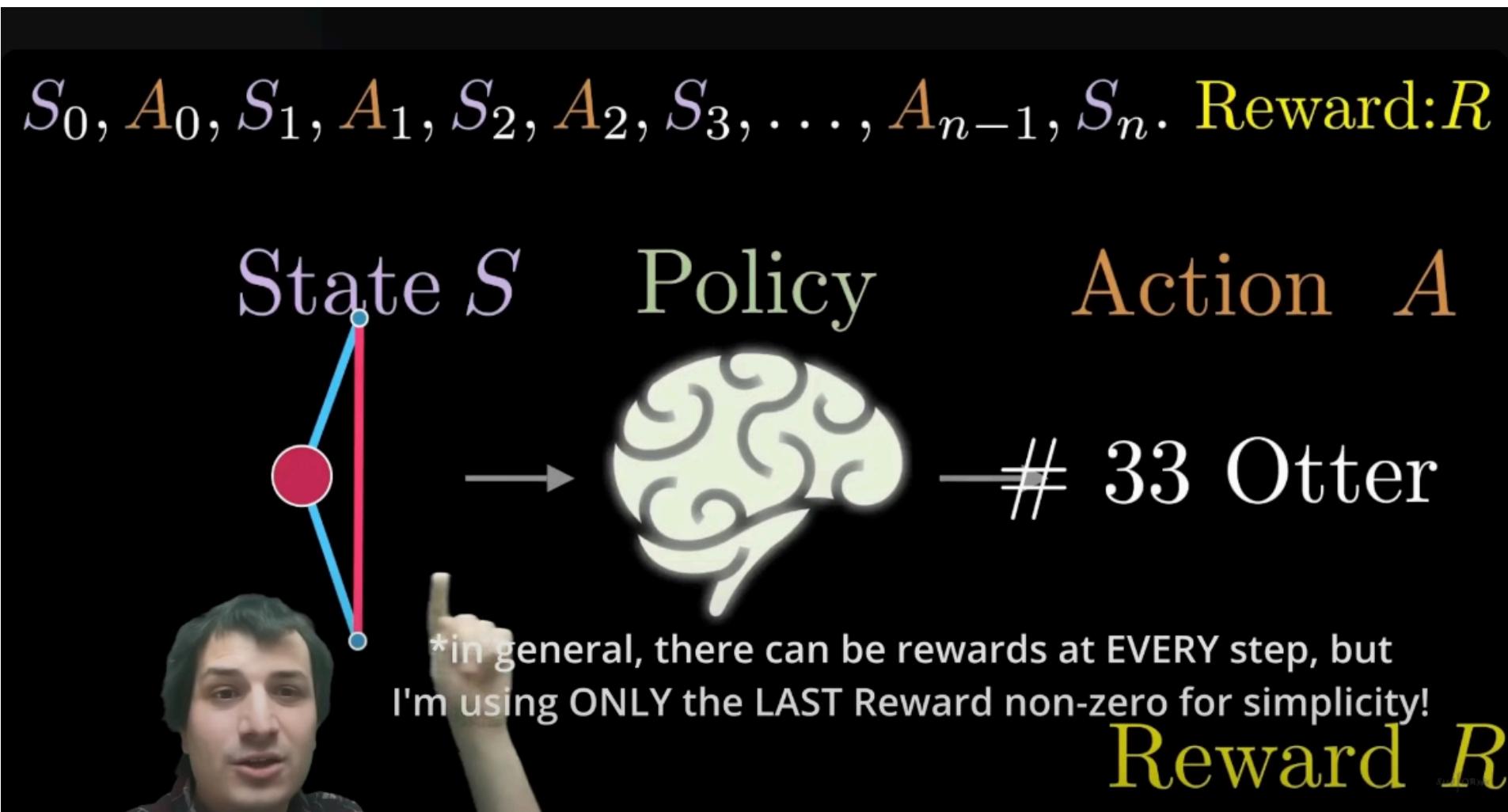
