# Multimodal Machine Learning

## Lecture 7.2: Reasoning 3
## Inference + Knowledge

Paul Liang

**Language Technologies Institute**

**Carnegie Mellon University**

# Midterm Project Report Instructions

- **Goal:** Evaluate state-of-the-art models on your dataset and identify key issues through a detailed error analysis
  - It will inform the design of your new research ideas
- **Report format:** 2 column (ICML template)
  - The report should follow a similar structure to a research paper
  - Teams of 3: 8 pages, Teams of 4: 8.5 pages, Teams of 5: 10 pages. Teams of 6: 10.5 pages
- **Number of SOTA models**
  - Teams of 3 or 4 should have at least two baseline models
  - Teams of 5 or 6 should have at least three baseline models
- **Error analysis**
  - This is one of the most important part of this report. You need to understand where previous models can be improved.

# Examples of Possible Error Analysis Approaches

- Dataset-based:
    - Split correct/incorrect by label
    - Manually inspect the samples that are incorrectly predicted
        - What are the commonalities?
        - What are differences with the correct ones?
    - Sub-dataset analysis: length of question, rare words, cluttered images, high frequency in signals?

# Examples of Possible Error Analysis Approaches

- Perturbation-based:
  - Make targeted changes to specific parts of the image.
  - Change one word/paraphrase/add redundant tokens.
  - See whether the model remains robust

# Examples of Possible Error Analysis Approaches

- Model-based:
  - Visualize feature attributions: LIME, $1^{st}/2^{nd}$ order gradients
  - Ablation studies to understand what model components are important
- Theory-based:
  - Write out the math! From optimization and learning perspective, does the model do what's expected?
  - Some useful tools: consider linear case/other simplest case and derive solution, do empirical sanity checks first.

[Liang et al., MultiViz: An Analysis Benchmark for Visualizing and Understanding Multimodal Models. arXiv 2022]

# Examples of Possible Error Analysis Approaches

# ON THE CONVERGENCE OF ADAM AND BEYOND

**Sashank J. Reddi, Satyen Kale & Sanjiv Kumar**
Google New York
New York, NY 10011, USA
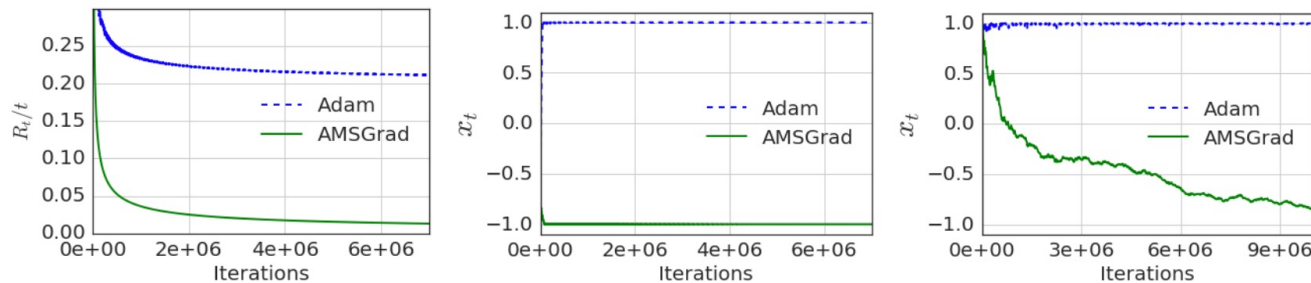{sashank,satyenkale,sanjivk}@google.com

Figure 1: Performance comparison of ADAM and AMSGRAD on synthetic example on a simple one dimensional convex problem inspired by our examples of non-convergence. The first two plots (left and center) are for the online setting and the the last one (right) is for the stochastic setting.

[Reddi et al., On the Convergence of Adam and Beyond. ICLR 2018]

# Examples of Possible Error Analysis Approaches

**Finding:** Image captioning models capture spurious correlations between gender and generated actions



You'll see more in today's reasoning lecture and in quantification lectures

[Hendricks et al., Women also Snowboard: Overcoming Bias in Captioning Models. ECCV 2018]

# Midterm Project Report Instructions

Main report sections:

- Abstract
- Introduction
- Related work
- Problem statement
- Multimodal baseline models
- Experimental methodology
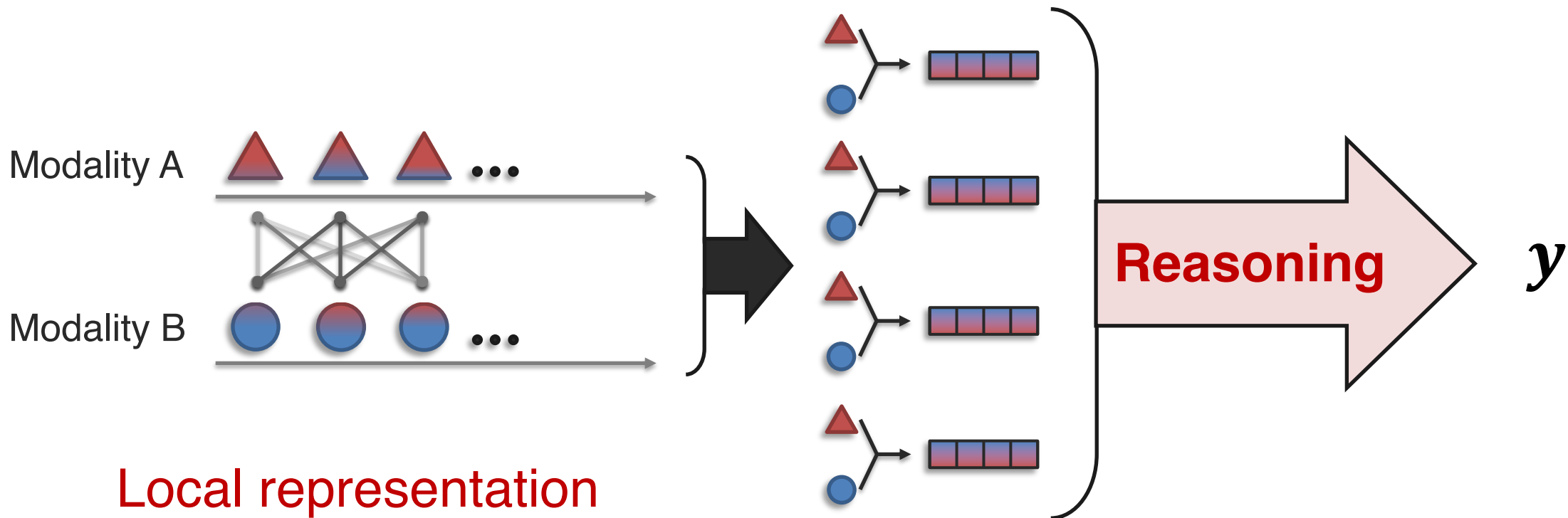- Results and discussion
- New research ideas

The structure is similar to a research paper submission ☺

# Upcoming Deadlines

- Monday October 31$^{st}$ 8pm: Midterm report deadline
- Tuesday and Thursday (11/1 and 11/3): midterm presentations
  - All students are expected to attend both presentation sessions in person
  - Each team will present either Tuesday or Thursday
  - The focus of these presentations is about your research ideas
  - Feedback will be given by all students, instructors and TAs

# Reasoning

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.



Modality A

Modality B

**Reasoning**

$y$

Local representation

+ Aligned representation

# Reasoning

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.
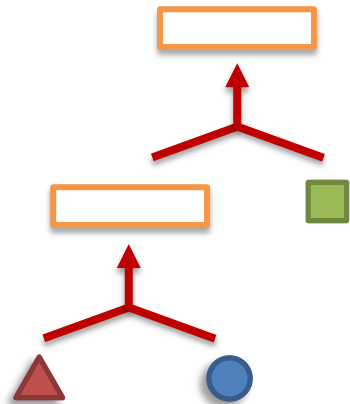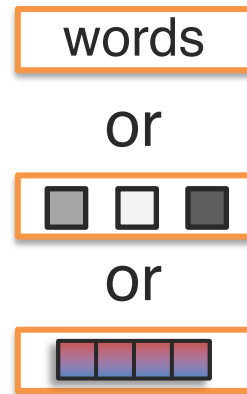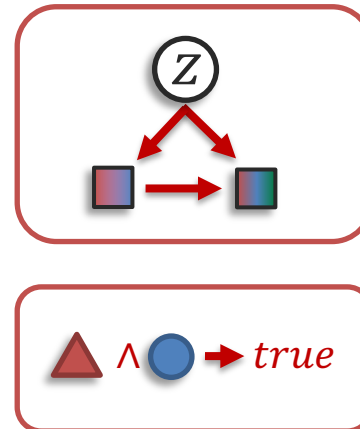
# Summary

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.

| | (A) Structure modeling | (B) Intermediate concepts | (C) Inference paradigm | (D) External knowledge |
|---|---|---|---|---|
| **Last Thursday** | Temporal Hierarchical | Continuous | | |
| **Tuesday** | Interactive | | | |
| **Today** | Discovery | Discrete | Causal Logical | Knowledge Commonsense |

# Sub-Challenge 3a: Structure Modeling

Carnegie Mellon University

# Interactive Structure

**Structure defined through interactive environment**
Main difference from temporal - actions taken at previous time steps affect future states

Integrates multimodality into the reinforcement learning framework



$a_1$  $a_2$  $a_3$  $a_T$

$s_1 =$  $s_2 =$  $s_3 =$  $s_T =$

Time

$t = 1$  $t = 2$  $t = 3$  ...  $t = T$

[Luketina et al., A Survey of Reinforcement Learning Informed by Natural Language. IJCAI 2019]
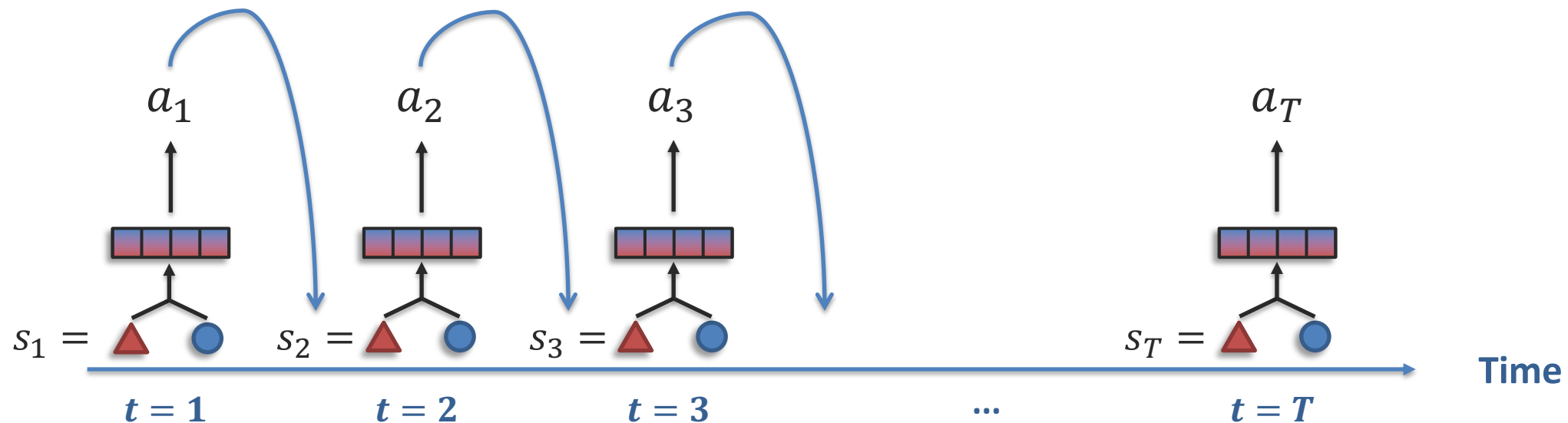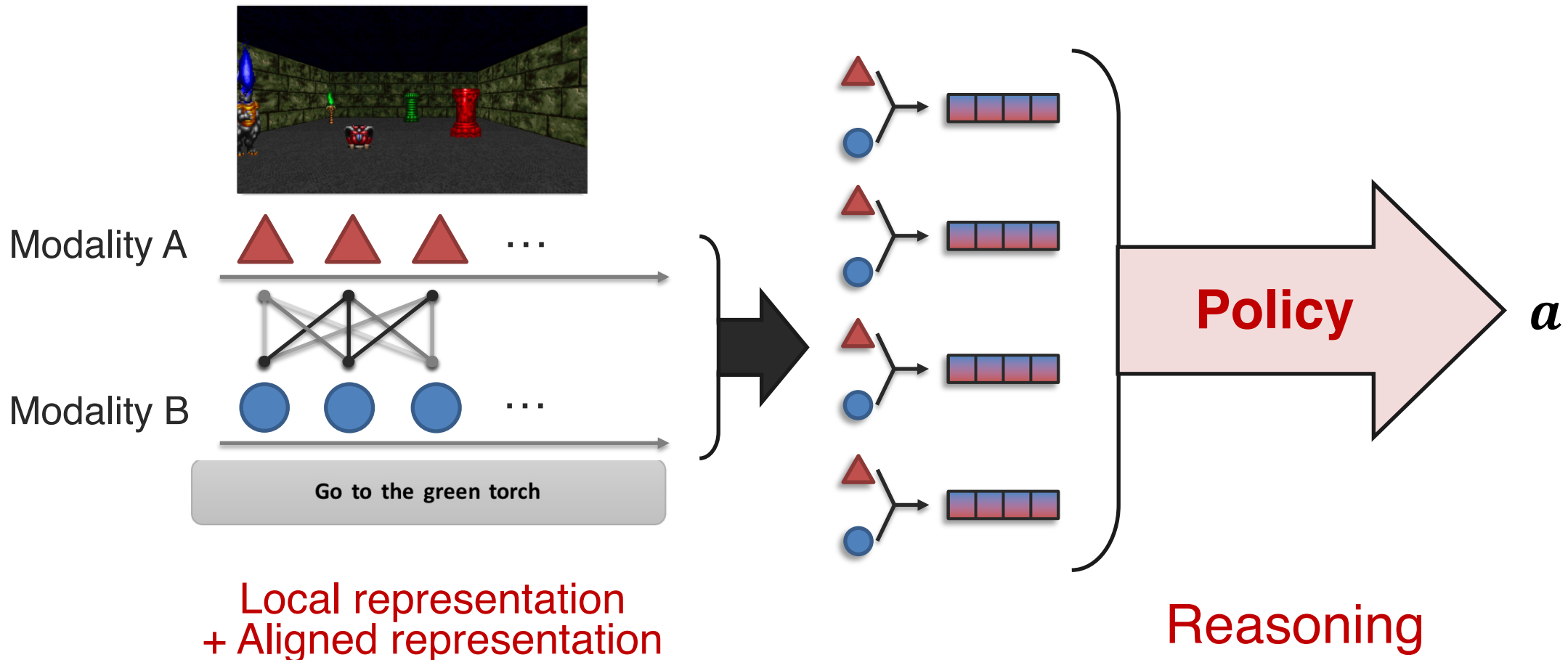
# Interactive Structure

