# Concept discovery from an image dataset
## Toward image representation with an emergent  language

Yuta Nakashima

Institute for Datability Science

Osaka University, Japan

June 28, 2023

# Question

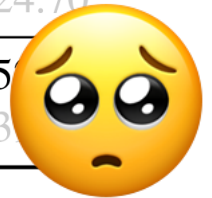🤔 Does a model really see an image/video?

# In visual question answering?
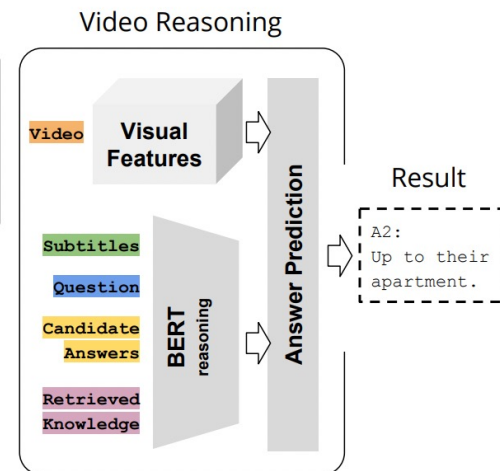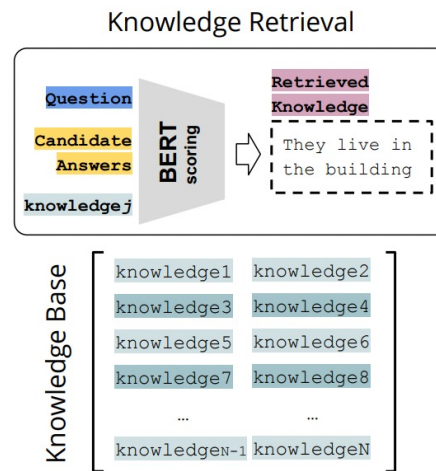


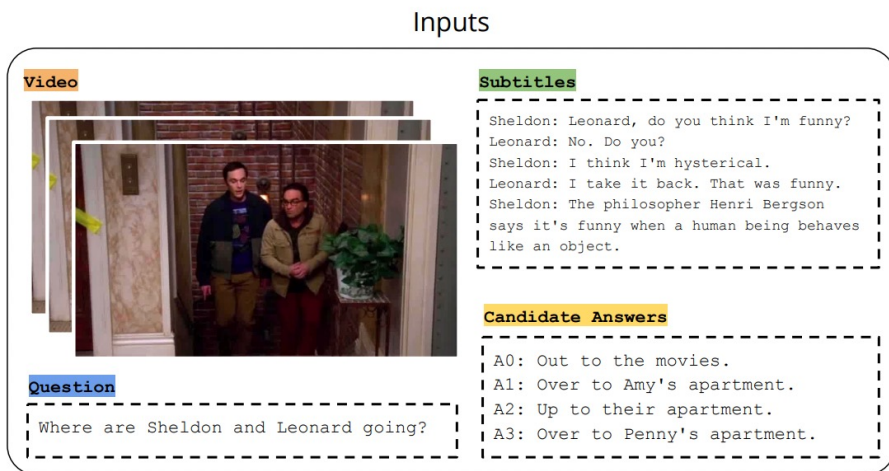What color are her eyes?
What is the mustache made of?

How many slices of pizza are there?
Is this a vegetarian pizza?

From: [Antol et al., "VQA: Visual question answering," ICCV 2015]

# In visual question answering?

|  | Model | Dataset | Overall | Yes/No | Number | Other |
|---|---|---|---|---|---|---|
| **Text-only** | per Q-type prior [5] | VQA v1 | 35.13 | 71.31 | 31.93 | 08.86 |
|  |  | VQA-CP v1 | 08.39 | 14.70 | 08.34 | 02.14 |
|  | d-LSTM Q [5] | VQA v1 | 48.23 | 79.05 | 33.70 | 28.81 |
|  |  | VQA-CP v1 | 20.16 | 35.72 | 11.07 | 08.34 |
| **Text and image** | d-LSTM Q + norm I [21] | VQA v1 | 54.40 | 79.82 | 33.87 | 40.54 |
|  |  | VQA-CP v1 | 23.51 | 34.53 | 11.40 | 17.42 |
|  | NMN [3] | VQA v1 | 54.83 | 80.39 | 33.45 | 41.07 |
|  |  | VQA-CP v1 | 29.64 | 38.85 | 11.23 | 27.88 |
|  | SAN [36] | VQA v1 | 55.86 | 78.54 | 33.46 | 44.51 |
|  |  | VQA-CP v1 | 26.88 | 35.34 | 11.34 | 24.70 |
|  | MCB [8] | VQA v1 | 60.97 | 81.62 | 34.56 | 5 |
|  |  | VQA-CP v1 | 34.39 | 37.96 | 11.80 | 3 |