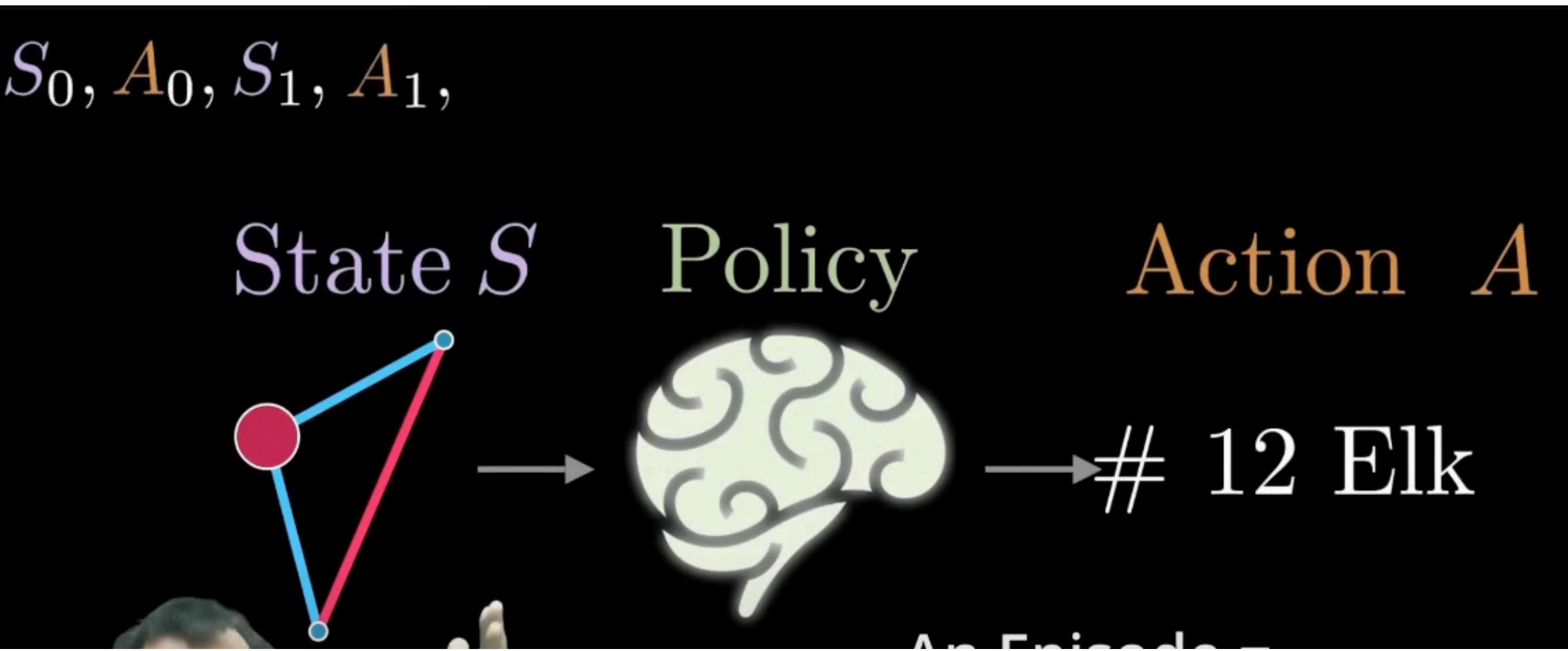