**Structure defined through interactive environment**
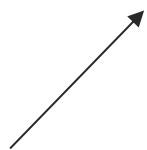Main difference from temporal - actions taken at previous time steps affect future states



Modality A

Modality B

Go to the green torch

Local representation
+ Aligned representation

**Policy**

$a$

Reasoning

# Summary: Exact Methods

Fully known MDP
states
transitions
rewards

Bellman
optimality
equations

$Q^*(s, a)$ Q-value iteration

$V^*(s)$ Value iteration



$$Q^*(s, a) = \mathbb{E}_{s'}\left[r(s, a, s') + \gamma V^*(s')\right]$$

# Summary: Exact Methods
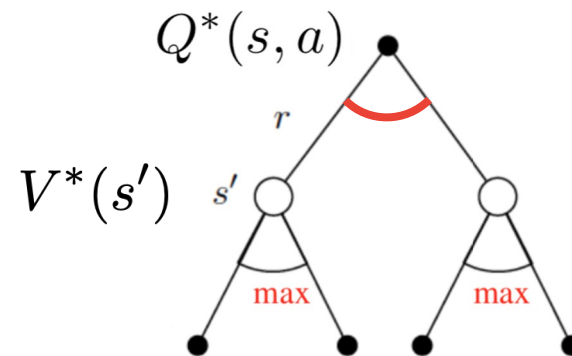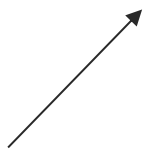
Fully known MDP
states
transitions
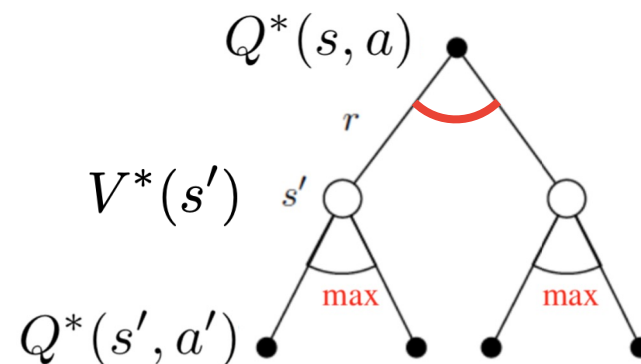rewards

Bellman
optimality
equations

$Q^*(s,a)$ **Q-value iteration**

$V^*(s)$ **Value iteration**

$$Q^*(s,a) = \mathbb{E}_{s'}\left[r(s,a,s') + \gamma V^*(s')\right]$$
$$= \mathbb{E}_{s'}\left[r(s,a,s') + \gamma \max_{a'} Q^*(s',a')\right]$$

# Summary: Exact Methods

Bellman
optimality
equations

$Q^*(s, a)$   **Q-value iteration**

$V^*(s)$   **Value iteration**

Fully known MDP
states
transitions
rewards

$Q^*$(s, a) = expected utility starting in s, taking action a, and (thereafter) acting optimally

Bellman Equation:

$$Q^*(s, a) = \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma \max_{a'} Q^*(s', a'))$$

Q-Value Iteration:

$$Q^*_{k+1}(s, a) \leftarrow \sum_{s'} P(s'|s, a)(R(s, a, s') + \gamma \max_{a'} Q^*_k(s', a'))$$

$Q^*(s, a)$

$r$

$V^*(s')$   $s'$

max   max

$Q^*(s', a')$

$$Q^*(s, a) = \mathbb{E}_{s'}\left[r(s, a, s') + \gamma V^*(s')\right]$$

$$= \mathbb{E}_{s'}\left[r(s, a, s') + \gamma \max_{a'} Q^*(s', a')\right]$$

$$= \sum_{s'} p(s'|s, a)\left(r(s, a, s') + \gamma \max_{a'} Q^*(s', a')\right)$$

# Summary: Exact Methods

Fully known MDP
states
transitions
rewards

Bellman
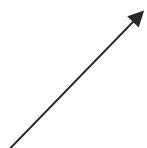optimality
equations

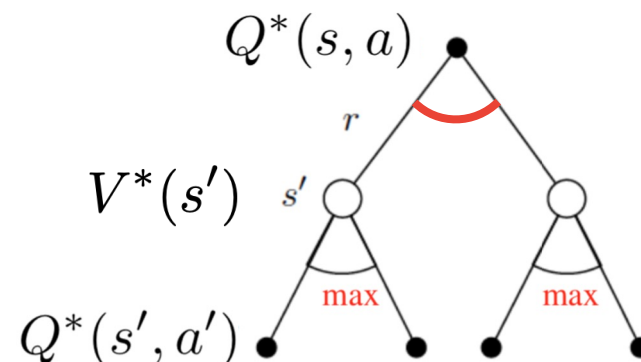$Q^*(s, a)$    **Q-value iteration**

$V^*(s)$    **Value iteration**

Bellman
expectation
equations

$Q^\pi(s, a)$    **Q-policy iteration**

$V^\pi(s)$    **Policy iteration**

**Repeat until policy converges. Guaranteed to converge to optimal policy.**

$Q^*(s, a)$

$V^*(s)$   $s'$   $r$

$Q^*(s', a')$   max   max

$$Q^*(s, a) = \mathbb{E}_{s'}\left[r(s, a, s') + \gamma V^*(s')\right]$$

$$= \mathbb{E}_{s'}\left[r(s, a, s') + \gamma \max_{a'} Q^*(s', a')\right]$$

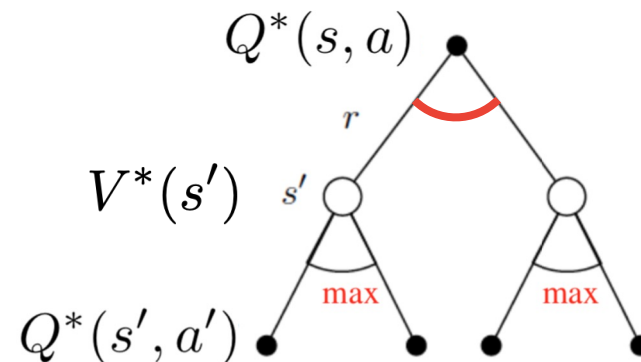$$= \sum_{s'} p(s'|s, a)\left(r(s, a, s') + \gamma \max_{a'} Q^*(s', a')\right)$$

**Limitations:**
**Iterate over and storage for all states and actions: requires small, discrete state and action space**
**Update equations require fully observable MDP and known transitions**

# Summary: Tabular Q-learning

MDP
with unknown
transitions

$\longrightarrow$

Bellman
optimality
equations

$\longrightarrow$

Replace true expectation
over transitions with
estimates

**Tabular Q-learning**

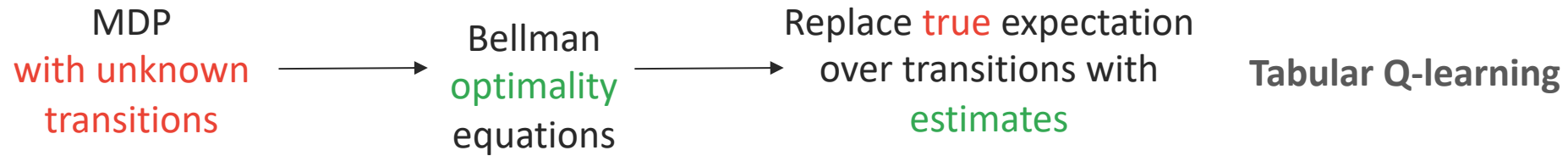$s' \sim P(s'|s,a)$    simulation and exploration, epsilon greedy is important!

Poor estimates of Q(s,a) at the start:

Bad initial estimates in the first few cases can drive policy into sub-optimal region, and never explore further.

$$\pi(s) = \begin{cases} \max_a \hat{Q}(s,a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{otherwise} \end{cases}$$

Gradually decrease epsilon as policy is learned.
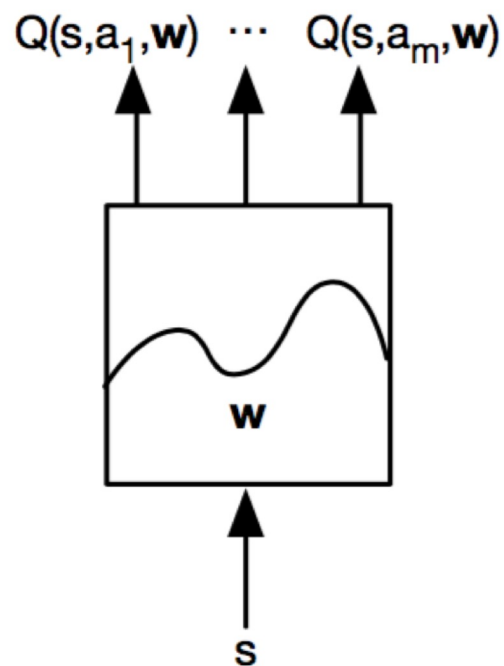
# Summary: Tabular Q-learning

MDP
with unknown
transitions

$\longrightarrow$

Bellman
optimality
equations

$\longrightarrow$

Replace true expectation
over transitions with
estimates

**Tabular Q-learning**

$$s' \sim P(s'|s,a)$$    simulation and exploration, epsilon greedy is important!

$$Q^*(s,a) = \mathbb{E}_{s'}\left[\underline{r(s,a,s') + \gamma \max_{a'} Q^*(s',a')}\right]$$

**old estimate**                          **target**

$$Q_{k+1}(s,a) \leftarrow Q_k(s,a) + \alpha\left(r(s,a,s') + \gamma \max_{a'} Q_k(s',a') - Q_k(s,a)\right)$$

**Tabular: keep a |S| x |A| table of Q(s,a)**
**Still requires small and discrete state and action space**
**How can we generalize to unseen states?**

# Summary: Deep Q-learning

$$Q^*(s,a) = \mathbb{E}_{s'} \left[ r(s,a,s') + \gamma \max_{a'} Q^*(s',a') \right]$$

**old estimate**           **target**

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \left( \underbrace{r + \gamma \max_{a'} Q(s',a'; w_i^-)}_{\text{Q-learning target}} - \underbrace{Q(s,a; w_i)}_{\text{Q-network}} \right)^2 \right]$$

Q(s,a₁,**w**) ··· Q(s,aₘ,**w**)

**w**

s

**Stochastic gradient descent + Experience replay + Fixed Q-targets**

**Works for high-dimensional state and action spaces**
**Generalizes to unseen states**

# Can we Directly Learn the Policy?

- Often $\pi$ can be simpler than Q or V

  - E.g., robotic grasp

  **Q(s,a) and V(s) very high-dimensional**
  **But policy could be just 'open/close hand'**

- V: doesn't prescribe actions

  - Would need dynamics model (+ compute 1 Bellman back-up)

