# RL: promises and challenges

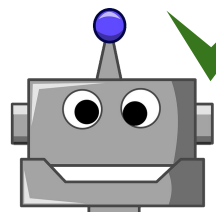# RL: promises and challenges
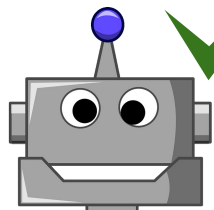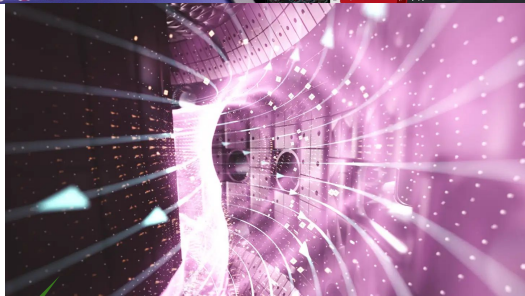


Reward is enough for intelligence?

David Silver 👤 ✉, Satinder Singh, Doina Precup, Richard S. Sutton
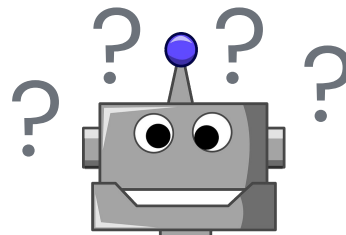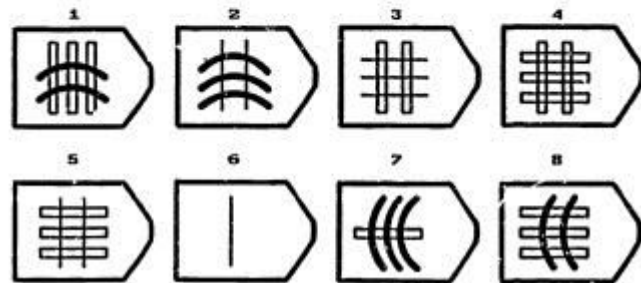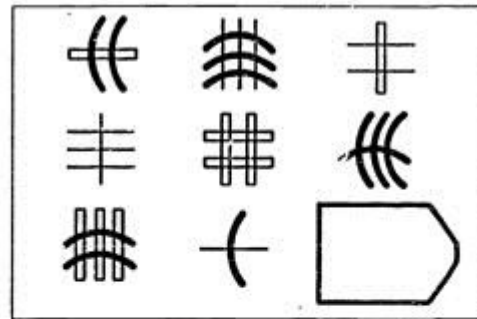
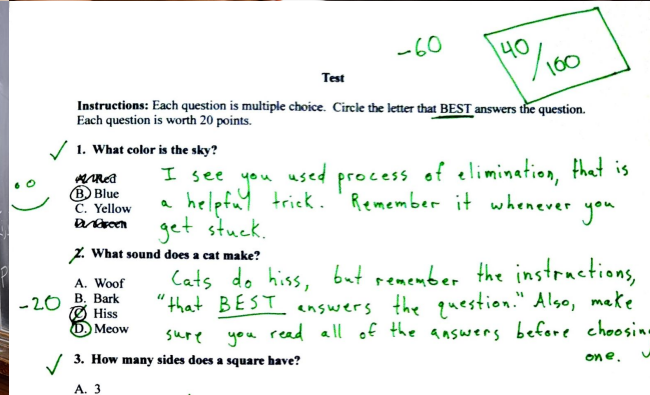# RL: promises and challenges



Reward is enough for intelligence?