# In video question answering?

## Inputs

**Video**

**Subtitles**

```
Sheldon: Leonard, do you think I'm funny?
Leonard: No. Do you?
Sheldon: I think I'm hysterical.
Leonard: I take it back. That was funny.
Sheldon: The philosopher Henri Bergson
says it's funny when a human being behaves
like an object.
```

**Question**

```
Where are Sheldon and Leonard going?
```

**Candidate Answers**

```
A0: Out to the movies.
A1: Over to Amy's apartment.
A2: Up to their apartment.
A3: Over to Penny's apartment.
```

## Knowledge Retrieval

**Question**

**Candidate Answers**

**knowledge$j$**

BERT scoring → **Retrieved Knowledge**

```
They live in
the building
```

**Knowledge Base**

| | |
|---|---|
| knowledge1 | knowledge2 |
| knowledge3 | knowledge4 |
| knowledge5 | knowledge6 |
| knowledge7 | knowledge8 |
| ... | ... |
| knowledge$N-1$ | knowledge$N$ |

## Video Reasoning

**Video** → **Visual Features**

**Subtitles**, **Question**, **Candidate Answers**, **Retrieved Knowledge** → BERT reasoning → Answer Prediction

**Result**

```
A2:
Up to their
apartment.
```

---

Penny: What are you doing at work these days?
Sheldon: Oh. I'm working on time-dependent backgrounds in string theory.
Specifically, quantum field theory in D-dimensional de Sitter space. in
D-dimensional de Sitter space. (...)

**What night is it?**

Wednesday
Monday
Friday
**Saturday** ✓

**R. Knowledge**

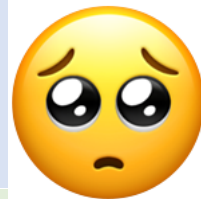Saturday is
Sheldon's
laundry night.

---

Leonard: You gotta admit, I'm delightful.
Penny: Why are you making this so difficult?
Leonard: It's not difficult for me, I'm having fun.
Penny: What do you want me to do? (...)

**Where is this taking place?**

A local Dance Center
**A bowling alley** ✓
Sheldon's bedroom
A blues dance club

**R. Knowledge**

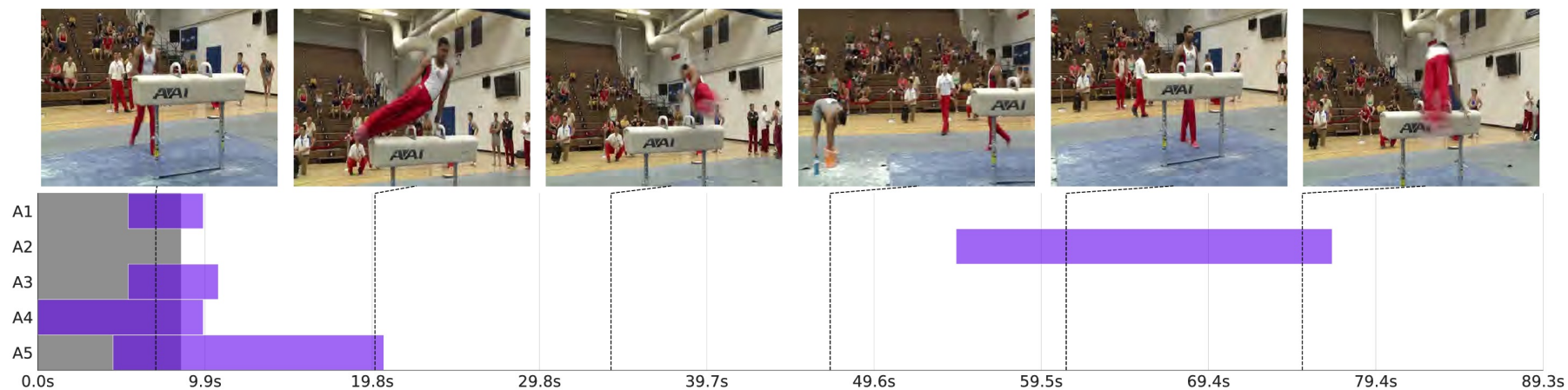It's the closest thing
to sports they like to
partake in.