$S_0, A_0,$

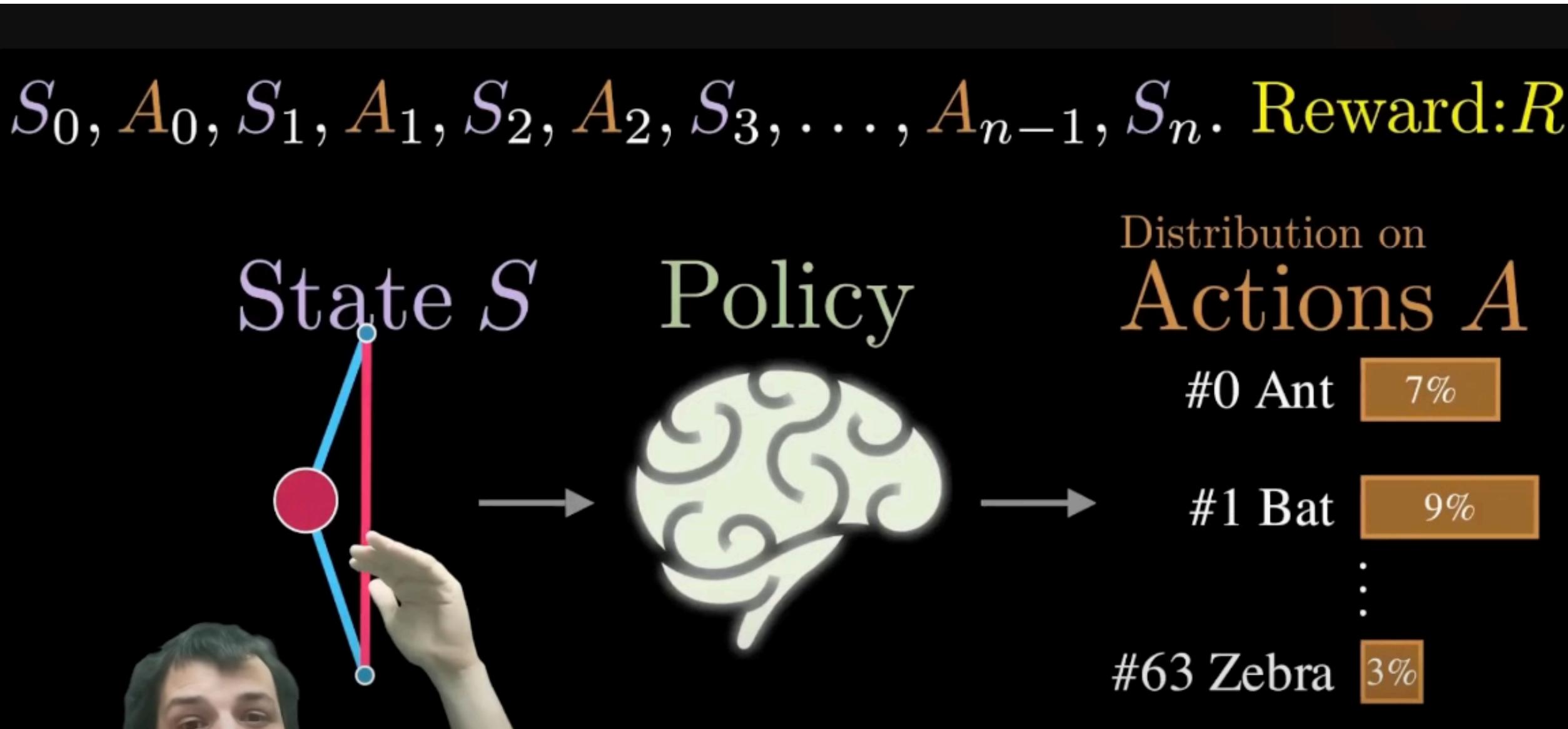


$S_0, A_0, S_1,$





$$\begin{aligned} \text{Reward } R \\ = \Delta X - |\Delta Y| \end{aligned}$$



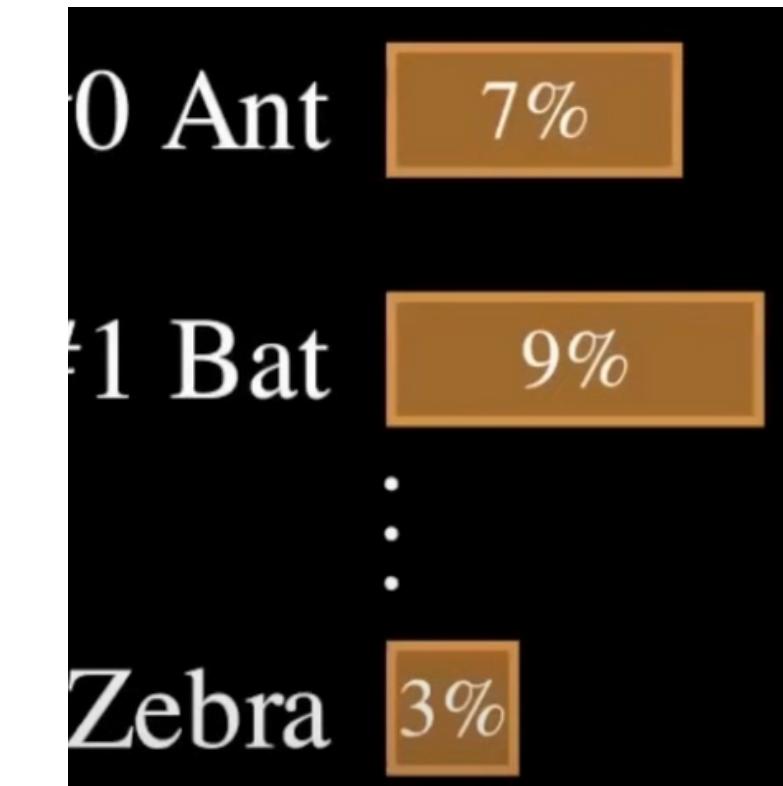
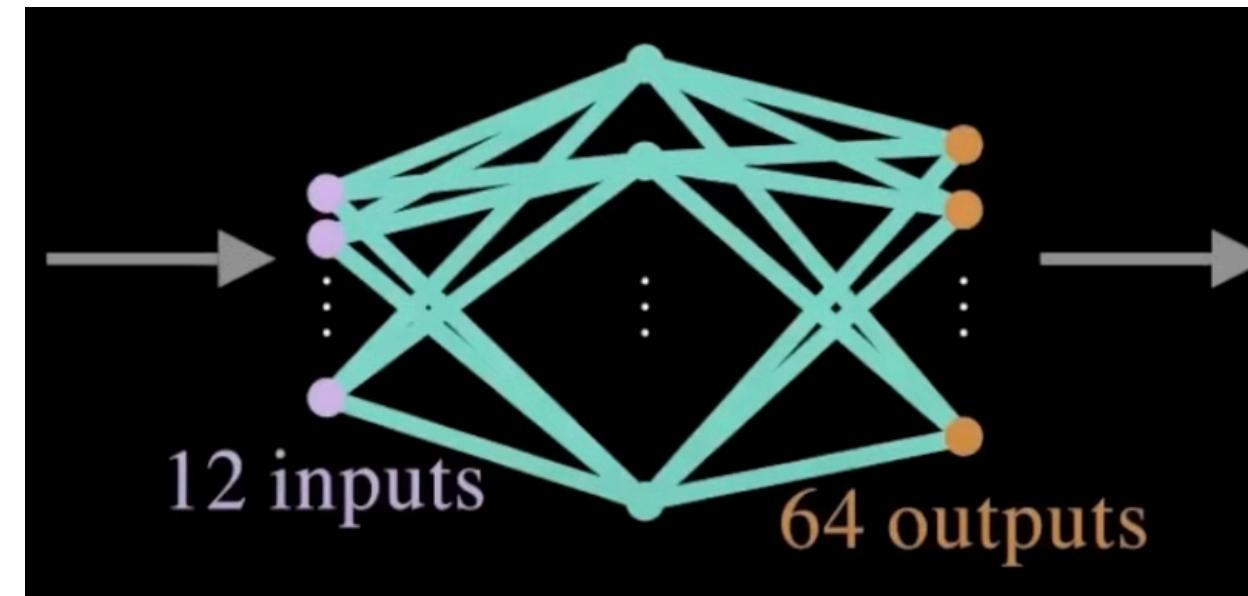
$$\begin{aligned} \text{Reward } R \\ = \Delta X - |\Delta Y| \end{aligned}$$

$S_0, A_0, S_1, A_1, S_2, A_2, S_3, \dots, A_{n-1}, S_n$. Reward: R

$\pi(A|S, \theta) = \mathbb{P}(\text{Choose } A \text{ when at } S, \text{ params} = \theta)$

Policy π

State S



↑
Params θ

Reward R
 $= \Delta X - |\Delta Y|$

The REINFORCE Algorithm

How to improve θ ?