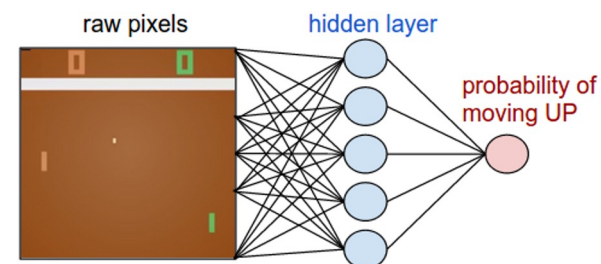- Q: need to be able to efficiently solve $\arg\max_a Q^*(s,a)$

  - Challenge for continuous / high-dimensional action spaces

$$\pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg\max_a \mathbb{E}_{s'}\left[r(s,a,s') + \gamma V^*(s')\right] \\ \epsilon, & \text{else} \end{cases} \qquad \pi^*(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg\max_a Q^*(s,a) \\ \epsilon, & \text{else} \end{cases}$$

[Slides from Fragkiadaki, 10-703 CMU]

# Summary: Policy Gradients

$\pi(a \mid s)$

raw pixels | hidden layer

probability of moving UP

1. Initialize a policy network at random
2. **Repeat Forever:**
3. Collect a bunch of rollouts with the policy  **epsilon greedy!**
4. Increase the probability of actions that worked well
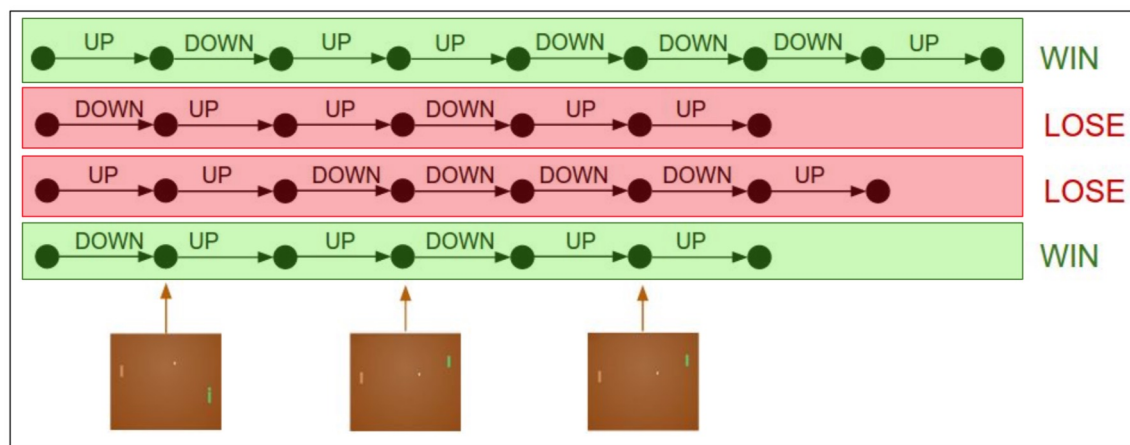
Pretend every action we took here was the correct label.

maximize: $\log p(y_i \mid x_i)$

Pretend every action we took here was the wrong label.

maximize: $(-1) * \log p(y_i \mid x_i)$

$$\sum_i A_i * \log p(y_i \mid x_i)$$

**Does not require transition probabilities**
**Does not estimate Q(), V()**
**Predicts policy directly**

UP → DOWN → UP → UP → DOWN → DOWN → DOWN → UP → WIN

DOWN → UP → UP → DOWN → UP → UP → LOSE

UP → UP → DOWN → DOWN → DOWN → DOWN → UP → LOSE

DOWN → UP → UP → DOWN → UP → UP → WIN

[Slides from Karpathy]

# Summary: Policy Gradients

Gradient estimator:

$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$$

**Interpretation:**

- If **r(trajectory)** is high, push up the probabilities of the actions seen
- If **r(trajectory)** is low, push down the probabilities of the actions seen

**REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)**

Input: a differentiable policy parameterization $\pi(a|s,\boldsymbol{\theta}), \forall a \in \mathcal{A}, s \in \mathcal{S}, \boldsymbol{\theta} \in \mathbb{R}^n$

Initialize policy weights $\boldsymbol{\theta}$

Repeat forever:

    Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$ following $\pi(\cdot|\cdot, \boldsymbol{\theta})$

    For each step of the episode $t = 0, \ldots, T-1$:

        $G_t \leftarrow$ return from step $t$

        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G_t \nabla_{\boldsymbol{\theta}} \log \pi(A_t | S_t, \boldsymbol{\theta})$

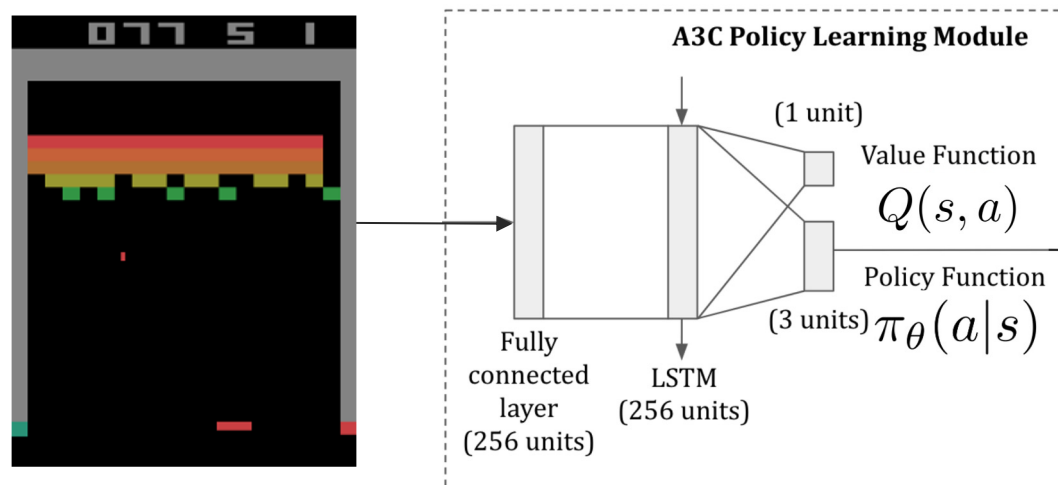**epsilon greedy**

# Summary: Actor-Critic Methods

**Problem:** The raw reward of a trajectory isn't necessarily meaningful. E.g. if all rewards are positive, you keep pushing up probabilities of all actions.
**What is important then?** Whether a reward is higher or lower than what you expect to get.

**Yes,** using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q function)

**Exploration + experience replay**
**Decorrelate samples**
**Critic: evaluates how good the action is** **Fixed targets**



**A3C Policy Learning Module**

(1 unit)
Value Function
$Q(s, a)$

(3 units) $\pi_\theta(a|s)$
Policy Function

Fully connected layer (256 units)
LSTM (256 units)

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s'\sim\mathcal{D}_i}\left[\left(r + \gamma \max_{a'} Q(s', a'; w_i^-) - Q(s, a; w_i)\right)^2\right]$$

Q-learning target          Q-network

$\pi_\theta(a|s)$

**Actor: decides what actions to take**

$$\nabla_\theta J(\theta) \approx \sum_{t\geq 0} \left(Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)\right) \nabla_\theta \log \pi_\theta(a_t|s_t)$$

**Variance reduction with a baseline**

[Minh et al., Asynchronous Methods for Deep Reinforcement Learning. ICML 2016]

# Summary: RL Methods

**Epsilon greedy + exploration**
**Experience replay**
**Decorrelate samples**
**Fixed targets**

**Value iteration**
**Policy iteration**
**(Deep) Q-learning**

**Policy gradients**

**Variance reduction with a baseline**

**Actor (policy)**
**Critic (Q-values)**

‣ Value Based
- Learned Value Function
- Implicit policy (e.g. ε-greedy)

‣ Policy Based
- No Value Function
- Learned Policy

‣ Actor-Critic
- Learned Value Function
- Learned Policy



Value Function    Policy

Value-Based    Actor Critic    Policy-Based

[Slides from Fragkiadaki, 10-703 CMU]

# Summary: Interactive Reasoning

## Instruction following



**Train**

Go to the short red torch
Go to the blue keycard
Go to the largest yellow object
Go to the green object

**Test**

Go to the tall green torch
Go to the red keycard
Go to the smallest blue object

Go to the green torch

## Embodied learning



Q: What color is the car?

A: Orange!

## Reward shaping



*"Jump over the skull while going to the left"*

## Domain knowledge



*The natural resources available where a population settles affects its ability to produce food and goods. Build your city on a plains or grassland square with a river running through it if possible.*

Figure 1: An excerpt from the user manual of the game Civilization II.

[Luketina et al., A Survey of Reinforcement Learning Informed by Natural Language. IJCAI 2019]

# Interactive Reasoning Challenges

## Learning from open-ended manuals



```
                    A L I E N
                 20th Century Fox
               Games of the Century
          (picture of the ALIEN movie poster)
         "In space no one can hear you scream"
                 Game Instructions
                  Fox Video Games


                    A L I E N

TO SET UP: Set up your video computer system and left joystick controller as
instructed in your manufacturer owner's manual.  Move the Color/B-W lever to
the correct setting.  Turn the power OFF and insert the Alien game cartridge.


(Screen shot of the ALIEN maze setup: Alien, Alien Egg, Human, Pulsar and
Play Level-demo mode only)


TO BEGIN: Turn the power ON.  Use the Game Select lever and Difficulty
Switches to choose a play level.  Press the Game Reset lever and get ready
to run for your life.

THE OBJECTIVE: Your job is to run through the hallways of your space ship
and crush all the Alien Eggs which have been placed there.  You must also
avoid or destroy the adult Aliens and snatch up as many prizes as possible.

THE CONTROLS: Tilt the joystick forward, backward, left and right to
maneuver through the hallways.  To smash Eggs, simply run over them.  You
may travel off one side of the maze and back into the other using the
"Hyperwarp Passage."  Each Human is equipped with a Flame Thrower that is
activated by the joystick button (see below).

SCREEN DISPLAY: The Play Level and Humans allowed per Play Level are
displayed in the bottom left corner of the screen when Alien is not in play.
During the game, the current score and Humans remaining are shown there.

LEVELS OF PLAY/DIFFICULTY SWITCHES/BONUS ROUNDS: Each game of Alien lasts
until you run out of Humans.  If you can clear all of the Eggs out of a
playing screen, you get the chance to earn extra points in a "Bonus Round"
and then are returned to a new and more difficult playing screen.  All
points and Humans remaining are carried over to the new screens.

Bonus Rounds: The object of the Bonus Round is to travel STRAIGHT UP to the
top of the screen and grab the prize shown there.  You have only eight
seconds to do so.  You do not lose a human if you fail, but you earn the
point value of the prize if you succeed.

Left Difficulty Switch A: Aliens travel in random order about the screen.

Left Difficulty Switch B: Aliend travel in fixed patterns about the screen.

Right Difficult Switch B: Capturing a Pulsar has standard effect on the Aliens.

Right Difficulty Switch A: Capturing a Pulsar has no effect on the Aliens.


(Screen shot of ALIEN maze: Flame Thrower, Prize, Hyperwarp Passages, Humans
Remaining and Current Score)


LEVEL 1 - NORMAL GAME PLAY: You begin with three Humans and receive a bonus
Human after successfully clearing the second screen.  Prizes appear in chart
order.
```

```
LEVEL 2 - ADVANCED GAME PLAY: You begin with two Humans and receive no bonus
Humans.  Prizes appear in chart order.

LEVEL 3 - FOR EXPERTS ONLY: You begin with three Humans and receive no bonus
Human after clearing the first screen.  All Prizes in Level 3 are Saturns.

LEVEL 4 - EASY PRACTICE GAME: You begin with six Humans and receive 1 bonus
Human after clearing the first sceen.  All Prizes in Level 4 are also Saturns.

OBJECTS/SCORING: Each time an Alien catches you, one Human is lost.  You
score points for smashing Eggs and frying Aliens with the aid of your Flame
Thrower or Pulsar.  In addition, you can gain points for picking up Prizes.
Be sure to record your high scores on the back of this booklet!

(Screen shot of the bonus round with the human at the bottom of the screen,
the prize at the top of the screen and the horizontal moving Aliens in the
centre portion -- similar to the road portion of Frogger.)

FLAME THROWER - 1 PER HUMAN: A spurt of flam from this contraption cause
Aliens to turn away from you or become immobilized for a short period of
time.  Use the Throwers carefully.  Each has only four secons of flame and
the Thrower will not operate in the extreme left or right areas of the
screen.  You can also use the Flame Thrower to run over a Pulsar without
picking it up, allowing you to save the Pulsar to use at a later time.

PULSARS - 3 PER MAZE: Capturing a Pulsar causes the Aliens to weaken and
turn blue.  Then, for a short period of time, you can destroy them by
running over and touching them.  The instant the Aliens return to their
original colr, however, they once again become deadly.

PRIZES - 2 PER MAZE: Prizes appear in all levels of play and in the Bonus
Rounds.

POINT CHART:

OBJECT                  POINTS  PRIZES          POINTS
Eggs                            10      Rocket          500
Pulsar                  100             Saturn          1,000
1st Alien                       500             Star Ship       2,000
2nd Alien                       1,000   1st Surprise    2,000-3,000
3rd Alien                       2,000   2nd Surprise    3,000
Completed Screen                1               3rd Surprise    5,000


HINTS FROM DALLAS NORTH...
A good playing strategy is to crush all of the Eggs in one area at a time,
keeping within easy readh of a Pulsar.  The best way to destroy Aliens is to
sit near a Pulsar until the Aliens are almost upon you.  Then grab that
Pulsar and go get 'em !

Use the Hyperwarp Passage to ditch Aliens.  Many times they won't follow you in.

If you're having trouble with the Bonus Rounds, try going between the Alien
pairs rather than around them.

SUPER SMASHERS (a place to enter your high scores)
Name                            Level           Score
```
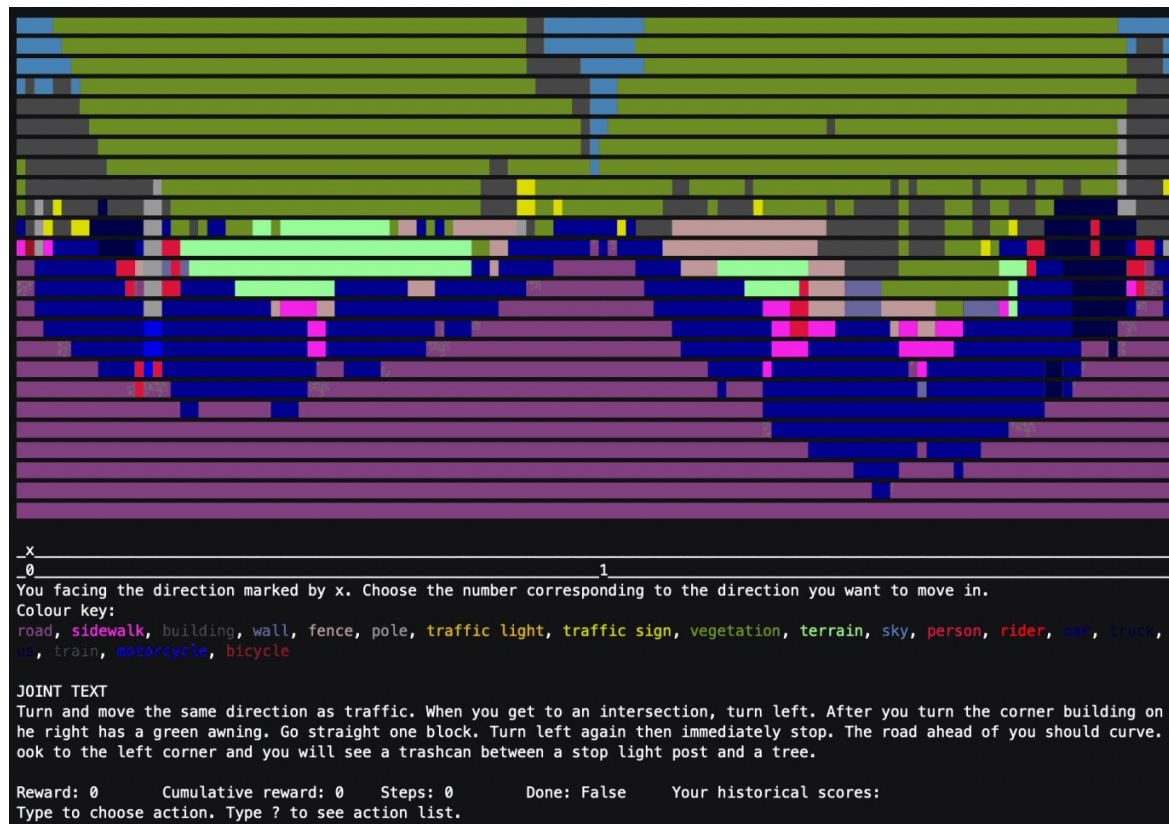
[Atari Learning Environment]