---

[Garcial et al., "KnowIT VQA: Answering Knowledge-Based Questions about Videos," AAAI 2020]

# In video question answering?

| | Model | Vis. | Text. | Temp. | Know. | All |
|---|---|---|---|---|---|---|
| | **Random** | 0.250 | 0.250 | 0.250 | 0.250 | 0.250 |
| Text-only — QA | word2vec sim | 0.108 | 0.163 | 0.151 | 0.180 | 0.161 |
| | BERT sim | 0.174 | 0.264 | 0.209 | 0.190 | 0.196 |
| | TFIDF | 0.434 | 0.377 | 0.488 | 0.485 | 0.461 |
| | LSTM Emb. | 0.444 | 0.428 | 0.512 | 0.515 | 0.489 |
| | LSTM BERT | 0.446 | 0.464 | 0.500 | 0.532 | 0.504 |
| | $\diamond$ ROCK$_{QA}$ | 0.542 | 0.475 | 0.547 | 0.535 | 0.530 |
| | **Humans** (Rookies, Blind) | 0.406 | 0.407 | 0.418 | 0.461 | 0.440 |
| Subs, QA | LSTM Emb. | 0.432 | 0.362 | 0.512 | 0.496 | 0.467 |
| | LSTM BERT | 0.452 | 0.446 | 0.547 | 0.530 | 0.504 |
| | TVQA$_{SQA}$ | 0.602 | 0.551 | 0.512 | 0.468 | 0.509 |
| | $\diamond$ ROCK$_{SQA}$ | 0.651 | 0.754 | 0.593 | 0.534 | 0.587 |
| | **Humans** (Rookies, Subs) | 0.618 | 0.837 | 0.453 | 0.498 | 0.562 |
| Text and image — Vis, Subs, QA | TVQA | 0.612 | 0.645 | 0.547 | 0.466 | 0.522 |
| | $\diamond$ ROCK$_{VSQA}$ Image | 0.643 | 0.739 | 0.581 | 0.539 | 0.587 |
| | $\diamond$ ROCK$_{VSQA}$ Concepts | 0.647 | 0.743 | 0.581 | 0.538 | 0.587 |
| | $\diamond$ ROCK$_{VSQA}$ Facial | 0.649 | 0.743 | 0.581 | 0.537 | 0.587 |
| | $\diamond$ ROCK$_{VSQA}$ Caption | 0.666 | 0.772 | 0.581 | 0.514 | 0.580 |
| | **Humans** (Rookies, Video) | 0.936 | 0.932 | 0.624 | 0.655 | 0.748 |

🥺

[Garcial et al., "KnowIT VQA: Answering Knowledge-Based Questions about Videos," AAAI 2020]

# In video moment retrieval?



As the walk continues,the cat stops and begins staring at a parked car with large red flames painted on the side.



A male gymnast walks up to a beam.

[Otani et al., "Uncovering hidden challenges in query-based video moment retrieval," BMVC 2020]

# In video moment retrieval?

## ActivityNet Captions



Bar chart showing results on ActivityNet Captions:
- Prior-Only Blind: 10.8
- Action-Aware Blind: 23.1
- CTRL: 29.0
- ACRN: 31.7
- TripNet: 32.2
- QSPN: 33.3
- SCDM: 35.9
- ABLR: 36.8
- Blind-TAN: 41.3
- 2D-TAN: 44.0

[Otani et al., "Uncovering hidden challenges in query-based video moment retrieval," BMVC 2020]

# Spurious correlation matters

## Visual question answering



From: [Agrawal et al., "Don't Just Assume; Look and Answer: Overcoming Priors for Visual Question Answering," CVPR 2018]