David Silver, Satinder Singh, Doina Precup, Richard S. Sutton

# What is different about human learning?

# What is different about human learning?

# What is different about human learning?
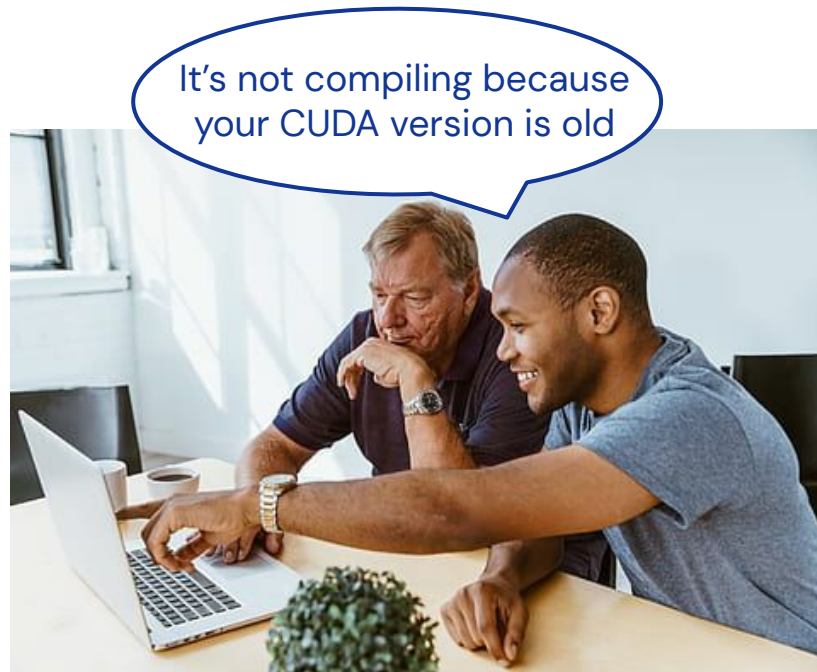


Explanation along with reward!
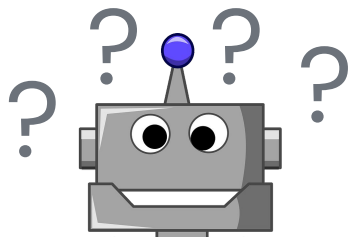
# What is an explanation?

Explanations are intended to *communicate* the links between:

- Concrete situation
- Abstract principles which are:
  - Generalizable
  - Causal

Explanations help us to learn and generalize abstract tasks!

(Lombrozo, 2006; Lombrozo & Carey, 2006; Woodward, 2003)

# Could agents learn and generalize better if trained with explanations?

# Odd-one-out task: abstraction and relational reasoning

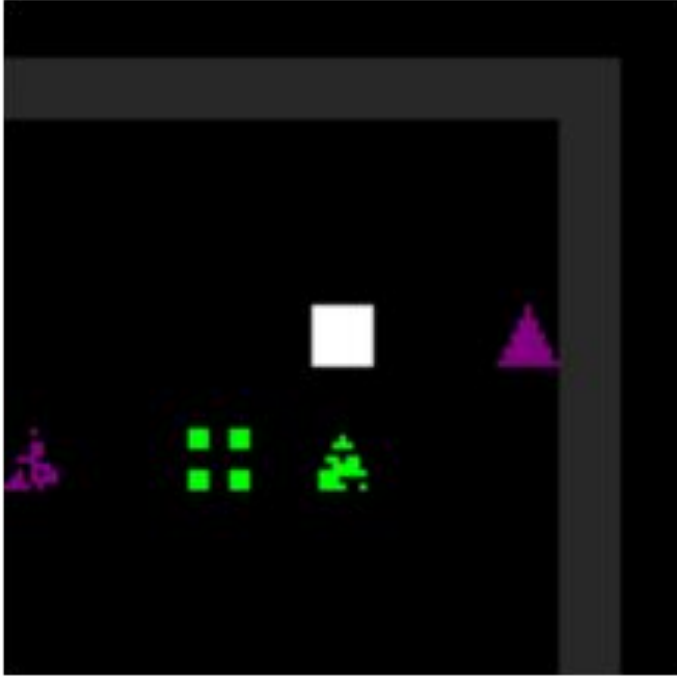# Odd-one-out task: abstraction and relational reasoning