- Observe $S_0, A_0, S_1, A_1, \dots, S_n$ and R
- If R is large, then tweak θ to make actions A_t more likely. i.e. $\uparrow \pi(A_t|S_t, \theta)$



State S
Params θ
Policy π
Actions A
 $\pi(A|S, \theta) = \mathbb{P}$
Reward R

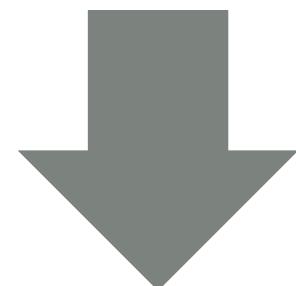
The REINFORCE Algorithm

How to improve θ ?

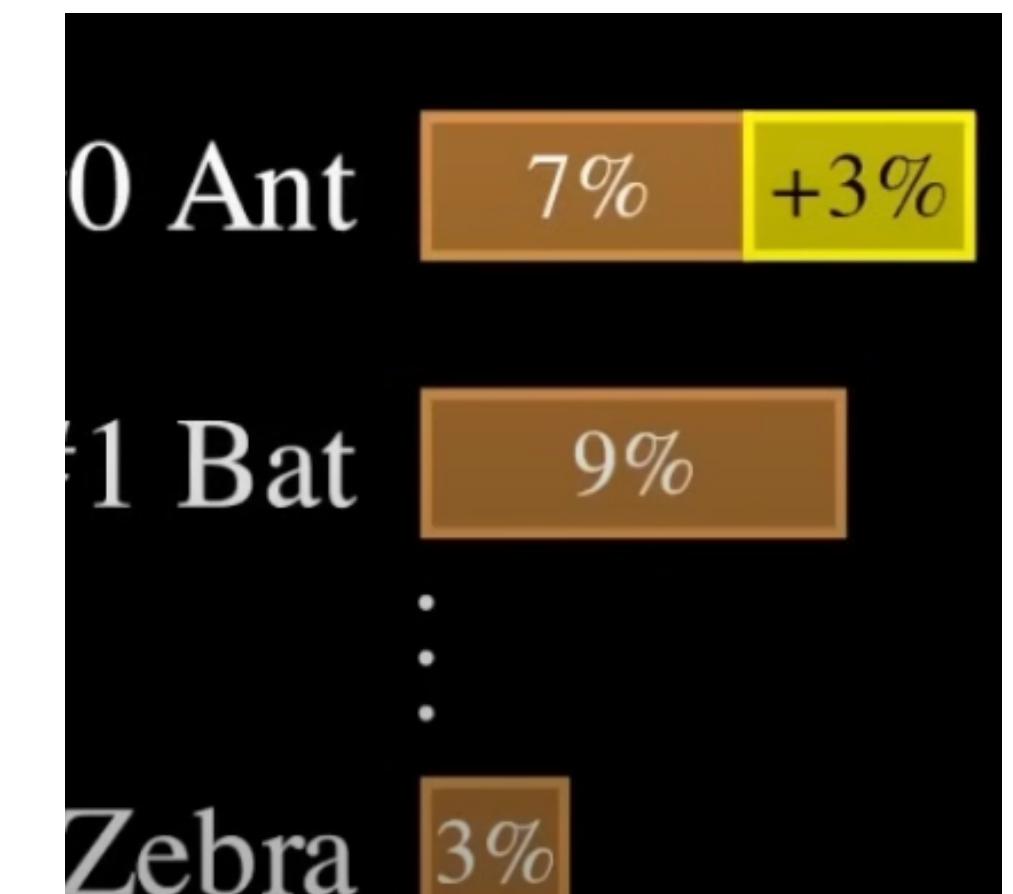
- Observe $S_0, A_0, S_1, A_1, \dots, S_n$ and R
- If R is large, then tweak θ to make actions A_t more likely. i.e. $\uparrow \pi(A_t | S_t, \theta)$