## Video moment retrieval



[Otani et al., "Uncovering hidden challenges in query-based video moment retrieval," BMVC 2020]

# A similar happens in visual-input-only tasks



(A) **Cow: 0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98

(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97

(C) No Person: 0.97, **Mammal: 0.96**, Water: 0.94, Beach: 0.94, Two: 0.94

From: [Beery et al., "Recognition in Terra Incognita" ECCV 2018]

# Mitigating the spurious correlation problem

- Considering certain attributes; for example:
  - [Burns et al., "Women also snowboard: Overcoming bias in captioning models," ECCV 2018]
  - [Agarwal et al., "Does data repair lead to fair models? Curating contextually fair data to reduce model bias," WACV 2022]

- Background matters; for example:
  - [Sagawa et al., "Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization," ICLR 2020]
  - [Taghanaki et al., "Robust representation learning via perceptual similarity metrics," NeurIPS 2022]

- Identifying confounders; for example:
  - [Liu et al., "Show, deconfound and tell: Image captioning with causal inference," CVPR 2022]

# Yet another (perhaps crazy) idea?

- Representing an image/video with a set or a sequence of discrete labels, but not with a continuous vector

# Discovering concepts in an image dataset

- SCOUTER [Li et al., ICCV 2021]
  - Better explainability with a *SINGLE* distinctive visual feature per class

$$7 = \boxed{7}$$



- BotCL [Wang et al., CVPR 2023]
  - Discovering concepts that describe image for a given classification task

$$7 = \{\boxed{-}, \boxed{'}, \boxed{7}, \boxed{/}\}$$

# SCOUTER 🧐 [Li et al. ICCV 2021]

- a

# SCOUTER 🧐 [Li et al. ICCV 2021]



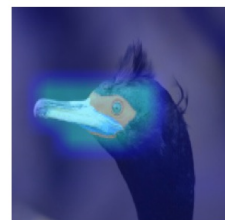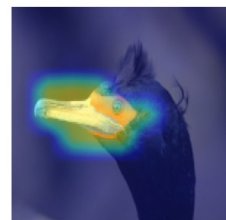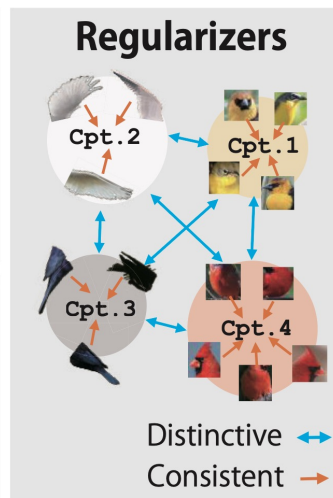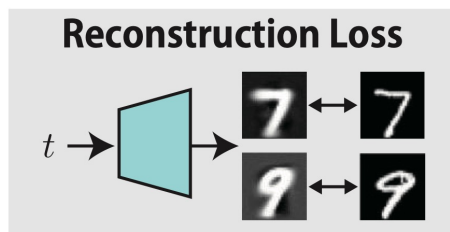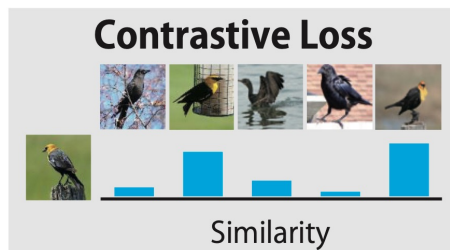| Input | SCOUTER (+) | | | SCOUTER (-) | | |
|-------|-------------|---|---|-------------|---|---|
| "7" | why "7" | why "1" | why "2" | why not "7" | why not "1" | why not "2" |
| loan | why loan | why tobacco | why cinema | why not loan | why not toba. | why not cine. |
| red-face cor. | why r-cor | why pel. cor. | why chat | why not r-cor | why not pel. | why not chat |