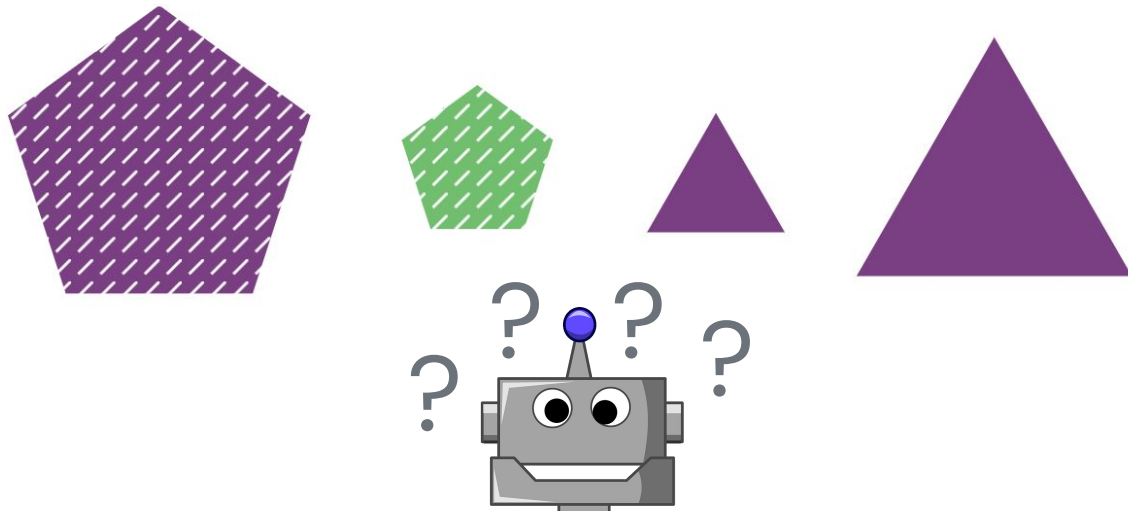Proper subsets don't reveal the answer!

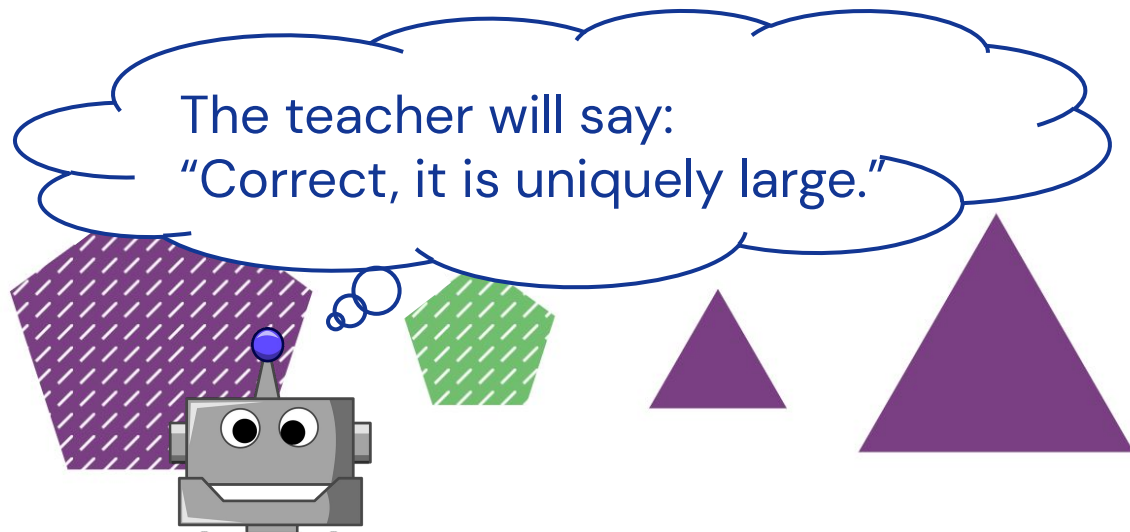A challenging credit assignment problem from reward alone.