**Open challenges**

## Learning from text-based games



[Zhong et al., SILG: The Multi-environment Symbolic Interactive Language Grounding Benchmark. NeurIPS 2021]

# Interactive Reasoning Challenges

**Learning from lots of offline data**



[Fan et al., MineDojo: Building Open-Ended Embodied Agents with Internet-Scale Knowledge. arXiv 2022]
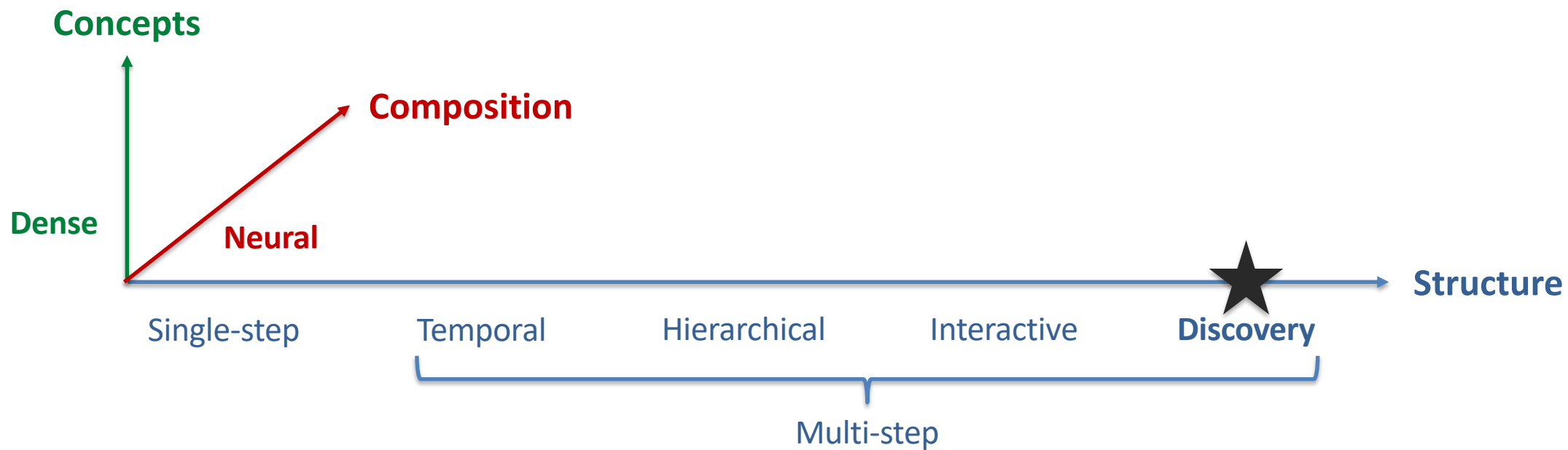
# Interactive Reasoning Challenges

**Hard to specify reward, but only final goal**
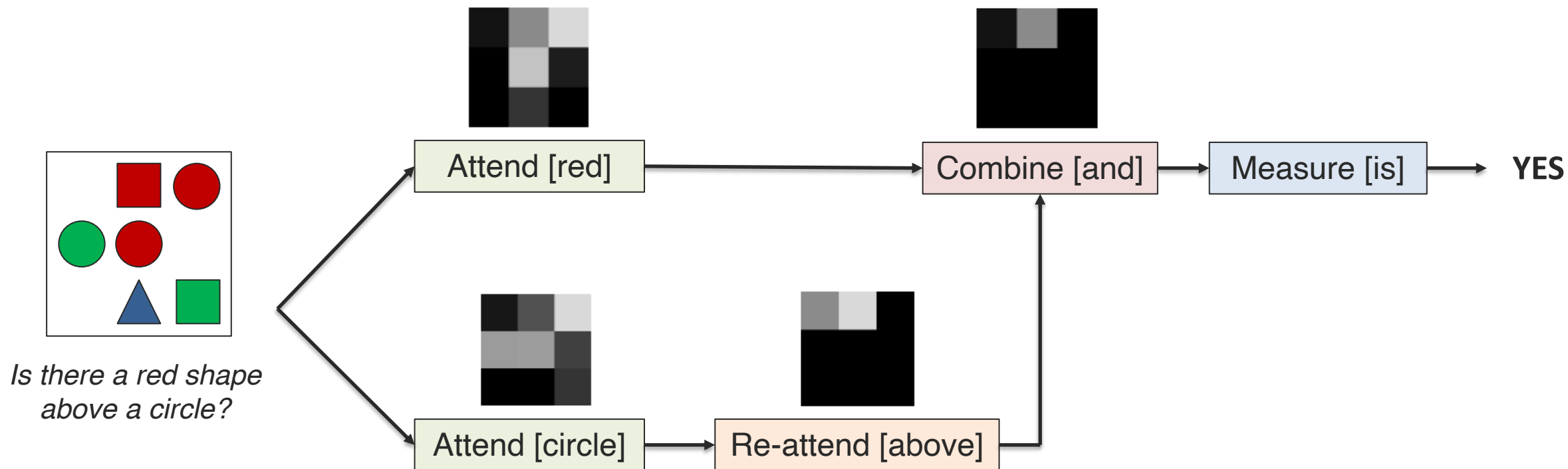


[Habitat Rearrangement Challenge 2022]

# Sub-Challenge 3a: Structure Modeling

# Structure Discovery

**End-to-end neural module networks**

Recall structure - leverage syntactic structure of language based on parsing



*Is there a red shape above a circle?*

[Andreas et al., Neural Module Networks. CVPR 2016]

# Structure Discovery

**End-to-end neural module networks**

Can we learn the structure end-to-end?

$$\tau \sim \pi_\theta(a|s)$$



$$\pi_\theta(a|s)$$

Attend [red]

Attend [circle]

Re-attend [above]

Combine [and]

Measure [is]

**?**

NMN

YES

$$r(\tau)$$

*Is there a red shape above a circle?*

[Hu et al., Learning to Reason: End-to-End Module Networks for Visual Question Answering. ICCV 2017]

# Stochastic Optimization

$$\max_{\theta} \mathbb{E}_{q_\theta(\mathbf{z})}[f(\mathbf{z})]$$

**RL**

$$\max_{\theta} J(\theta) \quad \text{Reward}$$

$$\max_{\theta} \mathbb{E}_{\tau \sim p(\tau;\theta)}[r(\tau)]$$

In RL (at least for discrete actions):     ???
- T is a sequence of discrete actions
- p(T; $\theta$ ) is not reparameterizable
- r(T) is a black box function
i.e. the environment

$$\pi_\theta(a|s)$$

r

a

s

**REINFORCE is a general-purpose solution!**

# Revisiting REINFORCE

$$\max_{\theta} \mathbb{E}_{q_{\theta}(\mathbf{z})}[f(\mathbf{z})]$$ <span style="color:red">(we will revisit this equation for generative models)</span>

We want to take gradients wrt $\theta$ of the term:

$$\nabla_{\theta} \mathbb{E}_{q_{\theta}(\mathbf{z})}[f(\mathbf{z})] = \mathbb{E}_{q_{\theta}(\mathbf{z})}[f(\mathbf{z}) \nabla_{\theta} \log q_{\theta}(\mathbf{z})]$$

We can now compute a Monte Carlo estimate:

Sample $\mathbf{z}^1, \mathbf{z}^2, ..., \mathbf{z}^K$ from $q_{\theta}(\mathbf{z})$ and estimate

$$\nabla_{\theta} \mathbb{E}_{q_{\theta}(\mathbf{z})}[f(\mathbf{z})] \approx \frac{1}{K} \sum_{k} [f(\mathbf{z}^k) \nabla_{\theta} \log q_{\theta}(\mathbf{z}^k)]$$

What we derived: sample trajectories and compute: $\boxed{\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)}$

<span style="color:green">- z can be discrete or continuous!</span>
<span style="color:green">- q(z) can be a discrete and continuous distribution!</span>
<span style="color:green">- q(z) must allow for easy sampling and be differentiable wrt $\theta$</span>
<span style="color:green">- f(z) can be a black box!</span>

# Structure Discovery

**End-to-end neural module networks**

Can we learn the structure end-to-end?