State S
Params θ
Policy π
Actions A
 $\pi(A|S, \theta) = \mathbb{P}$
Reward R

$$\nabla_{\theta} \log \pi(A_t | S_t, \theta)$$



$$R \nabla_{\theta} \log \pi(A_t | S_t, \theta)$$



```
1 def log_policy_pi(params_theta, state_S, action_A):
2     # "forward" is the neural network a.k.a. "the brain"
3     log_probs = forward(params_theta, state_S) #log-probs of ALL possible actions
4     return log_probs[action_A] #Plug in the action we took!
5
6 #Do the derivative using JAX. (Does gradient w.r.t. the 0th input by default)
7 grad_log_pi = jax.grad(log_policy_pi)
8
9 #Update the parameters!
10 params_theta += reward_R * grad_log_pi(params_theta, state_S, action_A)
```

$$\theta_{\text{new}} = \theta + R \nabla_{\theta} \log \pi(A_t | S_t, \theta)$$

Avg. Reward: $\Delta X - |\Delta Y|$ (Body Lengths Per Action)



Improvement: The GRPO Algorithm

How to improve θ ? Reward size R is relative!

- Observe $S_0, A_0, S_1, A_1, \dots, S_n$ and R
- If R is large, then tweak θ to make actions A_t more likely. i.e. $\uparrow \pi(A_t|S_t, \theta)$

$$\theta_{\text{new}} = \theta + R \nabla_{\theta} \log \pi(A_t|S_t, \theta)$$

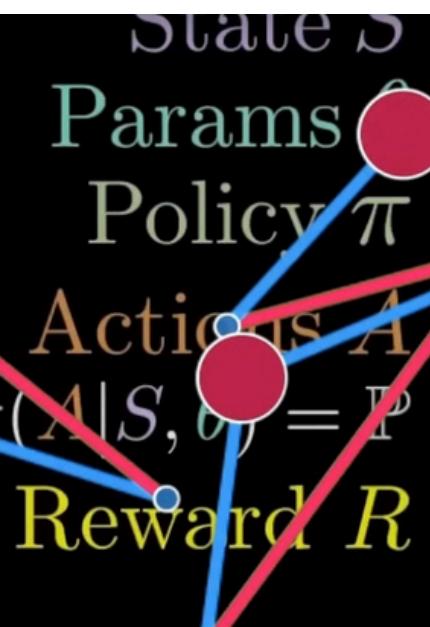
Improvement: The GRPO Algorithm

How to improve θ ? Reward size R is relative!

- Observe $S_0, A_0^{(1)}, S_1^{(1)}, A_1^{(1)}, \dots, S_n^{(1)}$ and $R^{(1)}$
- Observe $S_0, A_0^{(2)}, S_1^{(2)}, A_1^{(2)}, \dots, S_n^{(2)}$ and $R^{(2)}$
- Observe $S_0, A_0^{(3)}, S_1^{(3)}, A_1^{(3)}, \dots, S_n^{(3)}$ and $R^{(3)}$

$$\theta_{\text{new}} = \theta + R \nabla_{\theta} \log \pi(A_t|S_t, \theta)$$

State S
Params θ
Policy π
Actions A
 $\pi(A|S, \theta) = \mathbb{P}$
Reward R



Improvement: The GRPO Algorithm

How to improve θ ? Reward size R is relative!

- Observe $S_0, A_0^{(1)}, S_1^{(1)}, A_1^{(1)}, \dots, S_n^{(1)}$ and $R^{(1)}$
 - Observe $S_0, A_0^{(2)}, S_1^{(2)}, A_1^{(2)}, \dots, S_n^{(2)}$ and $R^{(2)}$
 - Observe $S_0, A_0^{(3)}, S_1^{(3)}, A_1^{(3)}, \dots, S_n^{(3)}$ and $R^{(3)}$
- $$\theta_{\text{new}} = \theta + \text{Advantage}^{(1)} \nabla_{\theta} \log \pi(A_t^{(1)} | S_t^{(1)}, \theta)$$

$$\text{Advantage}^{(1)} = \frac{R^{(1)} - \text{avg}\{R^{(1)}, R^{(2)}, R^{(3)}\}}{\text{std}\{R^{(1)}, R^{(2)}, R^{(3)}\}}$$