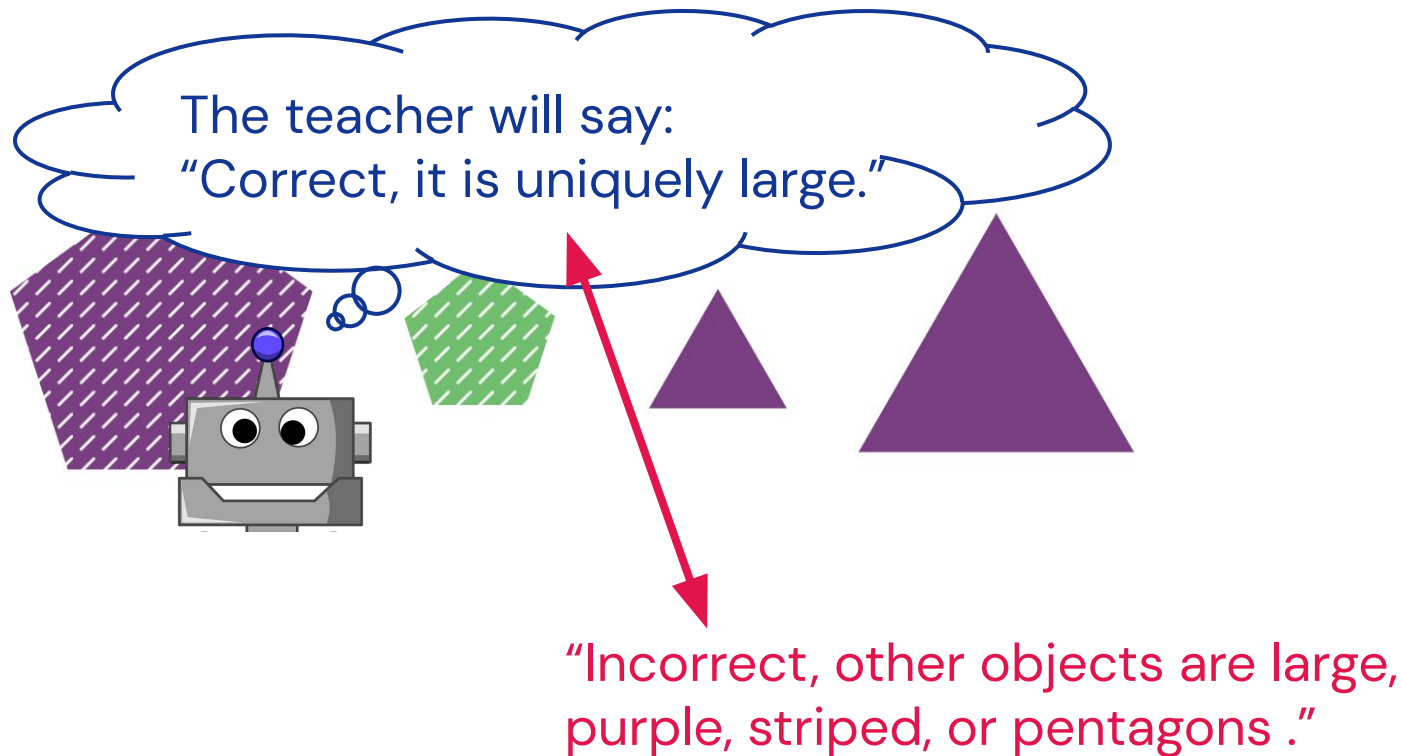# Instantiated in 2D and 3D environments

# Agents struggle to learn odd-one-out tasks from rewards alone: abstractions and relations are hard to infer!