$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$$



$$\tau \sim \pi_\theta(a | s)$$

$$\pi_\theta(a_t | s_t)$$

RNN

Attend [red]
Attend [circle]
Re-attend [above]
Combine [and]
Measure [is]

NMN → YES

$r(\tau)$

*Is there a red shape above a circle?*

[Hu et al., Learning to Reason: End-to-End Module Networks for Visual Question Answering. ICCV 2017]

# Structure Discovery

**Structure fully learned from optimization and data**

$\tau \sim \pi_\theta(a|s)$

1. Define basic representation building blocks

| ReLU | Layer norm | Conv | Self-attention |

2. Define basic fusion building blocks

| Concat fuse | Attention fuse | Add fuse |

3. Automatically search for composition using neural architecture search

$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_\theta \log \pi_\theta(a_t|s_t)$$

**Nice, but slow!**



$\pi_\theta(a_t|s_t)$

[Xu et al., MUFASA: Multimodal Fusion Architecture Search for Electronic Health Records. AAAI 2021]

# Continuous Structure Discovery

Biggest problem: discrete optimization is slow.
Differentiable optimization for structure learning:

1. Approximate selection with softmax:

$$o'(x) = \sum_i \frac{\exp(\alpha_i)}{\sum_i \exp(\alpha_i)} o_i(x)$$

2. Solve bi-level optimization problem

$$\min_\alpha \quad \mathcal{L}_{val}(w^*(\alpha), \alpha)$$
$$\text{s.t.} \quad w^*(\alpha) = \text{argmin}_w \ \mathcal{L}_{train}(w, \alpha)$$

$$\mathcal{L}_{val}(w^*(\alpha), \alpha)$$

(valid data)

$y$

| Concat fuse | Add fuse | Attention fuse |

$\alpha_1$ $\alpha_2$ $\alpha_3$

| Conv | Layer norm |

[Liu et al., DARTS: Differentiable Architecture Search. ICLR 2019]

# Continuous Structure Discovery

Biggest problem: discrete optimization is slow.
Differentiable optimization for structure learning:

1. Approximate selection with softmax:

$$o'(x) = \sum_i \frac{\exp(\alpha_i)}{\sum_i \exp(\alpha_i)} o_i(x)$$

2. Solve bi-level optimization problem

$$\min_\alpha \quad \mathcal{L}_{val}(w^*(\alpha), \alpha)$$

$$\text{s.t.} \quad w^*(\alpha) = \text{argmin}_w \ \mathcal{L}_{train}(w, \alpha)$$

3. Convert softmax to argmax

**Faster but still non-trivial**

[Liu et al., DARTS: Differentiable Architecture Search. ICLR 2019]

# Continuous Structure Discovery

In general, optimization over directed acyclic graphs (DAGs):

Graph **G**, Data **X**, Adjacency matrix **W:**

$$\min_{W} \ell(W; X) \qquad \overset{?}{\Longleftrightarrow} \qquad \min_{W} \ell(W; X)$$

$$s.t. \ G(W) \in DAG \qquad\qquad\qquad s.t. \ h(W) = 0$$

(combinatorial 😱)         (smooth 😎)

(d)

[Zheng et al., DAGs with NO TEARS: Continuous Optimization for Structure Learning. NeurIPS 2018]

# Continuous Structure Discovery

$$\min_{W} \ell(W; X)$$
$$\text{s.t.} \quad G(W) \in DAG$$

$$\overset{?}{\Longleftrightarrow}$$

$$\min_{W} \ell(W; X)$$
$$\text{s.t.} \quad h(W) = 0$$

In <u>our paper</u>, we showed that such a function $h$ exists,

$$h(W) = \text{tr}(e^{W \circ W}) - d,$$



(d)

[Zheng et al., DAGs with NO TEARS: Continuous Optimization for Structure Learning. NeurIPS 2018]

# Continuous Structure Discovery

$$\min_{W} \quad \ell(W; X)$$

$$s.t. \quad G(W) \in DAG$$

$$\overset{?}{\iff}$$

$$\min_{W} \quad \ell(W; X)$$

$$s.t. \quad h(W) = 0$$

In our paper, we showed that such a function $h$ exists,

$$h(W) = \operatorname{tr}(e^{W \circ W}) - d,$$

and that it has a simple gradient:

$$\nabla h(W) = (e^{W \circ W})^T \circ 2W.$$



(d)

[Zheng et al., DAGs with NO TEARS: Continuous Optimization for Structure Learning. NeurIPS 2018]

# Continuous Structure Discovery

$$\min_{W} \ell(W; X)$$
$$\text{s.t. } G(W) \in DAG$$

$$\overset{?}{\Longleftrightarrow}$$

$$\min_{W} \ell(W; X)$$
$$\text{s.t. } h(W) = 0$$

In our paper, we showed that such a function $h$ exists,

$$h(W) = \operatorname{tr}(e^{W \circ W}) - d,$$

and that it has a simple gradient:

$$\nabla h(W) = (e^{W \circ W})^T \circ 2W.$$

Here the $\circ$ is the element-wise product, $d$ is the size of the graph, $\operatorname{tr}$ is the trace of a matrix, and the matrix exponential is defined as the infinite power series

$$e^A = I + A + \frac{1}{2!}A^2 + \frac{1}{3!}A^3 + \cdots$$

(d)

[Zheng et al., DAGs with NO TEARS: Continuous Optimization for Structure Learning. NeurIPS 2018]

# Continuous Structure Discovery

$$\min_W \ell(W; X)$$

s.t. $G(W) \in DAG$

$\overset{?}{\Longleftrightarrow}$

$$\min_W \ell(W; X)$$

s.t. $h(W) = 0$

$$h(W) = \text{tr}(e^{W \circ W}) - d,$$

$$e^A = I + A + \frac{1}{2!}A^2 + \frac{1}{3!}A^3 + \cdots$$

- *K*-th power of adjacency matrix **W** counts the number of *k*-step paths from one node to another.
- If the diagonal of the matrix power is all zeros, there are no *k*-step cycles.
- Acyclic = check all *k = 1,2, …, size of graph.*

Can now do continuous optimization to solve for W, but **nonconvex**

(d)

[Zheng et al., DAGs with NO TEARS: Continuous Optimization for Structure Learning. NeurIPS 2018]

# Sub-Challenge 3b: Intermediate Concepts

**Definition:** The parameterization of individual multimodal concepts in the reasoning process.

# Discrete Concepts via Hard Attention

Hard attention 'gates' (0/1) rather than soft attention (softmax between 0-1)
- Can be seen as discrete layers in between differentiable neural net layers



$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$$

[Xu et al., Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. ICML 2015]
[Chen et al., Multimodal Sentiment Analysis with Word-level Fusion and Reinforcement Learning. ICMI 2017]

# Discrete Concepts via Hard Attention

Hard attention 'gates' (0/1) rather than soft attention (softmax between 0-1)
- Can be seen as discrete layers in between differentiable neural net layers

**Sentiment analysis, emotion recognition**



Reject          Pass          Reject

*Figure 3.* Visualization of the attention for each generated word. The rough visualizations obtained by upsampling the attention weights and smoothing. (top)"soft" and (bottom) "hard" attention (note that both models generated the same captions in this example).

**Image captioning**



A    bird    flying    over    a    body    of    water    .

[Xu et al., Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. ICML 2015]
[Chen et al., Multimodal Sentiment Analysis with Word-level Fusion and Reinforcement Learning. ICMI 2017]

# Discrete Concepts via Language

- Large language/video/audio models interacting with each other
- Each language model has its own distinct *domain knowledge*
- Interaction is scripted and zero-shot



Guided multimodal discussion

Combining domain knowledge

[Zeng et al., Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language. arXiv 2022]

# Discrete Concepts via Language

Image captioning



[Zeng et al., Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language. arXiv 2022]

# Discrete Concepts via Language

Robot perception and planning



[Zeng et al., Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language. arXiv 2022]

# Discrete Concepts via Language

Video reasoning



[Zeng et al., Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language. arXiv 2022]

**Open challenges**

**Many open directions**

Prompt engineering – what is going on???

## Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing

| Pengfei Liu | Weizhe Yuan | Jinlan Fu |
|---|---|---|
| Carnegie Mellon University | Carnegie Mellon University | National University of Singapore |
| pliu3@cs.cmu.edu | weizhey@cs.cmu.edu | jinlanjonna@gmail.com |
| **Zhengbao Jiang** | **Hiroaki Hayashi** | **Graham Neubig** |
| Carnegie Mellon University | Carnegie Mellon University | Carnegie Mellon University |
| zhengbaj@cs.cmu.edu | hiroakih@cs.cmu.edu | gneubig@cs.cmu.edu |

We'll see more of this in transference

[Liu et al., Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. arXiv 2021]

# Sub-Challenge 3c: Inference Paradigm

**Definition:** How increasingly abstract concepts are inferred from individual multimodal evidences.

Recall representation fusion:

Modality A

Modality B

Fusion

**Concepts**

**Inference**

Representation

**Structure**

Single-step

Multi-step

**Potential issues:**
- Models may capture spurious correlations
- Not robust to targeted manipulations
- Lack of interpretability/control

# Sub-Challenge 3c: Inference Paradigm

**Definition:** How increasingly abstract concepts are inferred from individual multimodal evidences.

**Towards explicit inference paradigms:**
1. Logical inference: given premises inferred from multimodal evidence, how can one derive **logical** conclusions?

# Logical Inference

**Inference through logical operators in question**



*Is there beer AND is there a WINE GLASS?*

*Is the man NOT wearing shoes AND is there beer?*

*Is there beer?*     *Is the man wearing shoes?*

Adversarial antonyms ✗

Logical connectives ✗

Basic premises ✓

Existing models struggle to capture logical connectives.
How can we make them more logical?

[Gokhale et al., VQA-LOL: Visual Question Answering Under the Lens of Logic. ECCV 2020]

# Logical Inference

**Inference through logical operators in question**

⬛ *AND* 🔺

*Differentiable **AND** composition operator!*



*Are they in a
restaurant **AND**
are they all boys?*

**Also applies to other logic connectives:
AND, OR, NOT**

*Are they in a
restaurant?*

*Are they all
boys?*

[Gokhale et al., VQA-LOL: Visual Question Answering Under the Lens of Logic. ECCV 2020]

# Soft Logical Operators

**Inference through logical operators in question**

■ *AND* ▲

*Differentiable **AND** composition operator!*

Fréchet inequalities:

$$max(0, p(A_1) + p(A_2) - 1) \leq p(A_1 \wedge A_2) \leq min(p(A_1), p(A_2)).$$

$$p(A_1 \wedge A_2) = p(A_1) + p(A_2) - p(A_1 \vee A_2) \qquad p(A_1 \vee A_2) \leq 1$$

$$p(A_1 \wedge A_2) \geq p(A_1) + p(A_2) - 1$$

[Gokhale et al., VQA-LOL: Visual Question Answering Under the Lens of Logic. ECCV 2020]

# Soft Logical Operators

**Inference through logical operators in question**

◼ *OR* ▲

*Differentiable **OR** composition operator!*

Fréchet inequalities:

$$max(0, p(A_1) + p(A_2) - 1) \leq p(A_1 \wedge A_2) \leq min(p(A_1), p(A_2)).$$

Differentiable, so you can now optimize wrt $p(A_1 \vee A_2)$ and $p(A_1 \vee A_2)$

[Gokhale et al., VQA-LOL: Visual Question Answering Under the Lens of Logic. ECCV 2020]

**Many open directions**



HasOfficeInCity(New York, Uber)
CityInCountry(USA, New York)

X = Uber

Y = USA

In which country Y does X have office?

HasOfficeInCountry(Y, X) ?

HasOfficeInCountry(Y, X) ← HasOfficeInCity(Z, X), CityInCountry(Y, Z)

X = Lyft

Y = France

HasOfficeInCity(Paris, Lyft)
CityInCountry(France, Paris)

Differentiable knowledge base reasoning

[Yang et al., Differentiable Learning of Logical Rules for Knowledge Base Reasoning. NeurIPS 2017]

# Sub-Challenge 3c: Inference Paradigm

**Definition:** How increasingly abstract concepts are inferred from individual multimodal evidences.

**Towards explicit inference paradigms:**
1. Logical inference
2. Causal inference: how can one determine the actual **causal** effect of a variable in a larger system?

# Causal Inference

**Association vs causation**

**Example:** How does class size impact student outcomes?

Why can't we just compare student outcomes among different class sizes?
- Poorer districts may have larger class sizes.
- Students in poorer districts may have access to fewer resources, more difficult family circumstances, etc.
- All of these factors may impact student outcomes.

Association describes how things are. Causation describes how things would have been under different circumstances.

(side note: correlation is a specific type of linear association)

[Slides from Victoria Lin]

# Causal Inference

**Association vs causation**

> ## Simple linear regression
>
> Consider the simple linear regression model
>
> $$Y = \beta_0 + \beta_1 X_1 + \ldots + \beta_p X_p + \epsilon$$

[Slides from Victoria Lin]

# Causal Inference

**Association vs causation**

> ## Simple linear regression
> Consider the simple linear regression model
>
> $$Y = \beta_0 + \beta_1 X_1 + \ldots + \beta_p X_p + \epsilon$$

How do we interpret the coefficient $\beta_1$?

- Commonly: The expected change in outcome $Y$ if covariate $X_1$ were increased by one, holding all other covariates constant
- Correctly: The expected difference in outcome for two data points who *happen to have* the same covariate values for $(X_2, \ldots, X_p)$ and whose values for $X_1$ *happen to differ* by one

*The first interpretation is in fact causal and requires extra assumptions!*

[Slides from Victoria Lin]

# Causal Inference

**Intervention**

Causal inference is reliant on the idea of interventions —what outcome might have occurred if X happened (an intervention), possibly contrary to observed data.