# BotCL [Wang et al., CVPR 2023]



**Concept Extractor**

Position Embedding

Concept Prototypes

$\{c_\kappa\}$

$P$

$F$

Slot Attention

Norm

$\{a_\kappa\}$

Sum

$t$

Classifier

$\hat{y}$

$\{v_\kappa\}$

Feature Aggregation

**Contrastive Loss**

Similarity

**Reconstruction Loss**

$t$

**Regularizers**

Cpt.2

Cpt.1

Cpt.3

Cpt.4

Distinctive

Consistent

# BotCL: Discovered concepts



Input    Cpt.1    Cpt.2    Cpt.3    Cpt.4    Cpt.5

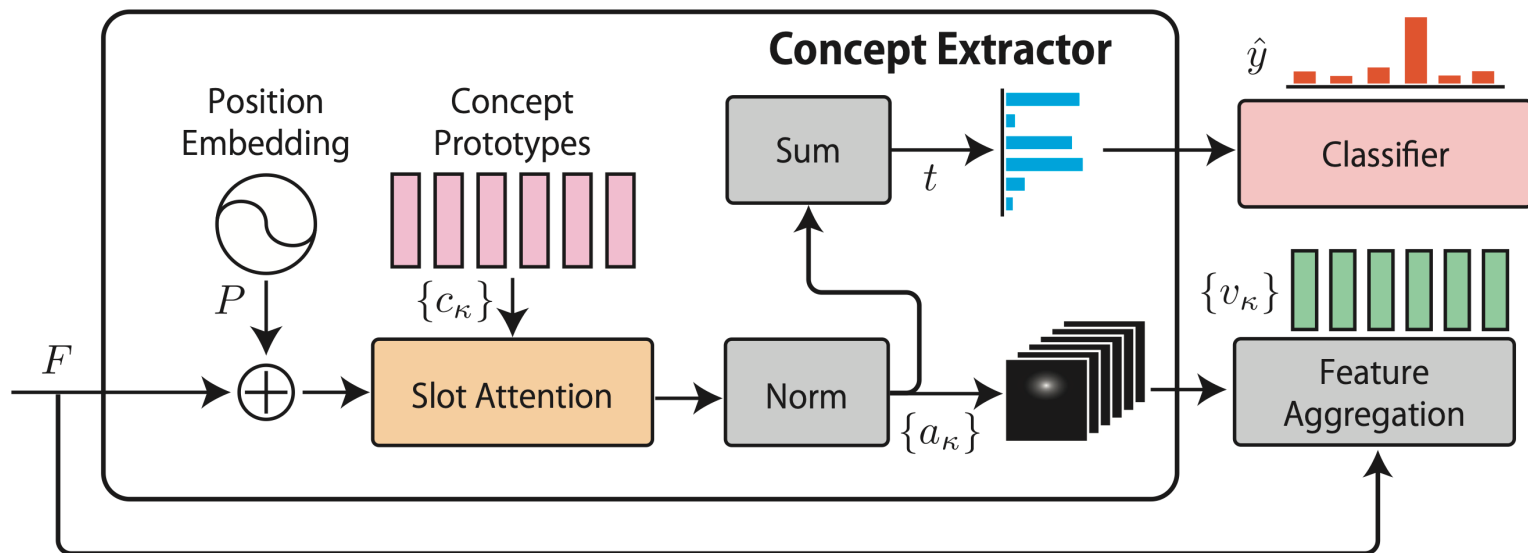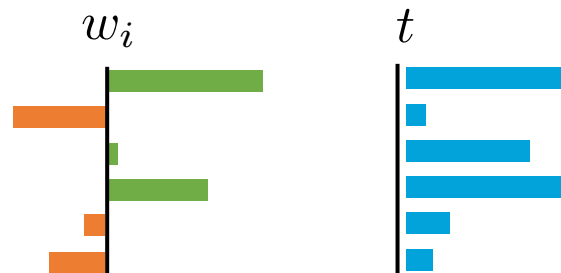Top-5 Activation

# BotCL: Evaluation of discovered concepts

# BotCL: Interpretation of classifier



- Our classifier is a single FC layer $y_i = w_i^\top t + b_i$
  - $w_i$ learns the correlations between each class and concept

# Next steps

- Efficient representation with a smaller set of concepts
    - More structured
    - Perhaps with grammar

- Without target task
    - Unsupervised (self-supervised) training for concept discovery

- Exploring some ways to identify spurious correlations
    - For vision and language tasks
    - For vision tasks