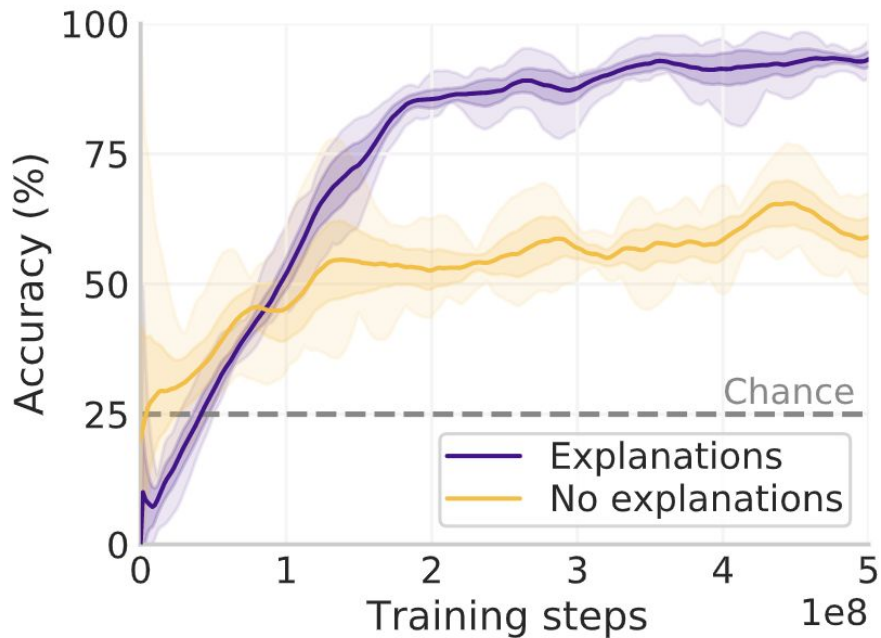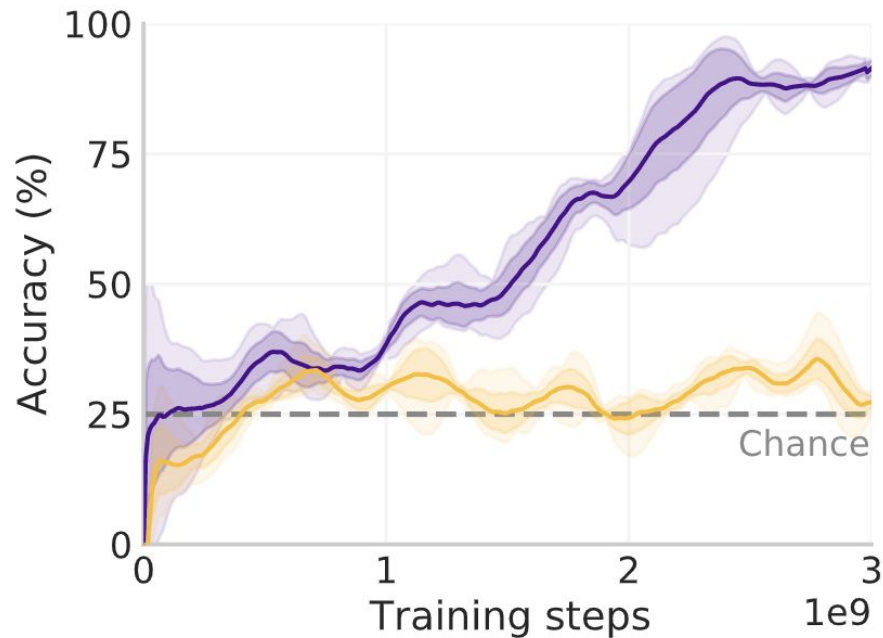
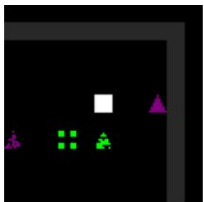# Predicting explanations during training