# Causal Inference

**Intervention**

Causal inference is reliant on the idea of interventions —what outcome might have occurred if X happened (an intervention), possibly contrary to observed data.



```
x = randn()
y = x + 1 + sqrt(3)*randn()
```
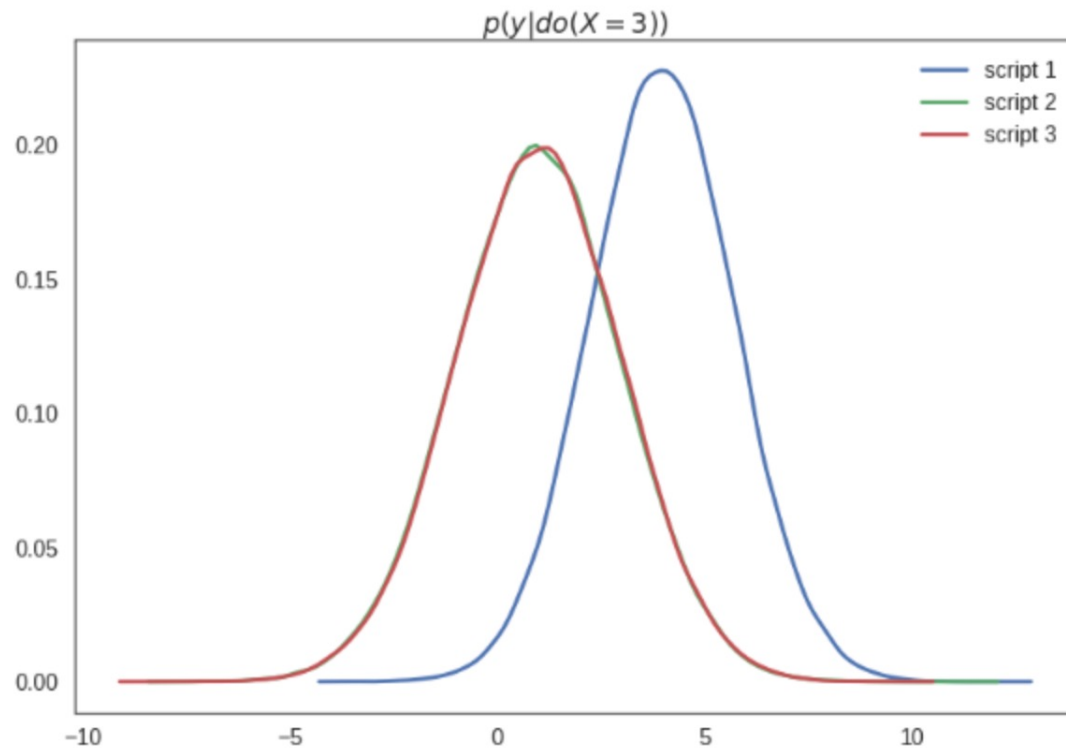
[Example from Ferenc Huszár: https://www.inference.vc/causal-inference-2-illustrating-interventions-in-a-toy-example/]

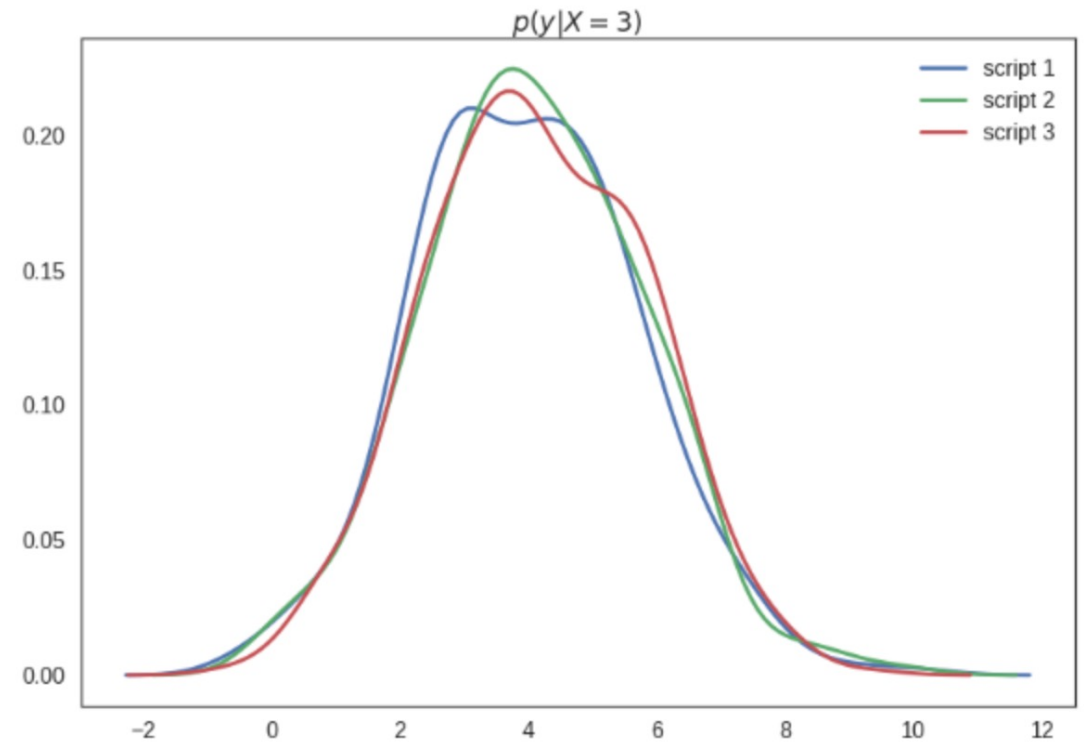# Causal Inference

**Intervention**

Causal inference is reliant on the idea of interventions —what outcome might have occurred if X happened (an intervention), possibly contrary to observed data.



[Example from Ferenc Huszár: https://www.inference.vc/causal-inference-2-illustrating-interventions-in-a-toy-example/]
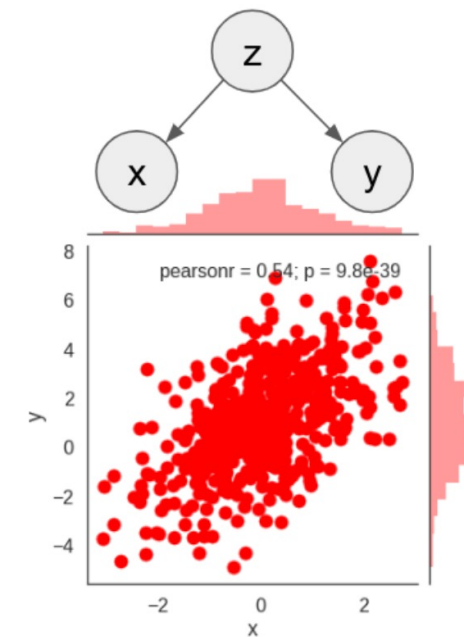
# Causal Inference

**Intervention**

Causal inference is reliant on the idea of interventions —what outcome might have occurred if X happened (an intervention), possibly contrary to observed data.



[Example from Ferenc Huszár: https://www.inference.vc/causal-inference-2-illustrating-interventions-in-a-toy-example/]
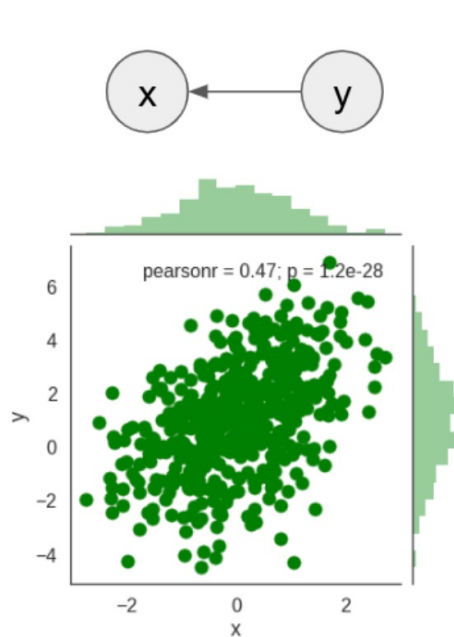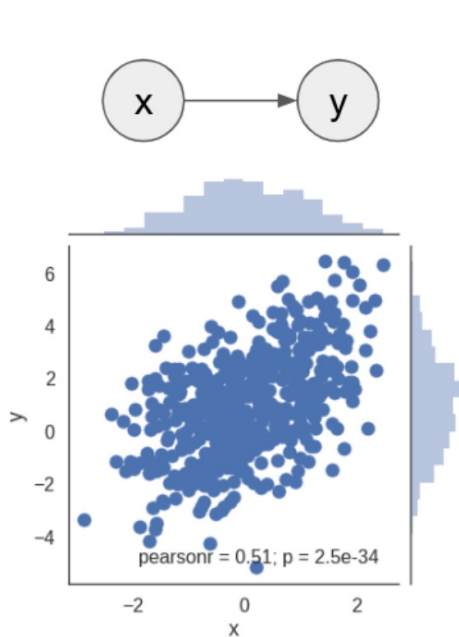
# Causal Inference

**Intervention**

Let's say I really want to set the value of *x* to 3.



```
x = randn()
x = 3
y = x + 1 + sqrt(3)*randn()
x = 3
```

```
y = 1 + 2*randn()
x = 3
x = (y-1)/4 + sqrt(3)*randn()/2
x = 3
```

```
z = randn()
x = 3
x = z
x = 3
y = z + 1 + sqrt(3)*randn()
x = 3
```

[Example from Ferenc Huszár: https://www.inference.vc/causal-inference-2-illustrating-interventions-in-a-toy-example/]

# Causal Inference

**Intervention**

Let's say I really want to set the value of *x* to 3. What happens to *y*?

# Causal Inference

**Intervention**

The marginal distribution of *y*: p(y | do(x=3)).          The marginal distribution of *y*: p(y | x=3).



The joint distribution of data alone is insufficient to predict behavior under interventions.

[Example from Ferenc Huszár: https://www.inference.vc/causal-inference-2-illustrating-interventions-in-a-toy-example/]

# Causal Inference

**Causal diagrams: arrow pointing from cause to effect.**

**Intervention** mutilates the graph by removing all edges that point into the variable on which intervention is applied (in this case *x*).



$$P(y|do(X)) = p(y|x) \qquad P(y|do(X)) = p(y) \qquad P(y|do(X)) = p(y)$$

[Example from Ferenc Huszár: https://www.inference.vc/causal-inference-2-illustrating-interventions-in-a-toy-example/]
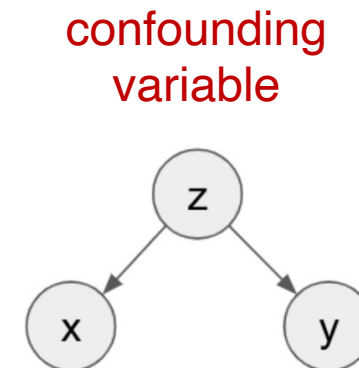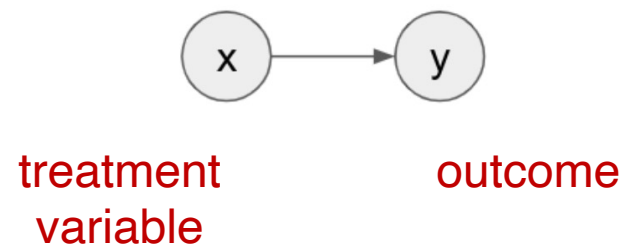
# Causal Inference

**Intervention in real-life is typically very hard!**

E.g., does treatment x treat disease y?

Can I estimate the intervention p(y|do(X=x))?
Requires answering: all else being equal, what would be the patient's outcome if they had not taken the treatment?



confounding
variable

treatment
variable

outcome

Lots of work, see Judea Pearl, The Book of Why
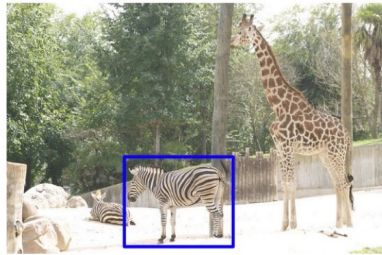
# Causal Inference

## Causal VQA: does my multimodal model capture causation or correlation?

**Covariant VQA**

Target object in question

Q: How many zebras are there in the picture?

A: 2



Baselines:                    **2**

i.e., treatment variable

zebras ⟶ prediction

BUT: correlation or causation?

[Agarwal et al., Towards Causal VQA: Revealing & Reducing Spurious Correlations by Invariant & Covariant Semantic Editing. CVPR 2020]

# Causal Inference

**Causal VQA: does my multimodal model capture causation or correlation?**

### Covariant VQA

Target object in question

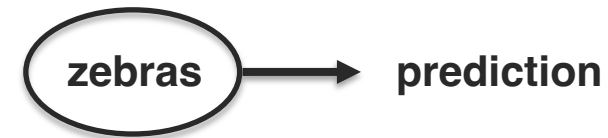Q: How many zebras are there in the picture?
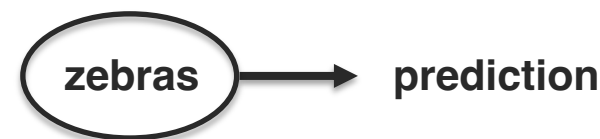A: 2          *zebra removed* A: 1



Baselines:          **2**                    **2**

*i.e., treatment variable*

zebras ⟶ **prediction**

**Interventional conditional:** $p(y|do(zebras = 1))$

Existing models struggle to adapt to targeted causal interventions.
How can we make them more robust to spurious correlations?

[Agarwal et al., Towards Causal VQA: Revealing & Reducing Spurious Correlations by Invariant & Covariant Semantic Editing. CVPR 2020]

# Causal Inference

**Causal VQA: does my multimodal model capture causation or correlation?**

**Invariant VQA**

Target irrelevant object

Q: What color is the balloon?

A: red



Baselines:            **pink**

umbrella

*i.e., confounding variable*

balloon ⟶ prediction

Is my model picking up irrelevant objects?

[Agarwal et al., Towards Causal VQA: Revealing & Reducing Spurious Correlations by Invariant & Covariant Semantic Editing. CVPR 2020]

# Causal Inference

**Causal VQA: does my multimodal model capture causation or correlation?**

**Invariant VQA**

Target irrelevant object
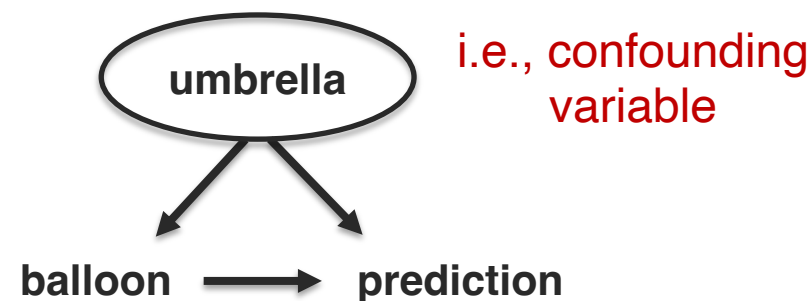
Q: What color is the balloon?

A: red          *umbrellas removed*; A: red



Baselines:          **pink**                    **red**



umbrella

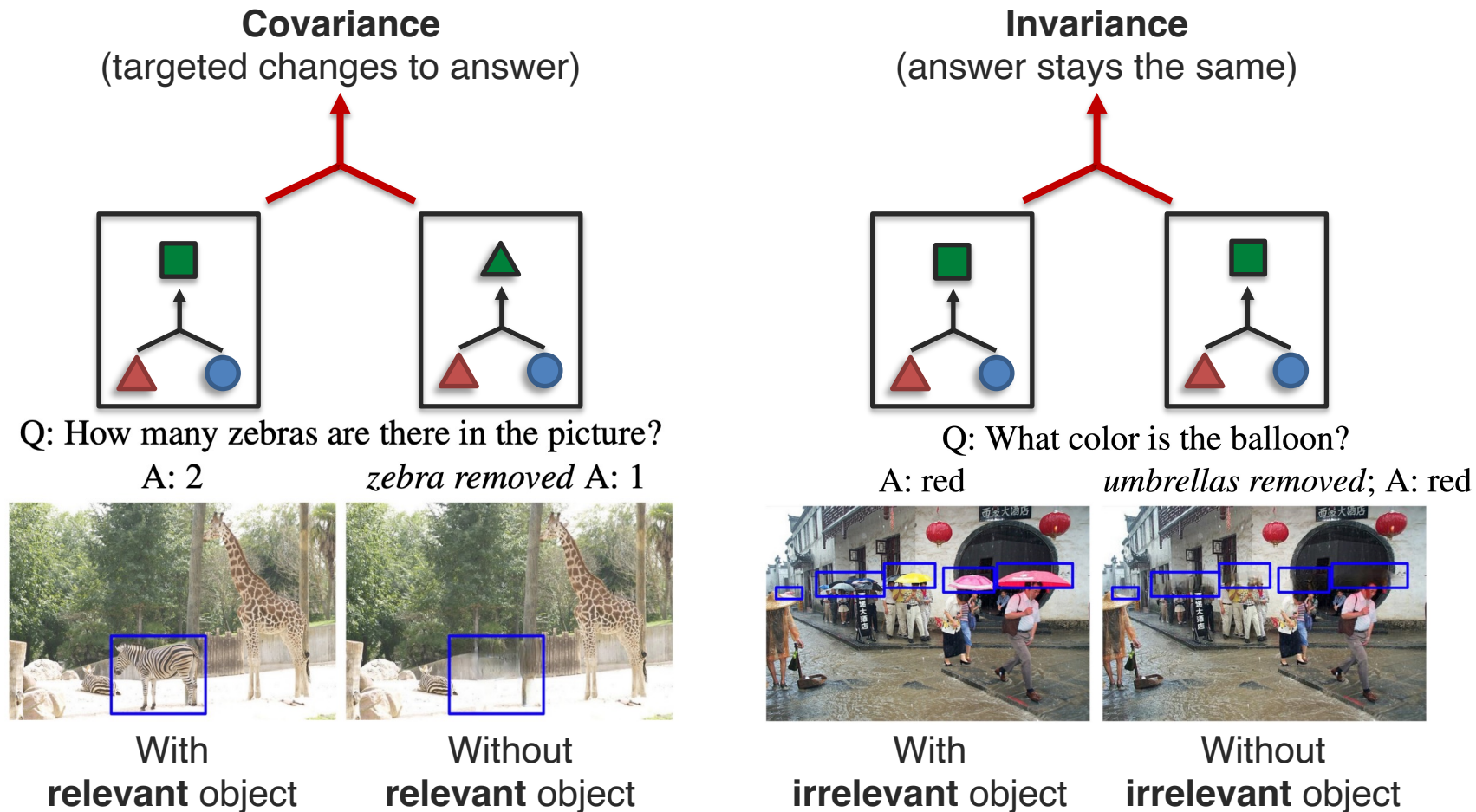i.e., confounding variable

balloon ⟶ prediction

**Interventional conditional:** $p(y|do(no\ umbrella))$

Existing models struggle to adapt to targeted causal interventions.
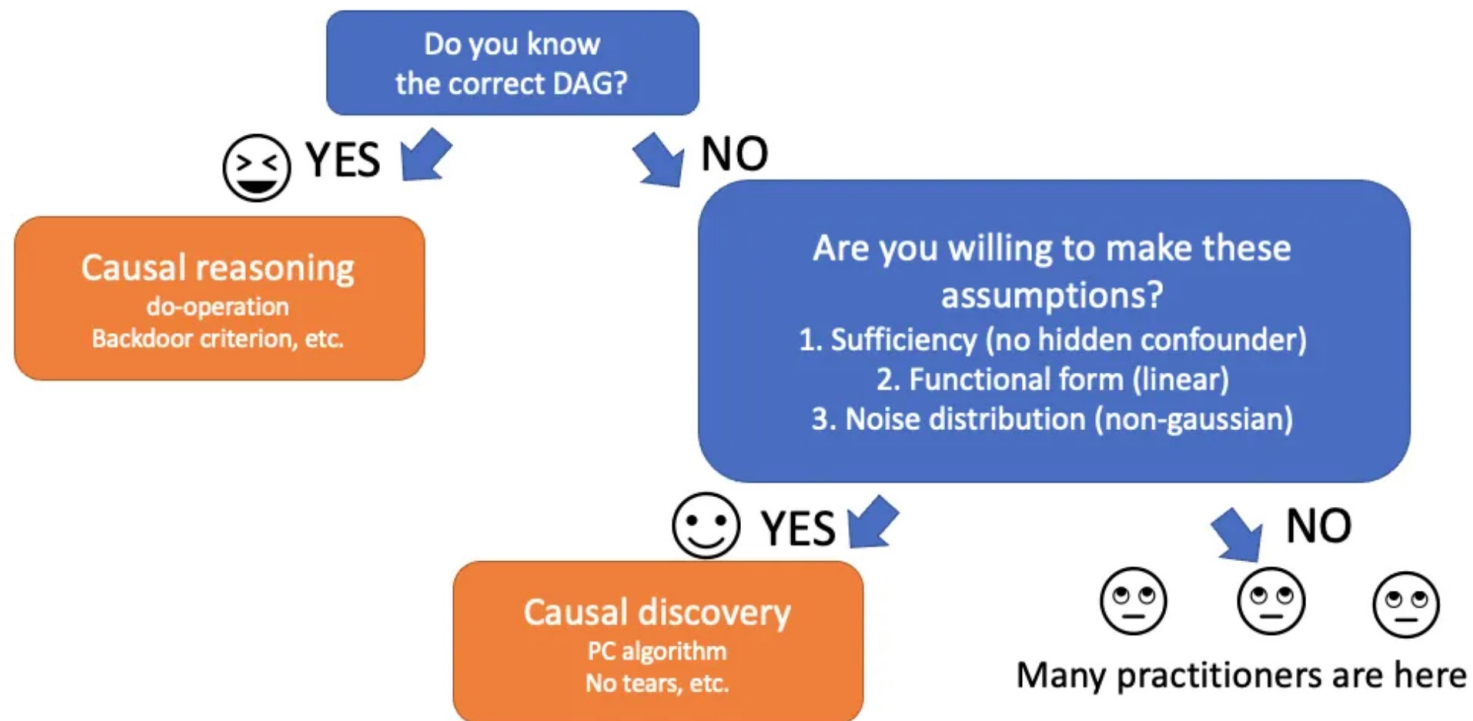How can we make them more robust to spurious correlations?

[Agarwal et al., Towards Causal VQA: Revealing & Reducing Spurious Correlations by Invariant & Covariant Semantic Editing. CVPR 2020]

# Causal Inference

**Causal inference via data augmentation**



**Covariance**
(targeted changes to answer)

Q: How many zebras are there in the picture?
A: 2          *zebra removed* A: 1

With
**relevant** object          Without
**relevant** object

**Invariance**
(answer stays the same)

Q: What color is the balloon?
A: red          *umbrellas removed*; A: red

With
**irrelevant** object          Without
**irrelevant** object

[Agarwal et al., Towards Causal VQA: Revealing & Reducing Spurious Correlations by Invariant & Covariant Semantic Editing. CVPR 2020]

**Many open directions**



Application of causality – current state

Causal deep learning, see https://www.vanderschaar-lab.com/causal-deep-learning/

**Many open directions**



**Ladder of causation**

- 3 – Counterfactual
- 2 – Intervention
- 1.5 – CDL
- 1 – Association

**The space between association and intervention**

Many interesting ML problems lie in Rung 1.5

- Robustness
  - Distribution shift
  - Adversarial attack
- Generalization
  - Domain adaptation
  - Transfer learning
  - Meta-learning
  - Few-shot learning
- Other potential areas
  - Fairness
  - Data augmentation
  - Etc.

1. Empirically verifiable
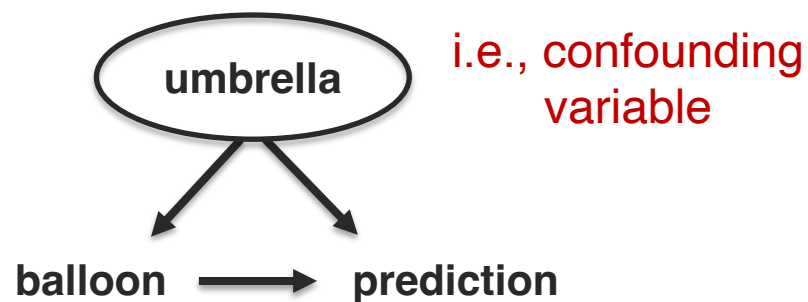2. "Good enough"

Causal deep learning, see https://www.vanderschaar-lab.com/causal-deep-learning/

# Sub-Challenge 3c: Inference Paradigm

**Definition:** How increasingly abstract concepts are inferred from individual multimodal evidences.
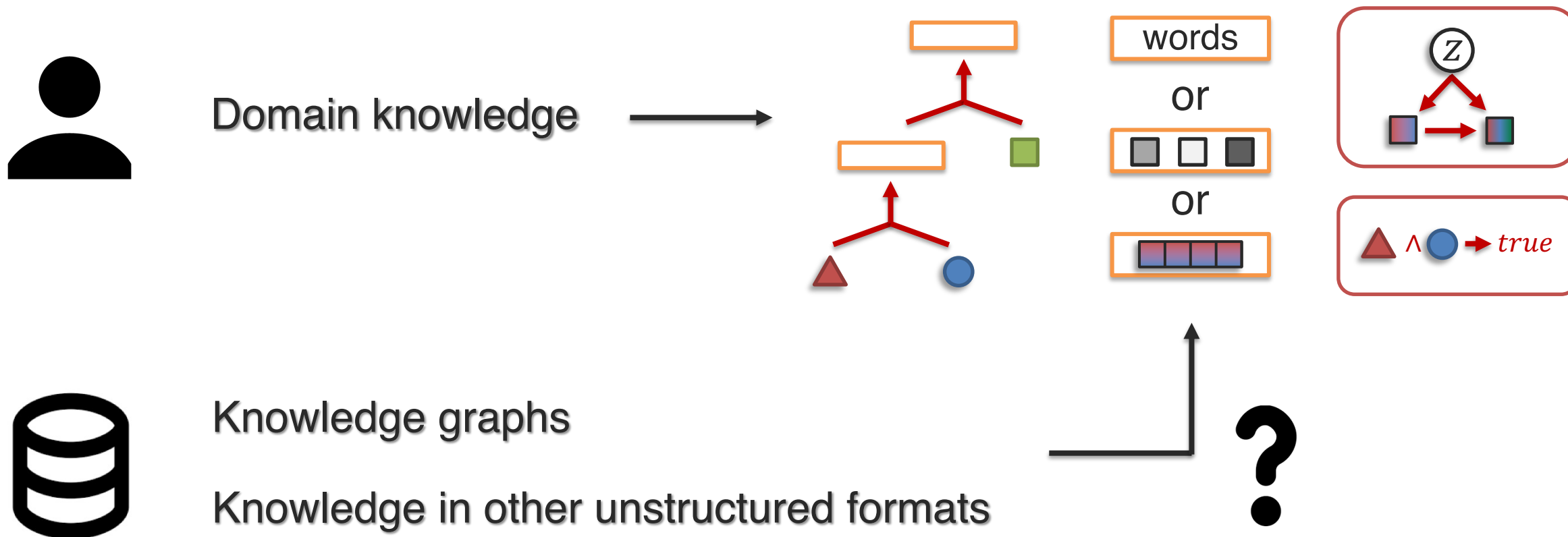
**Towards explicit inference paradigms:**
1. Logical inference
2. Causal inference