# Predicting explanations during training

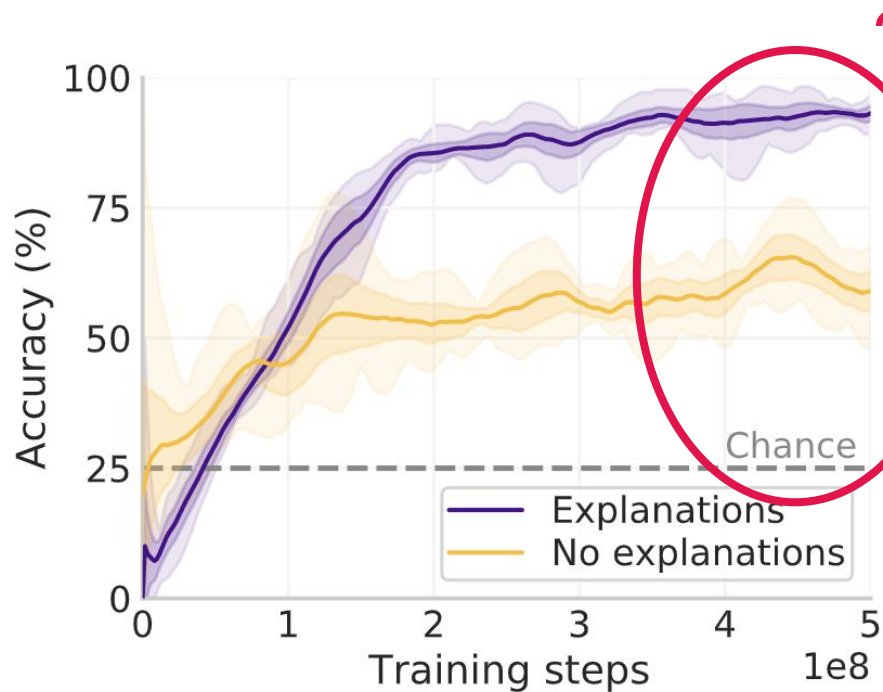# Explanations improve learning of odd-one-out tasks
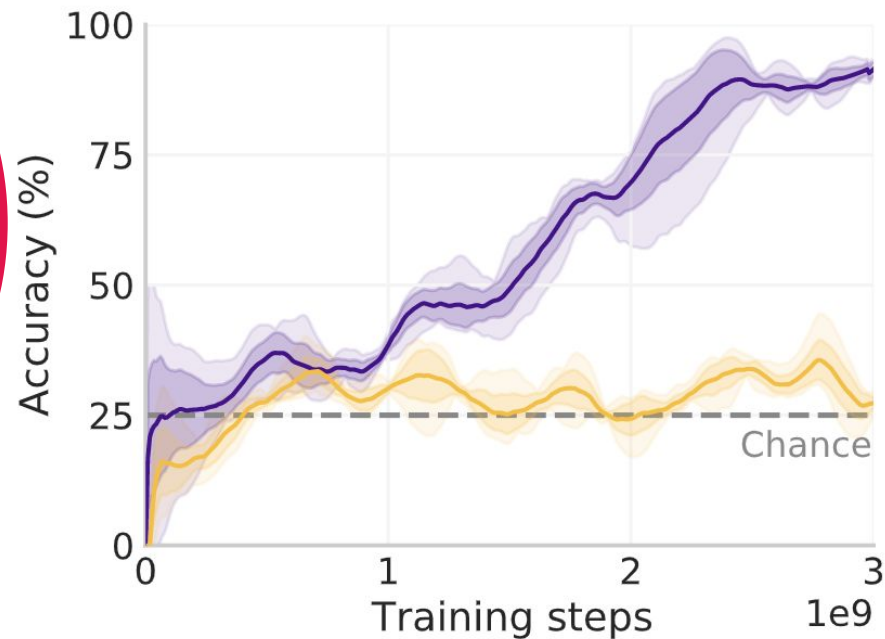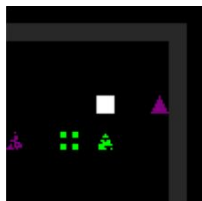


(b) 2D results.

(c) 3D results.

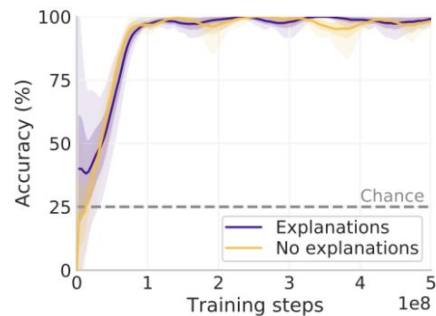# Explanations improve learning of odd-one-out tasks
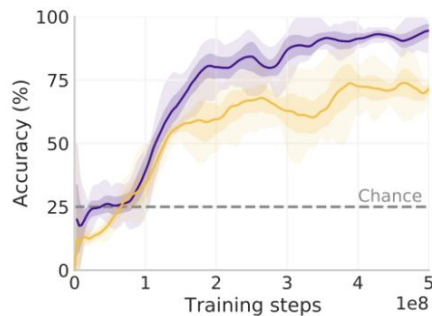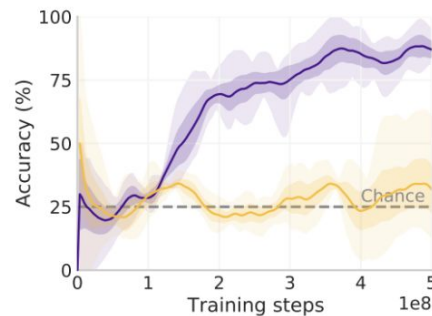


(b) 2D results.