**Nice, but you don't get these for free!**

i.e., confounding variable

umbrella

balloon ⟶ prediction

Inference

Causal

*true*

Logical

Representation

# Sub-Challenge 3d: Knowledge

**Definition:** The derivation of knowledge in the study of inference, structure, and reasoning.



Domain knowledge

words

or

or

Knowledge graphs

Knowledge in other unstructured formats

# External Knowledge: Multimodal Knowledge Graphs

**Knowledge can also be gained from external sources**
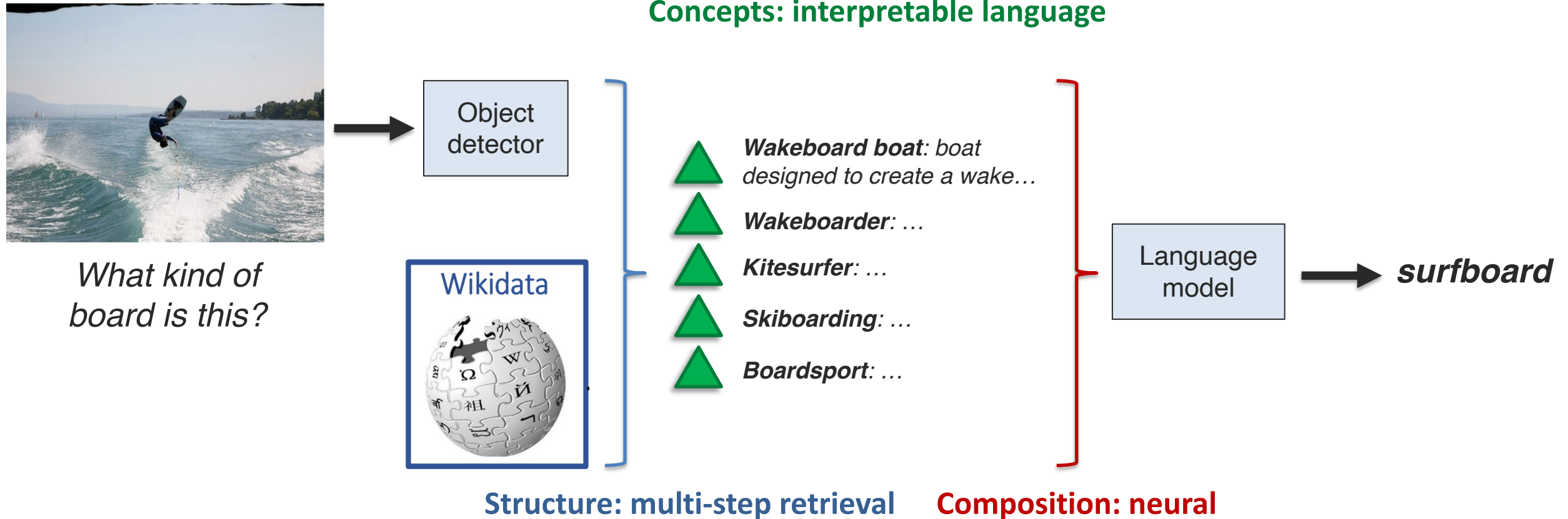


*What kind of board is this?*

*Requires knowledge of water sports, sports equipment, etc.*

Existing models struggle when external knowledge is needed. How can we leverage external knowledge?

[Marino et al., OK-VQA: A visual question answering benchmark requiring external knowledge. CVPR 2019]
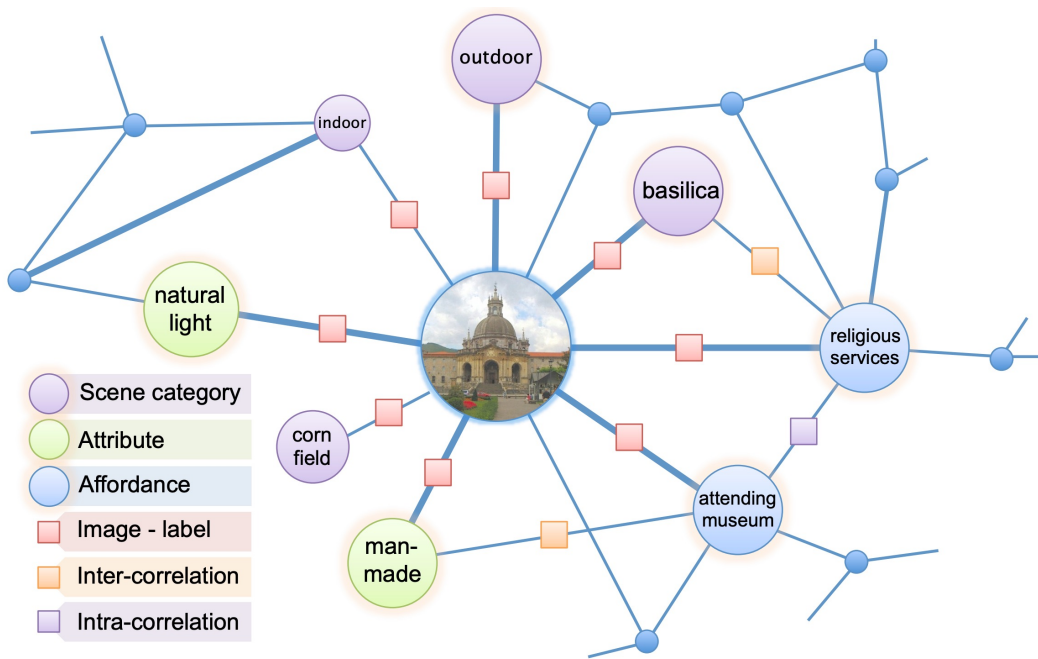
# External Knowledge: Multimodal Knowledge Graphs

**Knowledge can also be gained from external sources**



**Concepts: interpretable language**

**What kind of board is this?**

Object detector

Wikidata

*Wakeboard boat*: boat designed to create a wake…

*Wakeboarder*: …

*Kitesurfer*: …

*Skiboarding*: …

*Boardsport*: …

Language model → **surfboard**

**Structure: multi-step retrieval**     **Composition: neural**

[Gui et al., KAT: A Knowledge Augmented Transformer for Vision-and-Language. NAACL 2022]

# External Knowledge: Multimodal Knowledge Graphs

**Knowledge can also be gained from external sources**



**Concepts: interpretable**

**Structure: multi-step inference**

**Composition: graph-based**

**Class** — **auditorium**

**Affordances** — community and social work, taking class for personal interest, religious practices, waiting, attending the performing arts

**Attributes** — congregating, indoor lighting, spectating, enclosed area, glossy

[Zhu et al., Building a Large-scale Multimodal Knowledge Base System for Answering Visual Queries. arXiv 2015]

# External Knowledge Challenges

Atomic: If-then commonsense

[Sap et al., Atomic: An Atlas of Machine Commonsense for If-Then Reasoning. AAAI 2019]

Open challenges

Delphi: Moral commonsense

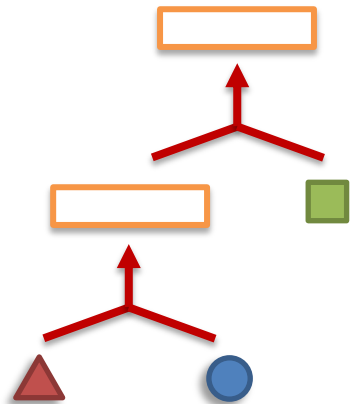Social Chemistry: Social commonsense

[Jiang et al., Can Machines Learn Morality? The Delphi Experiment. arXiv 2021]
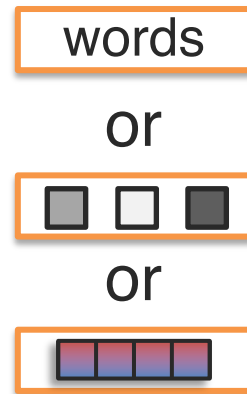[Forbes et al., Social Chemistry 101: Learning to Reason about Social and Moral Norms. EMNLP 2020]

# Summary: Reasoning

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.



(A) Structure modeling

(B) Intermediate concepts

words

or

or

(C) Inference paradigm

$z$

▲ ∧ ● → *true*

(D) External knowledge

# Summary: Reasoning

**Definition:** Combining knowledge, usually through multiple inferential steps, exploiting multimodal alignment and problem structure.
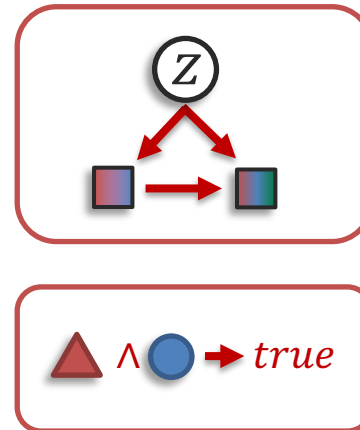
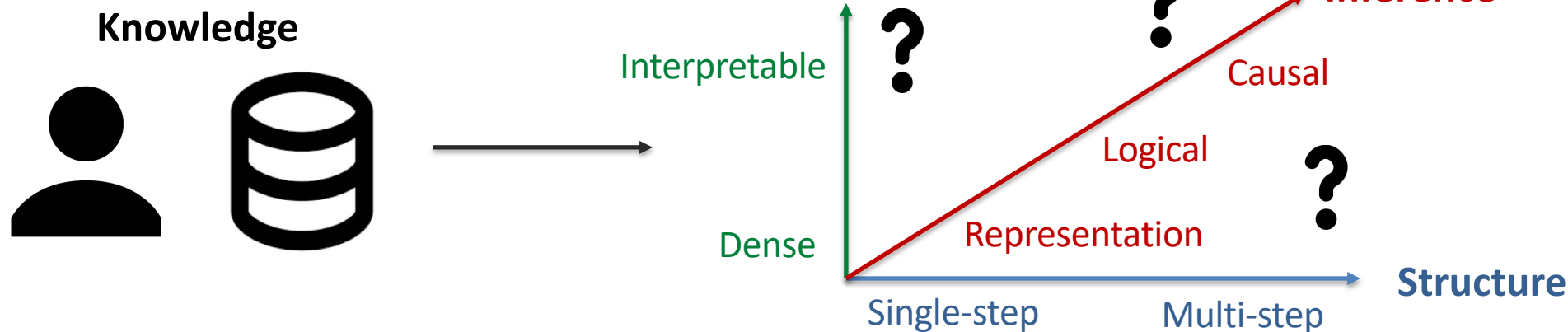| | (A) Structure modeling | (B) Intermediate concepts | (C) Inference paradigm | (D) External knowledge |
|---|---|---|---|---|
| **Last Thursday** | Temporal Hierarchical | Continuous | | |
| **Tuesday** | Interactive | | | |
| **Today** | Discovery | Discrete | Causal Logical | Knowledge Commonsense |

# More Reasoning

**Open challenges:**
- Structure: multi-step inference
- Concepts: interpretable + differentiable representations
- Composition: explicit, logical, causal…
- Knowledge: integrating explicit knowledge with pretrained models
- Probing pretraining models for reasoning capabilities