(c) 3D results.

# In 2D, explanations help agents move past focus on "shortcut" features to learn more difficult ones



(a) Position.  (b) Color.  (c) Shape.  (d) Texture.

Easier ──────────────────► Harder

For CNNs

(cf. Hermann & Lampinen, 2020; Geirhos et al., 2020)

# Explanations can help with other important challenges (see paper)

# Explanations can help with other important challenges (see paper)

Explanations can overcome ambiguous training to shape how an agent generalizes OOD!

**Train (confounded):**



**Evaluation (deconfounded):**

# Explanations can help with other important challenges (see paper)

Explanations can overcome ambiguous training to shape how an agent generalizes OOD!

Explanations allow agents to meta-learn how to perform experiments to identify causality!

# Explanations can help with other important challenges (see paper)

Explanations can overcome ambiguous training to shape how an agent generalizes OOD!
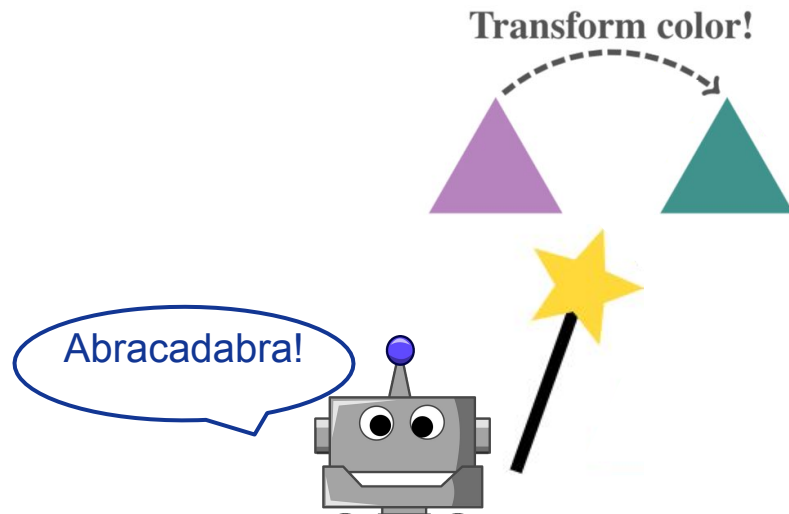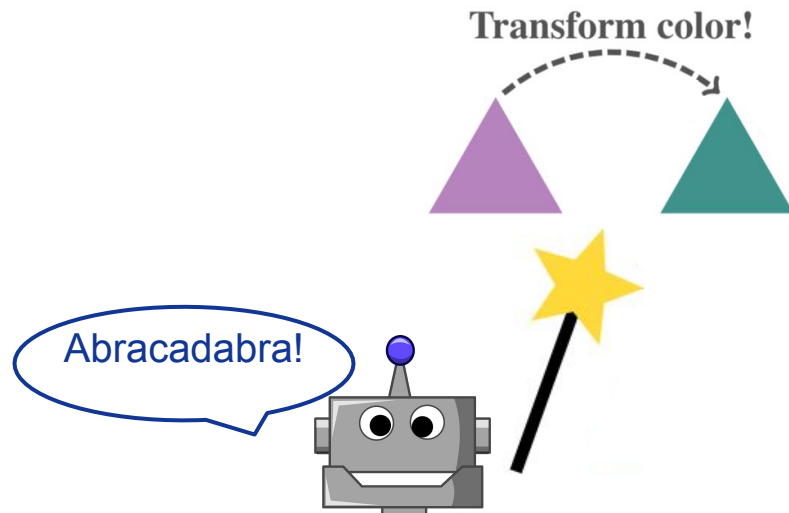
Explanations allow agents to meta-learn how to perform experiments to identify causality!



+ see paper for more analysis, control conditions like unsupervised auxiliary losses, ...